

## Graphical models, exponential families, and variational inference

Martin J. Wainwright<sup>1</sup> and Michael I. Jordan<sup>2</sup>

<sup>1</sup> *Department of Statistics, and Department of Electrical Engineering and Computer Science Berkeley 94720, USA, wainwrig@stat.berkeley.edu*

<sup>2</sup> *Department of Statistics, and Department of Electrical Engineering and Computer Science Berkeley 94720, USA, jordan@stat.berkeley.edu*

### Abstract

The formalism of probabilistic graphical models provides a unifying framework for capturing complex dependencies among random variables, and building large-scale multivariate statistical models. Graphical models have become a focus of research in many statistical, computational and mathematical fields, including bioinformatics, communication theory, statistical physics, combinatorial optimization, signal and image processing, information retrieval and statistical machine learning. Many problems that arise in specific instances—including the key problems of computing marginals and modes of probability distributions—are best studied in the general setting. Working with exponential family representations, and exploiting the conjugate duality between the cumulant function and the entropy for exponential families, we develop general variational representations of the problems of computing likelihoods, marginal probabilities and most probable configurations. We describe how a wide variety of algorithms—among them sum-product, cluster variational methods, expectation-

propagation, mean field methods, and max-product—can all be understood in terms of exact or approximate forms of these variational representations. The variational approach provides a complementary alternative to Markov chain Monte Carlo as a general source of approximation methods for inference in large-scale statistical models.

# Contents

---

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Background</b>	<b>5</b>
2.1	Probability distributions on graphs	5
2.1.1	Directed graphical models	6
2.1.2	Undirected graphical models	7
2.1.3	Factor graphs	8
2.2	Conditional independence	9
2.3	Statistical inference and exact algorithms	10
2.4	Applications	12
2.4.1	Hierarchical Bayesian models	12
2.4.2	Contingency table analysis	13
2.4.3	Constraint satisfaction and combinatorial optimization	14
2.4.4	Bioinformatics	15
2.4.5	Language and speech processing	18
2.4.6	Image processing and spatial statistics	20

2.4.7	Error-correcting coding	21
2.5	Exact inference algorithms	24
2.5.1	Message-passing on trees	25
2.5.2	Junction tree representation	28
2.6	Message-passing algorithms for approximate inference	33
<b>3</b>	<b>Graphical models as exponential families</b>	<b>36</b>
3.1	Exponential representations via maximum entropy	36
3.2	Basics of exponential families	38
3.3	Examples of graphical models in exponential form	41
3.4	Mean parameterization and inference problems	50
3.4.1	Mean parameter spaces and marginal polytopes	51
3.4.2	Role of mean parameters in inference problems	59
3.5	Properties of $A$	60
3.5.1	Derivatives and convexity	61
3.5.2	Forward mapping to mean parameters	62
3.6	Conjugate duality: Maximum likelihood and maximum entropy	65
3.6.1	General form of conjugate dual	65
3.6.2	Some simple examples	68
3.7	Challenges in high-dimensional models	71
<b>4</b>	<b>Sum-product, Bethe-Kikuchi, and expectation-propagation</b>	<b>74</b>
4.1	Sum-product and Bethe approximation	75
4.1.1	A tree-based outer bound to $\mathbb{M}(G)$	76
4.1.2	Bethe entropy approximation	80
4.1.3	Bethe variational problem and sum-product	82
4.1.4	Inexactness of Bethe and sum-product	87
4.1.5	Bethe optima and reparameterization	90
4.1.6	Bethe and loop series expansions	93
4.2	Kikuchi and hypertree-based methods	97
4.2.1	Hypergraphs and hypertrees	97
4.2.2	Kikuchi and related approximations	102

4.2.3	Generalized belief propagation	104
4.3	Expectation-propagation algorithms	106
4.3.1	Entropy approximations based on term decoupling	108
4.3.2	Optimality in terms of moment-matching	115
<b>5</b>	<b>Mean field methods</b>	<b>124</b>
5.1	Tractable families	125
5.2	Optimization and lower bounds	126
5.2.1	Generic mean field procedure	127
5.2.2	Mean field and Kullback-Leibler divergence	128
5.3	Examples of mean field algorithms	131
5.3.1	Naive mean field updates	131
5.3.2	Structured mean field	135
5.4	Non-convexity of mean field	141
<b>6</b>	<b>Variational methods in parameter estimation</b>	<b>146</b>
6.1	Estimation in fully observed models	146
6.1.1	Maximum likelihood for triangulated graphs	147
6.1.2	Iterative methods for computing MLEs	149
6.2	Partially observed models and expectation-maximization	151
6.2.1	Exact EM algorithm in exponential families	152
6.2.2	Variational EM	156
6.3	Variational Bayes	157
<b>7</b>	<b>Convex relaxations and upper bounds</b>	<b>162</b>
7.1	Generic convex combinations and convex surrogates	164
7.2	Variational methods from convex relaxations	167
7.2.1	Tree-reweighted sum-product and Bethe	167
7.2.2	Reweighted Kikuchi approximations	173
7.2.3	Convex forms of expectation-propagation	177
7.2.4	Planar graph decomposition	178
7.3	Other convex variational methods	179
7.3.1	Semidefinite constraints and log-determinant bound	179

7.3.2	Other entropy approximations and polytope bounds	182
7.4	Algorithmic stability	182
7.5	Convex surrogates in parameter estimation	185
7.5.1	Surrogate likelihoods	186
7.5.2	Optimizing surrogate likelihoods	186
7.5.3	Penalized surrogate likelihoods	188
<b>8</b>	<b>Max-product and LP relaxations</b>	<b>190</b>
8.1	Variational formulation of computing modes	190
8.2	Max-product and linear programming for trees	193
8.3	Max-product for Gaussians and other convex problems	197
8.4	First-order LP relaxation and reweighted max-product	199
8.4.1	Basic properties of first-order LP relaxation	200
8.4.2	Connection to max-product message-passing	205
8.4.3	Modified forms of message-passing	208
8.4.4	Examples of the first-order LP relaxation	211
8.5	Higher-order LP relaxations	220
<b>9</b>	<b>Moment matrices and conic relaxations</b>	<b>228</b>
9.1	Moment matrices and their properties	229
9.1.1	Multinomial and indicator bases	230
9.1.2	Marginal polytopes for hypergraphs	233
9.2	Semidefinite bounds on marginal polytopes	233
9.2.1	Lasserre sequence	234
9.2.2	Tightness of semidefinite outer bounds	236
9.3	Link to LP relaxations and graphical structure	240
9.4	Second-order cone relaxations	244
<b>10</b>	<b>Discussion</b>	<b>246</b>
<b>A</b>	<b>Background material</b>	<b>249</b>
A.1	Background on graphs and hypergraphs	249

A.2	Basics of convex sets and functions	251
A.2.1	Convex sets and cones	251
A.2.2	Convex and affine hulls	252
A.2.3	Affine hulls and relative interior	252
A.2.4	Polyhedra and their representations	253
A.2.5	Convex functions	253
A.2.6	Conjugate duality	255
<b>B</b>	<b>Proofs for exponential families and duality</b>	<b>257</b>
B.1	Proof of Theorem 1	257
B.2	Proof of Theorem 2	259
B.3	General properties of $\mathcal{M}$ and $A^*$	261
B.3.1	Properties of $\mathcal{M}$	261
B.3.2	Properties of $A^*$	262
B.4	Proof of Theorem 3(b)	263
B.5	Proof of Theorem 5	264
<b>C</b>	<b>Variational principles for multivariate Gaussians</b>	<b>265</b>
C.1	Gaussian with known covariance	265
C.2	Gaussian with unknown covariance	267
C.3	Gaussian mode computation	268
<b>D</b>	<b>Clustering and augmented hypergraphs</b>	<b>270</b>
D.1	Covering augmented hypergraphs	270
D.2	Specification of compatibility functions	274
<b>E</b>	<b>Miscellaneous results</b>	<b>276</b>
E.1	Möbius inversion	276
E.2	Pairwise Markov random fields	277
	<b>References</b>	<b>279</b>

# 1

---

## Introduction

---

Graphical models bring together graph theory and probability theory in a powerful formalism for multivariate statistical modeling. In various applied fields including bioinformatics, speech processing, image processing and control theory, statistical models have long been formulated in terms of graphs, and algorithms for computing basic statistical quantities such as likelihoods and score functions have often been expressed in terms of recursions operating on these graphs; examples include phylogenies, pedigrees, hidden Markov models, Markov random fields, and Kalman filters. These ideas can be understood, unified and generalized within the formalism of graphical models. Indeed, graphical models provide a natural tool for formulating variations on these classical architectures, as well as for exploring entirely new families of statistical models. Accordingly, in fields that involve the study of large numbers of interacting variables, graphical models are increasingly in evidence.

Graph theory plays an important role in many computationally-oriented fields, including combinatorial optimization, statistical physics and economics. Beyond its use as a language for formulating models, graph theory also plays a fundamental role in assessing computational



## 2 Introduction

complexity and feasibility. In particular, the running time of an algorithm or the magnitude of an error bound can often be characterized in terms of structural properties of a graph. This statement is also true in the context of graphical models. Indeed, as we discuss, the computational complexity of a fundamental method known as the *junction tree algorithm*—which generalizes many of the recursive algorithms on graphs cited above—can be characterized in terms of a natural graph-theoretic measure of interaction among variables. For suitably sparse graphs, the junction tree algorithm provides a systematic solution to the problem of computing likelihoods and other statistical quantities associated with a graphical model.

Unfortunately, many graphical models of practical interest are not “suitably sparse,” so that the junction tree algorithm no longer provides a viable computational framework. One popular source of methods for attempting to cope with such cases is the *Markov chain Monte Carlo* (MCMC) framework, and indeed there is a significant literature on the application of MCMC methods to graphical models [e.g., 26, 90]. Our focus in this paper is rather different: we present an alternative computational methodology for statistical inference that is based on *variational methods*. These techniques provide a general class of alternatives to MCMC, and have applications outside of the graphical model framework. As we will see, however, they are particularly natural in their application to graphical models, due to their relationships with the structural properties of graphs.

The phrase “variational” itself is an umbrella term that refers to various mathematical tools for optimization-based formulations of problems, as well as associated techniques for their solution. The general idea is to express a quantity of interest as the solution of an optimization problem. The optimization problem can then be “relaxed” in various ways, either by approximating the function to be optimized or by approximating the set over which the optimization takes place. Such relaxations, in turn, provide a means of approximating the original quantity of interest.

The roots of both MCMC methods and variational methods lie in statistical physics. Indeed, the successful deployment of MCMC methods in statistical physics motivated and predated their entry into statis-

tics. However, the development of MCMC methodology specifically designed for statistical problems has played an important role in sparking widespread application of such methods in statistics [85]. A similar development in the case of variational methodology would be of significant interest. In our view, the most promising avenue towards a variational methodology tuned to statistics is to build on existing links between variational analysis and the exponential family of distributions [4, 10, 40, 71]. Indeed, the notions of convexity that lie at the heart of the statistical theory of the exponential family have immediate implications for the design of variational relaxations. Moreover, these variational relaxations have particularly interesting algorithmic consequences in the setting of graphical models, where they again lead to recursions on graphs.

Thus, we present a story with three interrelated themes. We begin in Section 2 with a discussion of graphical models, providing both an overview of the general mathematical framework, and also presenting several specific examples. All of these examples, as well as the majority of current applications of graphical models, involve distributions in the exponential family. Accordingly, Section 3 is devoted to a discussion of exponential families, focusing on the mathematical links to convex analysis, and thus anticipating our development of variational methods. In particular, the principal object of interest in our exposition is a certain *conjugate dual* relation associated with exponential families. From this foundation of conjugate duality, we develop a general variational representation for computing likelihoods and marginal probabilities in exponential families. Subsequent sections are devoted to the exploration of various instantiations of this variational principle, both in exact and approximate forms, which in turn yield various algorithms for computing exact and approximate marginal probabilities, respectively. In Section 4, we discuss the connection between the Bethe approximation and the sum-product algorithm, including both its exact form for trees and approximate form for graphs with cycles, as well as the connection between Bethe-like approximations and other algorithms, including generalized sum-product and expectation-propagation. In Section 5, we discuss the class of mean field methods, which arise from a qualitatively different approximation to the exact variational principle, one

#### 4 Introduction

which generates lower bounds on the likelihood. In Section 6, we discuss the role of variational methods in parameter estimation, including both the fully observed and partially observed cases. Both Bethe-type and mean field methods are based on non-convex optimization problems, which typically have multiple solutions. In contrast, Section 7 discusses variational methods based on convex relaxations of the exact variational principle, many of which are also guaranteed to yield upper bounds on the log likelihood. Finally, in Section 8, we discuss the problem of computing modes, which corresponds to an integer program for the case of discrete random variables, and various relaxations of the exact variational principle, including those based on linear programming as well as other types of conic relaxations.

The scope of this paper is limited in the following sense: given a distribution represented as a graphical model, we are concerned with the problem of computing marginal probabilities (including likelihoods), as well as the problem of computing modes. We refer to such computational tasks as problems of “probabilistic inference,” or “inference” for short. As with presentations of MCMC methods, such a limited focus may appear to aim most directly at applications in Bayesian statistics. While Bayesian statistics is indeed a natural terrain for deploying many of the methods that we present here, we see these methods as having applications throughout statistics, within both the frequentist and Bayesian paradigms, and we indicate some of these applications at various junctures in the paper.

# 2

---

## Background

---

We begin with background on graphical models. The key idea is that of *factorization*: a graphical model consists of a collection of probability distributions that factorize according to the structure of an underlying graph. Here we are using the terminology “distribution” loosely; our notation  $p$  should be understood as a mass function (density with respect to counting measure) in the discrete case, and a density with respect to Lebesgue measure in the continuous case. Our focus in this section is the interplay between probabilistic notions such as conditional independence on one hand, and on the other hand, graph-theoretic notions such as cliques and separation.

### 2.1 Probability distributions on graphs

We begin by describing the various types of graphical formalisms that are useful. A graph  $G = (V, E)$  is formed by a collection of vertices  $V$ , and a collection of edges  $E \subset V \times V$ . Each edge consists of a pair vertices  $s, t \in E$ , and may either be *undirected*, in which case there is no distinction between edge  $(s, t)$  and edge  $(t, s)$ , or *directed*, in which case we write  $(s \rightarrow t)$  to indicate the direction. See Appendix A.1 for

## 6 Background

more background on graphs and their properties.

In order to define a graphical model, we associate with each vertex  $s \in V$  a random variable  $X_s$  taking values in some space  $\mathcal{X}_s$ . Depending on the application, this state space  $\mathcal{X}_s$  may either be continuous, (e.g.,  $\mathcal{X}_s = \mathbb{R}$ ) or discrete (e.g.,  $\mathcal{X}_s = \{0, 1, \dots, r - 1\}$ ). We use lower-case letters (e.g.,  $x_s \in \mathcal{X}_s$ ) to denote particular elements of  $\mathcal{X}_s$ , so that the notation  $\{X_s = x_s\}$  corresponds to the event that the random variable  $X_s$  takes the value  $x_s \in \mathcal{X}_s$ . For any subset  $A$  of the vertex set  $V$ , we define  $X_A := \{X_s \mid s \in A\}$ , with the notation  $x_A := \{x_s \mid s \in A\}$  corresponding to the analogous quantity for values of the random vector  $X_A$ . Similarly, we define  $\otimes_{s \in A} \mathcal{X}_s$  to be the Cartesian product of the state spaces for each of the elements of  $X_A$ .

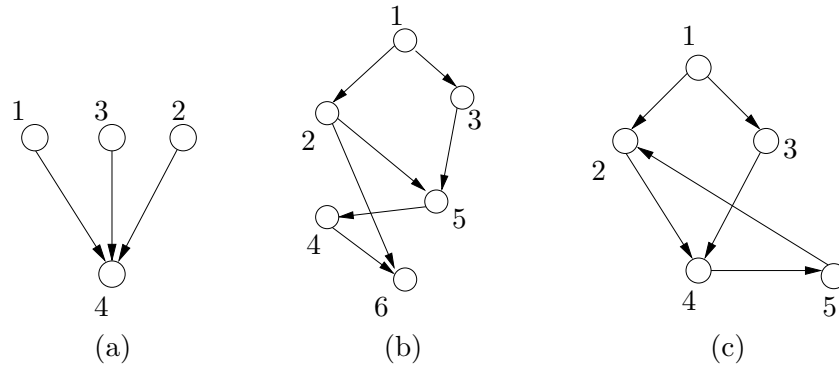
### 2.1.1 Directed graphical models

Given a directed graph with edges  $(s \rightarrow t)$ , we say that  $t$  is a child of  $s$ , or conversely, that  $s$  is a parent of  $t$ . For any vertex  $s \in V$ , let  $\pi(s)$  denote the set of all parents of given node  $s \in V$ . (If a vertex  $s$  has no parents, then the set  $\pi(s)$  should be understood to be empty.) A directed cycle is a sequence  $(s_1, s_2, \dots, s_k)$  such that  $(s_i \rightarrow s_{i+1}) \in E$  for all  $i = 1, \dots, k - 1$ , and  $(s_k \rightarrow s_1) \in E$ . See Figure 2.1 for an illustration of these concepts.

Now suppose that  $G$  is a directed acyclic graph (DAG), meaning that every edge is directed, and that the graph contains no directed cycles. For any such DAG, we can define a partial order on the vertex set  $V$  by the notion of ancestry: we say that node  $s$  is an ancestor of  $u$  if there is a directed path  $(s, t_1, t_2, \dots, t_k, u)$  (see Figure 2.1(b)). Given a DAG, for each vertex  $s$  and its parent set  $\pi(s)$ , let  $p_s(x_s \mid x_{\pi(s)})$  denote a non-negative function over the variables  $(x_s, x_{\pi(s)})$ , normalized such that  $\int p_s(x_s \mid x_{\pi(s)}) dx_s = 1$ . In terms of these local functions, a *directed graphical model* consists of a collection of probability distributions (densities or mass functions) that factorize in the following way:

$$p(x) = \prod_{s \in V} p_s(x_s \mid x_{\pi(s)}). \quad (2.1)$$

It can be verified that our use of notation is consistent, in that



**Fig. 2.1.** (a) A simple directed graphical model with four variables ( $X_1, X_2, X_3, X_4$ ). Vertices  $\{1, 2, 3\}$  are all parents of vertex 4, written  $\pi(4) = \{1, 2, 3\}$ . (b) A more complicated directed acyclic graph (DAG) that defines a partial order on its vertices. Note that vertex 6 is a child of vertex 2, and vertex 1 is an ancestor of 6. (c) A forbidden directed graph (non-acyclic) that includes the directed cycle ( $2 \rightarrow 4 \rightarrow 5 \rightarrow 2$ ).

$p_s(x_s \mid x_{\pi(s)})$  is, in fact, the conditional distribution of  $X_s = x_s$  given  $\{X_{\pi(s)} = x_{\pi(s)}\}$  for the global distribution  $p(x)$  defined by the factorization (2.1). This follows by an inductive argument that makes use of the normalization condition imposed on  $p_s(\cdot)$ , and the partial ordering induced by the ancestor relationship of the DAG.

### 2.1.2 Undirected graphical models

In the undirected case, the probability distribution factorizes according to functions defined on the *cliques* of the graph. A clique  $C$  is a fully-connected subset of the vertex set  $V$ , meaning that  $(s, t) \in E$  for all  $s, t \in C$ . Let us associate with each clique  $C$  a *compatibility function*  $\psi_C : (\otimes_{s \in C} \mathcal{X}_s) \rightarrow \mathbb{R}_+$ . Recall that  $\otimes_{s \in C} \mathcal{X}_s$  denotes the Cartesian product of the state spaces of the random vector  $X_C$ , such that the compatibility function  $\psi_C$  is a local quantity, defined only for elements  $x_C$  within the clique.

With this notation, an *undirected graphical model*—also known as a *Markov random field* (MRF), or a *Gibbs distribution*—is a collection

of distributions that factorize as follows:

$$p(x_1, x_2, \dots, x_m) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(x_C), \quad (2.2)$$

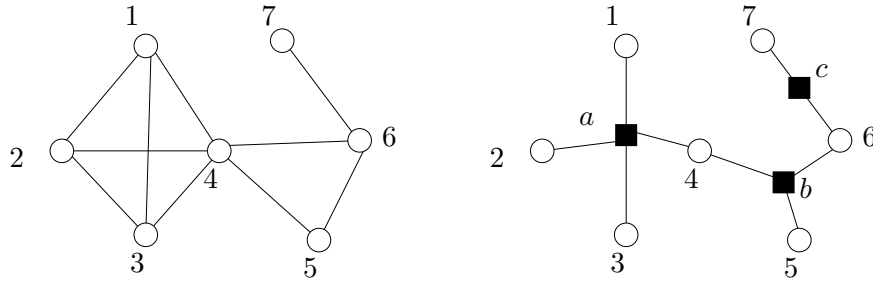
where  $Z$  is a constant chosen to ensure that the distribution is normalized. The set  $\mathcal{C}$  is often taken to be the set of all *maximal cliques* of the graph; i.e., the set of cliques that are not properly contained within any other clique. This condition can be imposed without loss of generality, because any representation based on non-maximal cliques can always be converted to one based on maximal cliques by redefining the compatibility function on a maximal clique to be the product over the compatibility functions on the subsets of that clique. However, there may be computational benefits to using a non-maximal representation—in particular, algorithms may be able to exploit specific features of the factorization special to the non-maximal representation. For this reason, we do not necessarily restrict compatibility functions to maximal cliques only, but instead define the set  $\mathcal{C}$  to contain all cliques. (Factor graphs, discussed in the following section, allow for a finer-grained specification of factorization properties.)

It is important to understand that for a general undirected graph the compatibility functions  $\psi_C$  need not have any obvious or direct relation to marginal distributions defined over the graph cliques. This is in contrast to the directed factorization (2.1) where the factors are marginal conditional probabilities.

### 2.1.3 Factor graphs

For large graphs, the factorization properties of a graphical model, whether undirected or directed, may be difficult to visualize from the usual depictions of graphs. The formalism of *factor graphs* provides an alternative graphical representation, one which emphasizes the factorization of the distribution [140, 154].

Let  $F$  represent an index set for the set of factors defining a graphical model distribution. In the undirected case, this set indexes the collection  $\mathcal{C}$ , while in the directed case  $F$  indexes the set of parent-child neighborhoods. We then consider a bipartite graph  $G' = (V, F, E')$ , where  $V$  is the original set of vertices, and  $E'$  is a new edge set, joining



**Fig. 2.2.** Illustration of undirected graphical models and factor graphs. (a) An undirected graph on  $m = 7$  vertices, with maximal cliques  $\{1, 2, 3, 4\}$ ,  $\{4, 5, 6\}$  and  $\{6, 7\}$ . (b) Equivalent representation of the undirected graph in (a) as a factor graph, assuming that we define compatibility functions only on the maximal cliques in (a). The factor graph is a bipartite graph with vertex set  $V = \{1, \dots, 7\}$  and factor set  $F = \{a, b, c\}$ , one for each of the compatibility functions of the original undirected graph.

only vertices  $s \in V$  to factors  $a \in F$ . In particular, edge  $(s, a) \in E'$  if and only if  $x_s$  participates in the factor indexed by  $a \in F$ . See Figure 2.2(b) for an illustration.

For undirected models, the factor graph representation is of particular value when  $\mathcal{C}$  consists of more than the maximal cliques. Indeed, the compatibility functions for the non-maximal cliques do not have an explicit representation in the usual representation of an undirected graph—however, the factor graph makes them explicit.

## 2.2 Conditional independence

Families of probability distributions as defined as in equations (2.1) or (2.2) also have a characterization in terms of conditional independencies among subsets of random variables—the *Markov properties* of the graphical model. We only touch upon this characterization here, as it is not needed in the remainder of the paper; for a full treatment, we refer the interested reader to Lauritzen [150].

For undirected graphical models, conditional independence is identified with the graph-theoretic notion of *reachability*. In particular, let  $A$ ,  $B$  and  $C$  be an arbitrary triple of mutually disjoint subsets of vertices. Let us stipulate that  $x_A$  be independent of  $x_B$  given  $x_C$  if there



is no path from a vertex in  $A$  to a vertex in  $B$  when we remove the vertices  $C$  from the graph. Ranging over all possible choices of subsets  $A$ ,  $B$  and  $C$  gives rise to a list of conditional independence assertions. It can be shown that this list is always consistent (i.e., there exist probability distributions that satisfy all of these assertions); moreover, the set of probability distributions that satisfy these assertions is exactly the set of distributions defined by (2.2) ranging over all possible choices of compatibility functions.

Thus, there are two equivalent characterizations of the family of probability distributions associated with an undirected graph. This equivalence is a fundamental mathematical result, linking an algebraic concept (factorization) and a graph-theoretic concept (reachability). This result also has algorithmic consequences, in that it reduces the problem of assessing conditional independence to the problem of assessing reachability on a graph, which is readily solved using simple breadth-first search algorithms [53].

An analogous result holds in the case of directed graphical models, with the only alteration being a different notion of reachability [150]. Once again, it is possible to establish an equivalence between the set of probability distributions specified by the directed factorization (2.1), and that defined in terms of a set of conditional independence assertions.

### 2.3 Statistical inference and exact algorithms

Given a probability distribution  $p$  defined by a graphical model, our focus will be solving one or more of the following *computational inference problems*:

- (a) computing the likelihood of observed data.
- (b) computing the marginal distribution  $p(x_A)$  over a particular subset  $A \subset V$  of nodes.
- (c) computing the conditional distribution  $p(x_A | x_B)$ , for disjoint subsets  $A$  and  $B$ , where  $A \cup B$  is in general a proper subset of  $V$ .
- (d) computing a mode of the density (i.e., an element  $\hat{x}$  in the set  $\arg \max_{x \in \mathcal{X}^m} p(x)$ ).

Clearly problem (a) is a special case of problem (b). The computation of a conditional probability in (c) is similar in that it also requires marginalization steps, an initial one to obtain the numerator  $p(x_A, x_B)$ , and a further step to obtain the denominator  $p(x_B)$ . In contrast, the problem of computing modes stated in (d) is fundamentally different, since it entails maximization rather than integration. Nonetheless, our variational development in the sequel will highlight some important connections between the problem of computing marginals and that of computing modes.

To understand the challenges inherent in these inference problems, consider the case of a discrete random vector  $x \in \mathcal{X}^m$ , where  $\mathcal{X}_s = \{0, 1, \dots, r-1\}$  for each vertex  $s \in V$ . A naive approach to computing a marginal at a single node—say  $p(x_s)$ —entails summing over all configurations of the form  $\{x' \in \mathcal{X}^m \mid x'_s = x_s\}$ . Since this set has  $r^{m-1}$  elements, it is clear that a brute force approach will rapidly become intractable. Even with binary variables ( $r = 2$ ) and a graph with  $m \approx 100$  nodes (a small size for many applications), this summation is beyond the reach of brute-force computation. Similarly, in this discrete case, computing a mode entails solving an integer programming problem over an exponential number of configurations. For continuous random vectors, the problems are no easier<sup>1</sup> and typically harder, since they require computing a large number of integrals.

For undirected graphs without cycles or for directed graphs in which each node has a single parent—known generically as *trees* in either case—these inference problems can be solved exactly by recursive “message-passing” algorithms of a dynamic programming nature, with a computational complexity that scales only linearly in the number of nodes. In particular, for the case of computing marginals, the dynamic programming solution takes the form of a general algorithm known as the *sum-product algorithm*, whereas for the problem of computing modes it takes the form of an analogous algorithm known as the *max-product algorithm*. We describe these algorithms in Section 2.5.1. More generally, as we discuss in Section 2.5.2, the *junction tree algorithm* provides a solution to inference problems for arbitrary graphs.

---

<sup>1</sup>The Gaussian case is an important exception to this statement.

The junction tree algorithm has a computational complexity that is exponential in a quantity known as the *treewidth* of the graph.

## 2.4 Applications

Before turning to these algorithmic issues, however, it is helpful to ground the discussion by considering some examples of applications of graphical models. We present examples of the use of graphical models in the general areas of Bayesian hierarchical modeling, contingency table analysis, and combinatorial optimization and satisfiability, as well as specific examples in bioinformatics, speech and language processing, image processing, spatial statistics, communication and coding theory.

### 2.4.1 Hierarchical Bayesian models

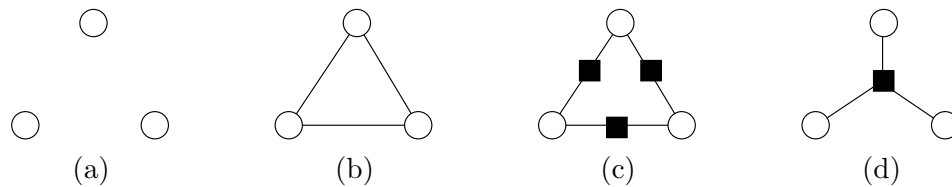
The Bayesian framework treats all model quantities —observed data, latent variables, parameters, nuisance variables—as random variables. Thus, in a graphical model representation of a Bayesian model, all such variables appear explicitly as vertices in the graph. The general computational machinery associated with graphical models applies directly to Bayesian computations of quantities such as marginal likelihoods and posterior probabilities of parameters. Although Bayesian models can be represented using either directed or undirected graphs, it is the directed formalism that is most commonly encountered in practice. In particular, in hierarchical Bayesian models, the specification of prior distributions generally involves additional parameters (i.e., “hyperparameters”), and the overall model is specified as a set of conditional probabilities linking hyperparameters, parameters and data. Taking the product of such conditional probability distributions defines the joint probability; this factorization is simply a particular instance of equation (2.1).

There are several advantages to treating a hierarchical Bayesian model as a directed graphical model. First, hierarchical models are often specified by making various assertions of conditional independence. These assertions imply other conditional independence relationships, and the reachability algorithms (mentioned in Section 2.1) provide a systematic method for investigating all such relationships. Second, the

visualization provided by the graph can be useful both for understanding the model (including the basic step of verifying that the graph is acyclic), as well for exploring extensions. Finally, general computational methods such as MCMC and variational inference algorithms can be implemented for general graphical models, and hence apply to hierarchical models in graphical form. These advantages and others have led to the development of general software programs for specifying and manipulating hierarchical Bayesian models via the directed graphical model formalism [90].

### 2.4.2 Contingency table analysis

Contingency tables are a central tool in categorical data analysis [148, 77, 1], dating back to the seminal work of Pearson, Yule and Fisher. A  $m$ -dimensional contingency table with  $r$  levels describes a probability distribution over  $m$  random variables  $(X_1, \dots, X_m)$ , each of which takes one of  $r$  possible values. For the special case  $m = 2$ , a contingency table with  $r$  levels is simply a  $r \times r$  matrix of non-negative elements summing to one, whereas for  $m > 2$ , it is a multi-dimensional array with  $r^m$  non-negative entries summing to one. Thus, the table fully specifies a probability distribution over a random vector  $(X_1, \dots, X_m)$ , where each  $X_s$  takes one of  $r$  possible values.



**Fig. 2.3.** Some simple graphical interactions in contingency table analysis. (a) Independence model. (b) General dependent model. (c) Pairwise interactions only. (d) Triplet interactions.

Contingency tables can be modeled within the framework of graphical models, and graphs provide a useful visualization of many natural questions. For instance, one central question in contingency table analysis is distinguishing different orders of interaction in the data. As a concrete example, given  $m = 3$  variables, a simple question is testing

whether or not the three variables are independent. From a graphical model perspective, this test amounts to distinguishing whether the associated interaction graph is completely disconnected or not (see panels (a) and (b) in Figure 2.3). A more subtle test is to distinguish whether the random variables interact in only a pairwise manner (i.e., with factors  $\psi_{12}$ ,  $\psi_{13}$ , and  $\psi_{23}$  in equation (2.2)), or if there is actually a three-way interaction (i.e., a term  $\psi_{123}$  in the factorization (2.2)). Interestingly, these two factorization choices cannot be distinguished using a standard undirected graph; as illustrated in Figure 2.3(b), the fully connected graph on 3 vertices does not distinguish between pairwise interactions without triplet interactions, versus triplet interaction (possibly including pairwise interaction as well). In this sense, the factor graph formalism is more expressive, since it can distinguish between pairwise and triplet interactions, as shown in panels (c) and (d) of Figure 2.3.

### 2.4.3 Constraint satisfaction and combinatorial optimization

Problems of constraint satisfaction and combinatorial optimization arise in a wide variety of areas, among them artificial intelligence [63, 188], communication theory [84], computational complexity theory [52], statistical image processing [86], and bioinformatics [190]. Many problems in satisfiability and combinatorial optimization [177] are defined in graph-theoretic terms, and thus are naturally recast in the graphical model formalism.

Let us illustrate by considering what is perhaps the best-known example of satisfiability, namely the 3-SAT problem [52]. It is naturally described in terms of a factor graph, and a collection of binary random variables  $(X_1, X_2, \dots, X_m) \in \{0, 1\}^m$ . For a given triplet  $\{s, t, u\}$  of vertices, let us specify some “forbidden” configuration  $(z_s, z_t, z_u) \in \{0, 1\}^3$ , and then define a triplet compatibility function

$$\psi_{stu}(x_s, x_t, x_u) = \begin{cases} 0 & \text{if } (x_s, x_t, x_u) = (z_s, z_t, z_u), \\ 1 & \text{otherwise.} \end{cases} \quad (2.3)$$

Each such compatibility function is referred to as a *clause* in the satisfiability literature, where such compatibility functions are

typically encoded in terms of logical operations. For instance, if  $(z_s, z_t, z_u) = (0, 0, 0)$ , then the function (2.3) can be written compactly as  $\psi_{stu}(x_s, x_t, x_u) = x_s \vee x_t \vee x_u$ , where  $\vee$  denotes the logical OR operation between two Boolean symbols.

As a graphical model, a single clause corresponds to the model in Figure 2.3(d); of greater interest are models over larger collections of binary variables, involving many clauses. One basic problem in satisfiability is to determine whether a given set of clauses  $F$  are satisfiable, meaning that there exists some configuration  $x \in \{0, 1\}^m$  such that  $\prod_{(s,t,u) \in F} \psi_{stu}(x_s, x_t, x_u) = 1$ . In this case, the factorization

$$p(x) = \frac{1}{Z} \prod_{(s,t,u) \in F} \psi_{stu}(x_s, x_t, x_u)$$

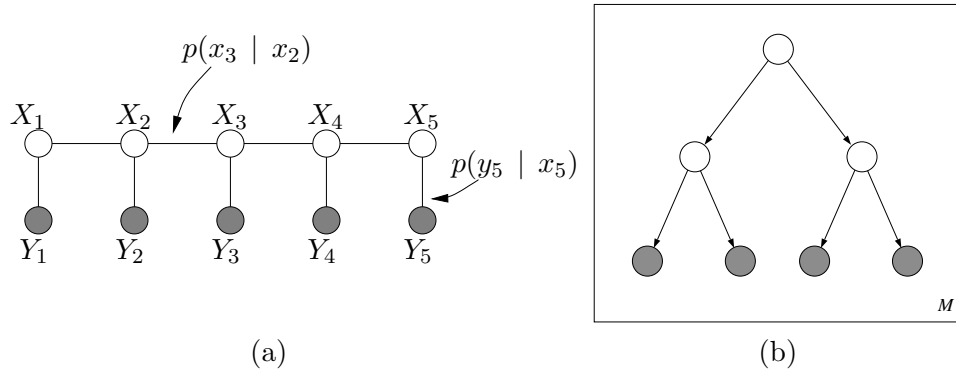
defines the uniform distribution over the set of satisfying assignments.

When instances of 3-SAT are randomly constructed—for instance, by fixing a clause density  $\alpha$ , and drawing  $q = \lceil \alpha m \rceil$  clauses over  $m$  variables uniformly from all triplets—the graphs tend to have a locally “tree-like” structure; see the factor graph in Figure 2.9(b) for an illustration. One major problem is determining the value  $\alpha^*$  at which this ensemble is conjectured to undergo a phase transition from being satisfiable (for small  $\alpha$ , hence with few constraints) to being unsatisfiable (for sufficiently large  $\alpha$ ). The survey propagation algorithm, developed in the statistical physics community [168, 167], turns out to be very successful in solving random instances of 3-SAT. Survey propagation turns out to be an instance of the sum-product or belief propagation algorithm, but as applied to an alternative graphical model for satisfiability problems [38, 159].

#### 2.4.4 Bioinformatics

Many classical models in bioinformatics are instances of graphical models, and the associated framework is often exploited in designing new models. In this section, we briefly review some instances of graphical models in both bioinformatics, both classical and recent. Sequential data play a central role in bioinformatics, and the workhorse underlying the modeling of sequential data is the *hidden Markov model* (HMM)

shown in Figure 2.4(a). The HMM is in essence a dynamical version of a



**Fig. 2.4.** (a) The graphical model representation of a generic hidden Markov model. The shaded nodes  $\{Y_1, \dots, Y_5\}$  are the observations and the unshaded nodes  $\{X_1, \dots, X_5\}$  are the hidden state variables. The latter form a Markov chain, in which  $X_s$  is independent of  $X_u$  conditional on  $X_t$ , where  $s < t < u$ . (b) The graphical model representation of a phylogeny on four extant organisms and  $M$  sites. The tree encodes the assumption that there is a first speciation event and then two further speciation events that lead to the four extant organisms. The box around the tree (a “plate”) is a graphical model representation of replication, here representing the assumption that the  $M$  sites evolve independently.

finite mixture model, in which observations are generated conditionally on a underlying latent (“hidden”) state variable. The state variables, which are generally taken to be discrete random variables following a multinomial distribution, form a Markov chain. The graphical model in Figure 2.4(a) is also a representation of the state-space model underlying Kalman filtering and smoothing, where the state variable is a Gaussian vector. These models thus also have a right to be referred to as “hidden Markov models,” but the terminology is most commonly used to refer to models in which the state variables are discrete.

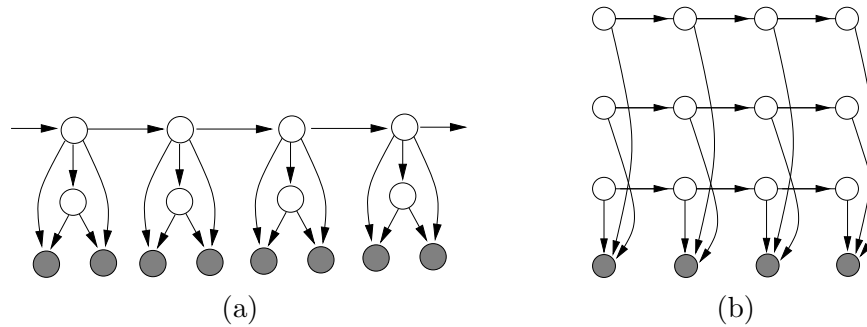
Applying the junction tree formalism to the HMM yields an algorithm that passes messages in both directions along the chain of state variables, and computes the marginal probabilities  $p(x_t, x_{t+1} | y)$  and  $p(x_t | y)$ . In the HMM context, this algorithm is often referred to as the *forward-backward algorithm* [192]. These marginal probabilities are often of interest in and of themselves, but are also important in their

role as expected sufficient statistics in an expectation-maximization (EM) algorithm for estimating the parameters of the HMM. Similarly, the maximum a posteriori state sequence can also be computed by the junction tree algorithm (with summation replaced by maximization)—in the HMM context the resulting algorithm is generally referred to as the *Viterbi algorithm* [79].

Gene-finding provides an important example of the application of the HMM [69]. To a first order of approximation, the genomic sequence of an organism can be segmented into regions containing genes and intergenic regions (that separate the genes), where a gene is defined as a sequence of nucleotides that can be further segmented into meaningful intragenic structures (exons and introns). The boundaries between these segments are highly stochastic and hence difficult to find reliably. HMMs have been the methodology of choice for attacking this problem, with designers choosing states and state transitions to reflect biological knowledge concerning gene structure [41]. Hidden Markov models are also used to model certain aspects of protein structure. For example, membrane proteins are specific kinds of proteins that embed themselves in the membranes of cells, and play important roles in the transmission of signals in and out of the cell. These proteins loop in and out of the membrane many times, alternating between hydrophilic amino acids (which prefer the environment of the membrane) and hydrophobic amino acids (which prefer the environment inside or outside the membrane). These and other biological facts are used to design the states and state transition matrix of the *transmembrane hidden Markov model*, an HMM for modeling membrane proteins [138].

Tree-structured models also play an important role in bioinformatics and language processing. For example, phylogenetic trees can be treated as graphical models. As shown in Figure 2.4(b), a phylogenetic tree is a tree-structured graphical model in which a set of observed nucleotides (or other biological characters) are assumed to have evolved from an underlying set of ancestral species. The conditional probabilities in the tree are obtained from evolutionary substitution models, and the computation of likelihoods are achieved by a recursion on the tree known as “pruning” [76]. This recursion is a special case of the junction tree algorithm.





**Fig. 2.5.** Variations on the HMM used in bioinformatics. (a) A phylogenetic HMM. (b) The factorial HMM.

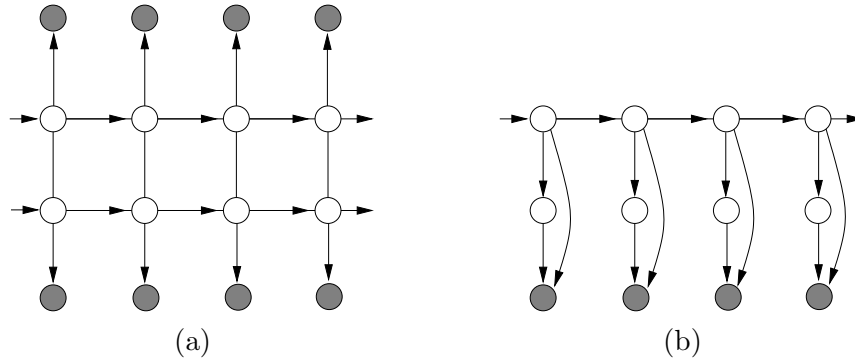
Figure 2.5 provides examples of more complex graphical models that have been explored in bioinformatics. Figure 2.5(a) shows a *hidden Markov phylogeny*, an HMM in which the observations are sets of nucleotides related by a phylogenetic tree [162, 189, 213]. This model has proven useful for gene-finding in the human genome based on data for multiple primate species [162]. Figure 2.5(b) shows a *factorial HMM*, in which multiple chains are coupled by their links to a common set of observed variables [89]. This model captures the problem of *multi-locus linkage analysis* in genetics, where the state variables correspond to phase (maternal or paternal) along the chromosomes in meiosis [227].

#### 2.4.5 Language and speech processing

In language problems, HMMs also play a fundamental role. An example is the part-of-speech problem, in which words in sentences are to be labeled as to their part of speech (noun, verb, adjective, etc). Here the state variables are the parts of speech, and the transition matrix is estimated from a corpus via the EM algorithm [160]. The result of running the Viterbi algorithm on a new sentence is a tagging of the sentence according to the hypothesized parts of speech of the words in the sentence. Moreover, essentially all modern speech recognition systems are built on the foundation of HMMs [114]. In this case the observations are generally a sequence of short-range speech spectra, and the states correspond to longer-range units of speech such as phonemes or pairs

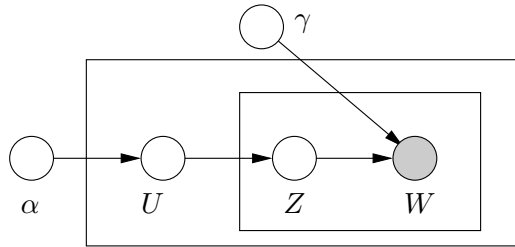
of phonemes. Large-scale systems are built by composing elementary HMMs into larger graphical models.

The graphical model shown in Figure 2.6(a) is a *coupled HMM*, in which two chains of state variables are coupled via links between the chains; this model is appropriate for fusing pairs of data streams such as audio and lip-reading data in speech recognition [203]. Figure 2.6(b) shows an HMM variant in which the state-dependent observation distribution is a finite mixture model. This variant is widely used in speech recognition systems [114].



**Fig. 2.6.** Extensions of HMMs used in language and speech processing. (a) The coupled HMM. (b) An HMM with mixture-model emissions.

Another model class that is widely studied in language processing are so-called “bag-of-words” models, which are of particular interest for modeling large-scale corpora of text documents. The terminology “bag-of-words” means that the order of words in a document is ignored—i.e., an assumption of exchangeability is made. The goal of such models is often that of finding latent “topics” in the corpus, and using these topics to cluster or classify the documents. An example “bag-of-words” model is the *latent Dirichlet allocation* model [29], in which a topic defines a probability distribution on words, and a document defines a probability distribution on topics. In particular, as shown in Figure 2.7, each document in a corpus is assumed to be generated by sampling a Dirichlet variable with hyperparameter  $\alpha$ , and then repeatedly selecting a topic according to these Dirichlet probabilities, and choosing a word



**Fig. 2.7.** Graphical illustration of the latent Dirichlet allocation model. The variable  $U$ , which is distributed as Dirichlet with parameter  $\alpha$ , specifies the parameter for the multinomial “topic” variable  $Z$ . The “word” variable  $W$  is also multinomial conditioned on  $Z$ , with  $\gamma$  specifying the word probabilities associated with each topic. The rectangles, known as *plates*, denote conditionally-independent replications of the random variables inside the rectangle.

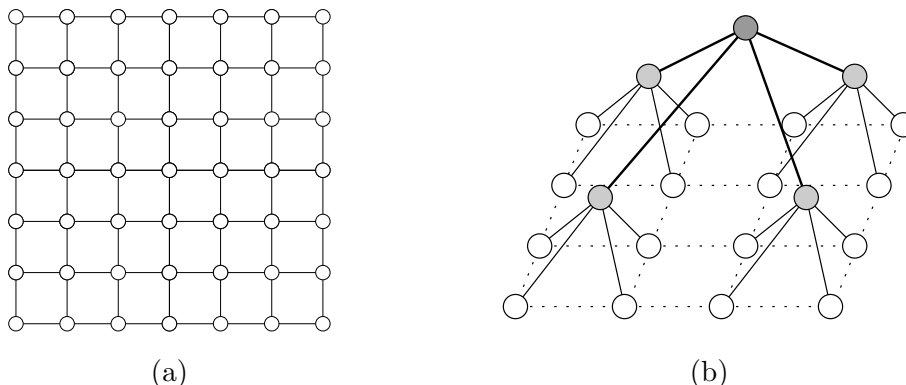
from the distribution associated with the selected topic.<sup>2</sup>

#### 2.4.6 Image processing and spatial statistics

For several decades, undirected graphical models (also known as Markov random fields) have played an important role in image processing [e.g., 258, 104, 55, 86], as well as in spatial statistics more generally [23, 24, 197, 25]. For modeling an image, the simplest use of a Markov random field model is in the pixel domain, where each pixel in the image is associated with a vertex in an underlying graph. More structured models are based on feature vectors at each spatial location, where each feature could be a linear multiscale filter (e.g., a wavelet), or a more complicated nonlinear operator.

For image modeling, one very natural choice of graphical structure is a 2D lattice, such as the 4-nearest neighbor variety shown in Figure 2.8(a). The potential functions on the edges between adjacent pixels (or more generally, features) are typically chosen to enforce local smoothness conditions. Various tasks in image processing, including denoising, segmentation, and super-resolution, require solving an inference problem on such a Markov random field. However, exact inference

<sup>2</sup>This model is discussed in more detail in Example 7 of Section 3.3.



**Fig. 2.8.** (a) The 4-nearest neighbor lattice model in 2D is often used for image modeling. (b) A multiscale quad tree approximation to a 2D lattice model [257]. Nodes in the original lattice (drawn in white) lie at the finest scale of the tree. The middle and top scales of the tree consist of auxiliary nodes (drawn in gray), introduced to model the fine scale behavior.

for large-scale lattice models is intractable, which necessitates the use of approximate methods. Markov chain Monte Carlo methods are often used [86], but they can be too slow and computationally intensive for many applications. More recently, message-passing algorithms such as sum-product and tree-reweighted max-product have become popular as an approximate inference method for image processing and computer vision problems [e.g., 81, 82, 134, 165, 222].

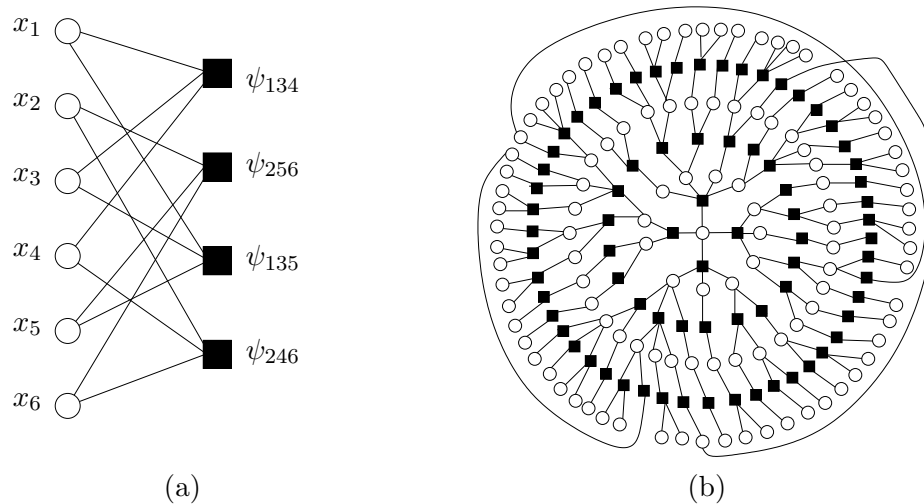
An alternative strategy is to sidestep the intractability of the lattice model by replacing it with a simpler—albeit approximate—model. For instance, multiscale quad trees, such as that illustrated in Figure 2.8(b), can be used to approximate lattice models [257]. The advantage of such a multiscale model is in permitting the application of efficient tree algorithms to perform exact inference. The trade-off is that the model is imperfect, and can introduce artifacts into image reconstructions.

#### 2.4.7 Error-correcting coding

A central problem in communication theory is that of transmitting information, represented as a sequence of bits, from one point to another. Examples include transmission from a personal computer over a net-

work, or from a satellite to a ground position. If the communication channel is noisy, then some of the transmitted bits may be corrupted. In order to combat this noisiness, a natural strategy is to add redundancy to the transmitted bits, thereby defining codewords. In principle, this coding strategy allows the transmission to be decoded perfectly even in the presence of some number of errors.

Many of the best codes in use today, including turbo codes and low-density parity check codes [e.g., 84, 164], are based on graphical models. Figure 2.9(a) provides an illustration of a very small parity check code,



**Fig. 2.9.** (a) A factor graph representation of a parity check code of length  $m = 6$ . Circular nodes  $\circ$  on the left represent bits in the code, whereas black squares  $\blacksquare$  on the right represent the associated factors, or parity checks. This particular code is a  $(2, 3)$  code, since each bit is connected to two parity variables, and each parity relation involves three bits. (b) A large factor graph with a locally “tree-like” structure. Random constructions of factor graphs on  $m$  vertices with bounded degree have cycles of typical length  $\asymp \log m$ ; this tree-like property is essential to the success of the sum-product algorithm for approximate decoding [196, 166].

represented here in the factor graph formalism [139]. (A somewhat larger code is shown in Figure 2.9(b)). The six white nodes on the left represent the bits that comprise the codewords (i.e., binary sequences of length six); each of the four gray nodes on the right corresponds to

a factor  $\psi_{stu}$  that represents the parity of the triple  $\{x_s, x_t, x_u\}$ . This parity relation, expressed mathematically as  $x_s \oplus x_t \oplus x_u \equiv z_{stu}$  in modulo two arithmetic, can be represented as an undirected graphical model using a compatibility function of the form

$$\psi_{stu}(x_s, x_t, x_u) := \begin{cases} 1 & \text{if } x_s \oplus x_t \oplus x_u = 1 \\ 0 & \text{otherwise.} \end{cases}$$

For the code shown in Figure 2.9, the parity checks range over the set of triples  $\{1, 3, 4\}$ ,  $\{1, 3, 5\}$ ,  $\{2, 4, 6\}$  and  $\{2, 5, 6\}$ .

The decoding problem entails estimating which codeword was transmitted on the basis of a vector  $y$  of noisy observations. With the specification of a model for channel noise, this decoding problem can be cast as an inference problem. Depending on the loss function, optimal decoding is based either on computing the marginal probability  $p(x_s = 1|y)$  at each node, or computing the most likely codeword (i.e., the mode of the posterior). For the simple code of Figure 2.9(a), optimal decoding is easily achievable via the junction tree algorithm. Of interest in many applications, however, are much larger codes in which the number of bits is easily several thousand. The graphs underlying these codes are not of low treewidth, so that the junction tree algorithm is not viable. Moreover, MCMC algorithms have not been deployed successfully in this domain.

For many graphical codes, the most successful decoder is based on applying the sum-product algorithm, described in Section 2.6. Since the graphical models defining good codes invariably have cycles, the sum-product algorithm is not guaranteed to compute the correct marginals, nor even to converge. Nonetheless, the behavior of this approximate decoding algorithm is remarkably good for a large class of codes. The behavior of sum-product algorithm is well understood in the asymptotic limit (as the code length  $m$  goes to infinity), where martingale arguments can be used to prove concentration results [195, 157]. For intermediate code lengths, in contrast, its behavior is not as well understood.

## 2.5 Exact inference algorithms

In this section, we turn to a description of the basic exact inference algorithms for graphical models. In computing a marginal probability, we must sum or integrate the joint probability distribution over one or more variables. We can perform this computation as a sequence of operations by choosing a specific ordering of the variables (and making an appeal to Fubini’s theorem). Recall that for either directed or undirected graphical models, the joint probability is a factored expression over subsets of the variables. Consequently, we can make use of the distributive law to move individual sums or integrals across factors that do not involve the variables being summed or integrated over. The phrase “exact inference” refers to the (essentially symbolic) problem of organizing this sequential computation, including managing the intermediate factors that arise. Assuming that each individual sum or integral is performed exactly, then the overall algorithm yields an exact numerical result.

To obtain the marginal probability of a single variable,  $p(x_s) = \mathbb{P}[X_s = x_s]$ , it suffices to choose a specific ordering of the remaining variables and to “eliminate”—that is, sum or integrate out—variables according to that order. Repeating this operation for each individual variable would yield the full set of marginals; this approach, however, is wasteful because it neglects to share intermediate terms in the individual computations. The sum-product and junction tree algorithms are essentially dynamic programming algorithms based on a calculus for sharing intermediate terms. The algorithms involve “message-passing” operations on graphs, where the messages are exactly these shared intermediate terms. Upon convergence of the algorithms, we obtain marginal probabilities for all cliques of the original graph.

Both directed and undirected graphical models involve factorized expressions for joint probabilities, and it should come as no surprise that exact inference algorithms treat them in an essentially identical manner. Indeed, to permit a simple unified treatment of inference algorithms, it is convenient to convert directed models to undirected models and to work exclusively within the undirected formalism. We do this by observing that the factors in (2.1) are not necessarily de-

fined on cliques, since the parents of a given vertex are not necessarily connected. We thus transform a directed graph to an undirected *moral graph*, in which all parents of each child are linked, and all edges are converted to undirected edges. On the moral graph, the factors in (2.1) are all defined on cliques, and (2.1) is a special case of the undirected representation in (2.2). Throughout the rest of the paper, we assume that this transformation has been carried out.

### 2.5.1 Message-passing on trees

We now turn to a description of message-passing algorithms for exact inference on trees. Our treatment is brief; further details can be found in various sources [2, 62, 139, 122, 151]. We begin by observing that the cliques of a tree-structured graph  $T = (V, E(T))$  are simply the individual nodes and edges. As a consequence, any tree-structured graphical model has the following factorization:

$$p(x) = \frac{1}{Z} \prod_{s \in V} \psi_s(x_s) \prod_{(s,t) \in E(T)} \psi_{st}(x_s, x_t). \quad (2.4)$$

Here we describe how the sum-product algorithm computes the marginal distribution

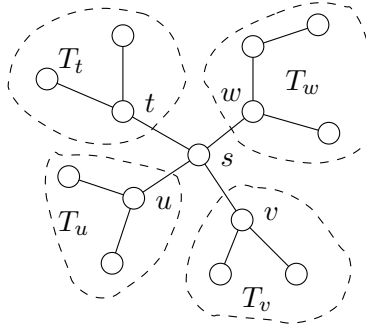
$$\mu_s(x_s) := \sum_{\{x' \mid x'_s = x_s\}} p(x) \quad (2.5)$$

for every node of a tree-structured graph. We will focus in detail on the case of discrete random variables, with the understanding that the computations carry over (at least in principle) to the continuous case by replacing sums with integrals.

**Sum-product algorithm:** The sum-product algorithm is a form of non-serial dynamic programming [17], which generalizes the usual serial form of deterministic dynamic programming [18] to arbitrary tree-structured graphs. The essential principle underlying dynamic programming (DP) is that of divide and conquer: we solve a large problem by breaking it down into a sequence of simpler problems. In the context of graphical models, the tree itself provides a natural way to break down the problem.



For an arbitrary  $s \in V$ , consider the set of its neighbors  $N(s) = \{u \in V \mid (s, u) \in E\}$ . For each  $u \in N(s)$ , let  $T_u = (V_u, E_u)$  be the subgraph formed by the set of nodes (and edges joining them) that can be reached from  $u$  by paths that *do not* pass through node  $s$ . The key property of a tree is that each such subgraph  $T_u$  is again a tree, and  $T_u$  and  $T_v$  are vertex-disjoint for  $u \neq v$ . In this way, each vertex  $u \in N(s)$  can be viewed as the root of a subtree  $T_u$ , as illustrated in Figure 2.10.



**Fig. 2.10.** Decomposition of a tree, rooted at node  $s$ , into subtrees. Each neighbor (e.g.,  $u$ ) of node  $s$  is the root of a subtree (e.g.,  $T_u$ ). Subtrees  $T_u$  and  $T_v$ , for  $t \neq u$ , are disconnected when node  $s$  is removed from the graph.

For each subtree  $T_t$ , we define  $x_{V_t} := \{x_u \mid u \in V_t\}$ . Now consider the collection of terms in equation (2.4) associated with vertices or edges in  $T_t$ . We collect all of these terms into the following product:

$$p(x_{V_t}; T_t) \propto \prod_{u \in V_t} \psi_u(x_u) \prod_{(u,v) \in E_t} \psi_{uv}(x_u, x_v). \quad (2.6)$$

With this notation, the conditional independence properties of a tree allow the computation of the marginal at node  $s$  to be broken down into a product of subproblems, one for each of the subtrees in the set  $\{T_t, t \in N(s)\}$ , in the following way:

$$\mu_s(x_s) = \kappa \psi_s(x_s) \prod_{t \in N(s)} M_{ts}^*(x_s) \quad (2.7a)$$

$$M_{ts}^*(x_s) := \sum_{x'_{V_t}} \psi_{st}(x_s, x'_t) p(x'_{V_t}; T_t) \quad (2.7b)$$

In these equations,  $\kappa$  denotes a positive constant chosen to ensure that  $\mu_s$  normalizes properly. For fixed  $x_s$ , the subproblem defining  $M_{ts}^*(x_s)$  is again a tree-structured summation, albeit involving a subtree  $T_t$  smaller than the original tree  $T$ . Therefore, it too can be broken down recursively in a similar fashion. In this way, the marginal at node  $s$  can be computed by a series of recursive updates.

Rather than applying the procedure described above to each node separately, the *sum-product algorithm* computes the marginals for all nodes simultaneously and in parallel. At each iteration, each node  $t$  passes a “message” to each of its neighbors  $u \in N(t)$ . This message, which we denote by  $M_{tu}(x_u)$ , is a function of the possible states  $x_u \in \mathcal{X}_u$  (i.e., a vector of length  $|\mathcal{X}_u|$  for discrete random variables). On the full graph, there are a total of  $2|E|$  messages, one for each direction of each edge. This full collection of messages is updated, typically in parallel, according to the following recursion:

$$M_{ts}(x_s) \leftarrow \kappa \sum_{x'_t} \left\{ \psi_{st}(x_s, x'_t) \psi_t(x'_t) \prod_{u \in N(t)/s} M_{ut}(x'_t) \right\}, \quad (2.8)$$

where  $\kappa > 0$  is a normalization constant. It can be shown [188] that for tree-structured graphs, iterates generated by the update (2.8) will converge to a unique fixed point  $M^* = \{M_{st}^*, M_{ts}^*, (s, t) \in E\}$  after a finite number of iterations. Moreover, component  $M_{ts}^*$  of this fixed point is precisely equal, up to a normalization constant, to the subproblem defined in equation (2.7b), which justifies our abuse of notation post hoc. Since the fixed point  $M^*$  specifies the solution to all of the subproblems, the marginal  $\mu_s$  at every node  $s \in V$  can be computed easily via equation (2.7a).

**Max-product algorithm:** Suppose that the summation in the update (2.8) is replaced by a maximization. The resulting *max-product algorithm* solves the problem of finding a mode of a tree-structured distribution  $p(x)$ . In this sense, it represents a generalization of the Viterbi algorithm [79] from chains to arbitrary tree-structured graphs. More specifically, the max-product updates will converge to another unique fixed point  $M^*$ —distinct, of course, from the sum-product fixed point. This fixed point can be used to compute the *max-marginal*

$\nu_s(x_s) := \max_{\{x' \mid x'_s = x_s\}} p(x')$  at each node of the graph, via the analog of equation (2.6). Given these max-marginals, it is straightforward to compute a mode  $\hat{x} \in \arg \max_x p(x)$  of the distribution; see the papers [62, 239] for further details. More generally, updates of this form apply to arbitrary *commutative semirings* on tree-structured graphs [232, 211, 62, 2]. The pairs “sum-product” and “max-product” are two particular examples of such an algebraic structure.

### 2.5.2 Junction tree representation

We have seen that inference problems on trees can be solved exactly by recursive message-passing algorithms. Given a graph with cycles, a natural idea is to cluster its nodes so as to form a *clique tree*—that is, an acyclic graph whose nodes are formed by the maximal cliques of  $G$ . Having done so, it is tempting to simply apply a standard algorithm for inference on trees. However, the clique tree must satisfy an additional restriction so as to ensure correctness of these computations. In particular, since a given vertex  $s \in V$  may appear in multiple cliques (say  $C_1$  and  $C_2$ ), what is required is a mechanism for enforcing consistency among the different appearances of the random variable  $x_s$ . It turns out that the following property is necessary and sufficient to enforce such consistency:

**Definition 1.** A clique tree has the *running intersection property* if for any two clique nodes  $C_1$  and  $C_2$ , all nodes on the unique path joining them contain the intersection  $C_1 \cap C_2$ . A clique tree with this property is known as a *junction tree*.

For what type of graphs can one build junction trees? An important result in graph theory asserts that a graph  $G$  has a junction tree if and only if it is *triangulated*.<sup>3</sup> (See Lauritzen [150] for a proof.) This result underlies the *junction tree algorithm* [151] for exact inference on arbitrary graphs:

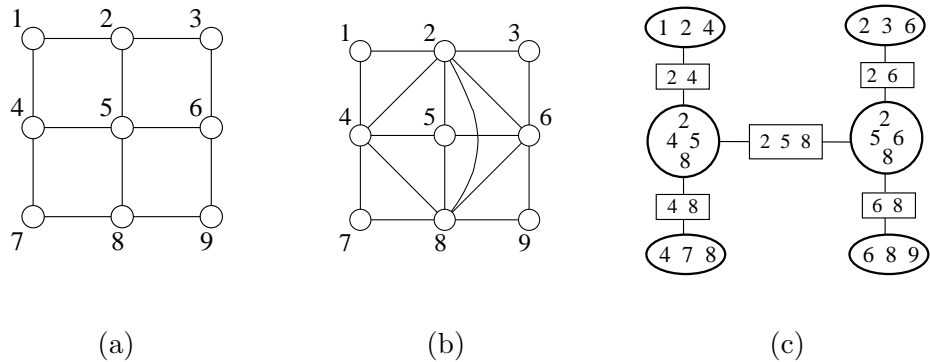
- (1) Given a graph with cycles  $G$ , triangulate it by adding edges as necessary.

---

<sup>3</sup>A triangulated graph is a graph in which every cycle of length four or longer has a chord.

- (2) Form a junction tree associated with the triangulated graph  $\tilde{G}$ .
- (3) Run a tree inference algorithm on the junction tree.

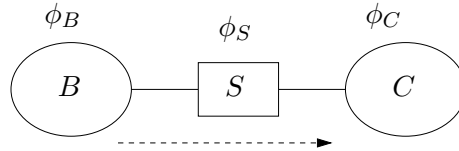
**Example 1 (Junction tree).** To illustrate the junction tree construction, consider the  $3 \times 3$  grid shown in Figure 2.11(a). The first step is to form a triangulated version  $\tilde{G}$ , as shown in Figure 2.11(b). Note



**Fig. 2.11.** Illustration of junction tree construction. (a) Original graph is a  $3 \times 3$  grid. (b) Triangulated version of original graph. Note the two 4-cliques in the middle. (c) Corresponding junction tree for triangulated graph in (b), with maximal cliques depicted within ellipses, and separator sets within rectangles.

that the graph would not be triangulated if the additional edge joining nodes 2 and 8 were not present. Without this edge, the 4-cycle  $(2 - 4 - 8 - 6 - 2)$  would lack a chord. As a result of this additional edge, the junction tree has two 4-cliques in the middle, as shown in Figure 2.11(c). These cliques grow larger quickly as the grid size is increased. ♣

In principle, the inference in the third step of the junction tree algorithm can be performed over an arbitrary commutative semiring (as mentioned in our earlier discussion of tree algorithms). We refer the reader to Dawid [62] for an extensive discussion of the max-product version of the junction tree algorithm. For concreteness, we limit our discussion here to the sum-product version of junction tree updates. There is an elegant way to express the basic algebraic operations in



**Fig. 2.12.** A message-passing operation between cliques  $B$  and  $C$  via the separator set  $S$ .

a junction tree inference algorithm that involves introducing potential functions not only on the cliques in the junction tree, but also on the *separators* in the junction tree—the intersections of cliques that are adjacent in the junction tree (the rectangles in Figure 2.11). Let  $\phi_C(x_C)$  denote a potential on a clique  $C$ , and let  $\phi_S(x_S)$  denote a potential on a separator  $S$ . We initialize the clique potentials by assigning each compatibility function in the original graph to (exactly) one clique potential and taking the product over these compatibility functions. The separator potentials are initialized to unity. Given this set-up, the basic message-passing step of the junction tree algorithm can be written in the following form:

$$\phi_S(x_S) \leftarrow \sum_{x_{B \setminus S}} \phi_B(x_B) \quad (2.9a)$$

$$\phi_C(x_C) \leftarrow \frac{\phi_S^*(x_S)}{\phi_S(x_S)} \phi_C(x_C), \quad (2.9b)$$

where in the continuous case the summation is replaced by a suitable integral. We refer to this pair of operations as “passing a message from clique  $B$  to clique  $C$ ” (see Figure 2.12). It can be verified that if a message is passed from  $B$  to  $C$ , and subsequently from  $C$  to  $B$ , then the resulting clique potentials are consistent with each other; that is, they agree with respect to the vertices  $S$ .

After a round of message passing on the junction tree, it can be shown that the clique potentials are proportional to marginal probabilities throughout the junction tree. Specifically, letting  $\mu_C(x_C)$  denote the marginal probability of  $x_C$ , we have  $\mu_C(x_C) \propto \phi_C^*(x_C)$  for all cliques  $C$ . This equivalence can be established by a suitable generalization of the proof of correctness of the sum-product algorithm presented previ-

ously (see also Lauritzen [150]). Note that achieving local consistency between pairs of cliques is obviously a necessary condition if the clique potentials are to be proportional to marginal probabilities. Moreover, the significance of the running intersection property is now apparent; namely, it ensures that local consistency implies global consistency.

An important by-product of the junction tree algorithm is an alternative representation of a distribution  $p$ . Let  $\mathcal{C}$  denote the set of all maximal cliques in  $\tilde{G}$  (i.e., nodes in the junction tree), and let  $\mathcal{S}$  represent the set of all separator sets (i.e., intersections between cliques that are adjacent in the junction tree). Note that a given separator set  $S$  may occur multiple times in the junction tree. For each separator set  $S \in \mathcal{S}$ , let  $d(S)$  denote the number of maximal cliques to which it is adjacent. The junction tree framework guarantees that the distribution  $p$  factorizes in the form

$$p(x) = \frac{\prod_{C \in \mathcal{C}} \mu_C(x_C)}{\prod_{S \in \mathcal{S}} [\mu_S(x_S)]^{d(S)-1}}, \quad (2.10)$$

where  $\mu_C$  and  $\mu_S$  are the marginal distributions over the cliques and separator sets respectively. Observe that unlike the representation of equation (2.2), the decomposition of equation (2.10) is directly in terms of marginal distributions, and does not require a normalization constant (i.e.,  $Z = 1$ ).

**Example 2 (Markov chain).** Consider the Markov chain  $p(x_1, x_2, x_3) = p(x_1) p(x_2 | x_1) p(x_3 | x_2)$ . The cliques in a graphical model representation are  $\{1, 2\}$  and  $\{2, 3\}$ , with separator  $\{2\}$ . Clearly the distribution cannot be written as the product of marginals involving only the cliques. It can, however, be written in terms of marginals if we include the separator:

$$p(x_1, x_2, x_3) = \frac{p(x_1, x_2) p(x_2, x_3)}{p(x_2)}.$$

Moreover, it can be easily verified that these marginals result from a single application of equation (2.9), given the initialization  $\phi_{\{1,2\}}(x_1, x_2) = p(x_1)p(x_2 | x_1)$  and  $\phi_{\{2,3\}}(x_2, x_3) = p(x_3 | x_2)$ . ♣

To anticipate part of our development in the sequel, it is helpful to consider the following “inverse” perspective on the junction tree

representation. Suppose that we are given a set of functions  $\tau_C(x_C)$  and  $\tau_S(x_S)$  associated with the cliques and separator sets in the junction tree. What conditions are necessary to ensure that these functions are valid marginals for some distribution? Suppose that the functions  $\{\tau_S, \tau_C\}$  are *locally consistent* in the following sense:

$$\sum_{x_S} \tau_S(x_S) = 1 \quad \text{normalization} \quad (2.11a)$$

$$\sum_{\{x'_C \mid x'_S=x_S\}} \tau_C(x'_C) = \tau_S(x_S) \quad \text{marginalization} \quad (2.11b)$$

The essence of the junction tree theory described above is that such local consistency is both necessary and sufficient to ensure that these functions are valid marginals for some distribution. For the sake of future reference, we state this result in the following:

**Proposition 1.** *A candidate set of local marginals  $\{\tau_S, \tau_C\}$  on the separator sets and cliques of a junction tree is globally consistent if and only if it is locally consistent in the sense of equation (2.11). Moreover, any such locally consistent quantities are the marginals of the probability distribution defined by equation (2.10).*

This particular consequence of the junction tree representation will play a fundamental role in our development in the sequel.

Finally, let us turn to the key issue of the computational complexity of the junction tree algorithm. Inspection of equation (2.9) reveals that the computational costs grow exponentially in the size of the maximal clique in the junction tree. Clearly then, it is of interest to control the size of this clique. The size of the maximal clique over all possible triangulations of a graph is an important graph-theoretic quantity known as the *treewidth* of the graph.<sup>4</sup> Thus, the complexity of the junction tree algorithm is exponential in the treewidth.

For certain classes of graphs, including chains and trees, the treewidth is small and the junction tree algorithm provides an effective solution to inference problems. Such families include many well-known graphical model architectures, and the junction tree algorithm

---

<sup>4</sup>To be more precise, the treewidth is one less than the size of this largest clique [see 30].

subsumes the classical recursive algorithms, including the pruning and peeling algorithms from computational genetics [76], the forward-backward algorithms for hidden Markov models [192], and the Kalman filtering-smoothing algorithms for state-space models [123]. On the other hand, there are many graphical models, including several of the examples treated in Section 2.4, for which the treewidth is infeasibly large. Coping with such models requires leaving behind the junction tree framework, and turning to approximate inference algorithms.

## 2.6 Message-passing algorithms for approximate inference

It is the goal of the remainder of the paper to develop a general theoretical framework for understanding and analyzing *variational methods* for computing approximations to marginal distributions or the partition function, as well as for solving integer programming problems. Doing so requires mathematical background on convex analysis and exponential families, which we provide starting in Section 3. Historically, many of these algorithms have been developed without this background, but rather via intuition or on the basis of analogies to exact or Monte Carlo algorithms. In this section, we give a high-level description of this flavor for two variational inference algorithms, with the goal of highlighting their simple and intuitive nature.

The first variational algorithm that we consider is a so-called “loopy” form of the sum-product algorithm, also referred to as the *belief propagation* algorithm. Recall that the sum-product algorithm is an exact inference algorithm for trees. From an algorithmic point of view, however, there is nothing to prevent one from running the procedure on a graph with cycles. More specifically, the message updates (2.8) can be applied at a given node while ignoring the presence of cycles—essentially pretending that any given node is embedded in a tree. Intuitively, such an algorithm might be expected to work well if the graph is sparse, such that the effect of messages propagating around cycles is appropriately diminished, or if suitable symmetries are present. As discussed in Section 2.4, this algorithm is in fact successfully used in various applications. Also, an analogous form of the max-product algorithm is used for computing approximate modes in



graphical models with cycles.

A second variational algorithm is the so-called *naive mean field* algorithm. For concreteness, here we describe this algorithm in application to the Ising model of statistical physics. The Ising model is a Markov random field involving a binary random vector  $x \in \{0, 1\}^m$ , in which pairs of adjacent nodes are coupled with a weight  $\theta_{st}$ , and each node has an observation weight  $\theta_s$ . (See Example 3 of Section 3.3 for a more detailed description of this model.) Consider now the Gibbs sampler for such a model. The basic step of a Gibbs sampler is to choose a node  $s \in V$  randomly, and then to update the state of the associated random variable according to the conditional probability with neighboring states fixed. More precisely, denoting by  $N(s)$  the neighbors of a node  $s \in V$ , and letting  $x_{N(s)}^{(p)}$  denote the state of the neighbors of  $s$  at iteration  $p$ , the Gibbs update for  $x_s$  takes the following form:

$$x_s^{(p+1)} = \begin{cases} 1 & \text{if } u \leq \{1 + \exp[-(\theta_s + \sum_{t \in N(s)} \theta_{st} x_t^{(p)})]\}^{-1} \\ 0 & \text{otherwise} \end{cases}, \quad (2.12)$$

where  $u$  is a sample from a uniform distribution  $\mathcal{U}(0, 1)$ .

In a dense graph, such that the cardinality of  $N(s)$  is large, we might attempt to invoke a law of large numbers or some other concentration result for  $\sum_{t \in N(s)} \theta_{st} x_t^{(p)}$ . To the extent that such sums are concentrated, it might make sense to replace sample values with expectations. That is, letting  $\mu_s$  denote an estimate of the marginal probability  $\mathbb{P}[x_s = 1]$  at each vertex  $s \in V$ , we might consider the following averaged version of equation (2.12):

$$\mu_s \leftarrow \left\{ 1 + \exp \left[ - \left( \theta_s + \sum_{t \in N(s)} \theta_{st} \mu_t \right) \right] \right\}^{-1}. \quad (2.13)$$

Thus, rather than flipping the random variable  $x_s$  with a probability that depends on the state of its neighbors, we update a parameter  $\mu_s$  deterministically that depends on the corresponding parameters at its neighbors. Equation (2.13) defines the naive mean field algorithm for the Ising model. As with the sum-product algorithm, the mean field algorithm can be viewed as a message-passing algorithm, in which the right-hand side of (2.13) represents the “message” arriving at vertex  $s$ .

At first sight, message-passing algorithms of this nature might seem rather mysterious, and do raise some questions. Do the updates have fixed points? Do the updates converge? What is the relation between the fixed points and the exact quantities? The goal of the remainder of this paper is to shed some light on such issues. Ultimately, we will see that a broad class of message-passing algorithms, including the mean field updates, the sum-product and max-product algorithms, as well as various extensions of these methods, can all be understood as solving either exact or approximate versions of variational problems. Exponential families and convex analysis, which are the subject of the following section, provide the appropriate framework in which to develop these variational principles in a unified manner.

# 3

---

## Graphical models as exponential families

---

In this section, we describe how many graphical models are naturally viewed as exponential families, a broad class of distributions that have been extensively studied in the statistics literature [5, 10, 40, 71]. Taking the perspective of exponential families illuminates some fundamental connections between inference algorithms and the theory of convex analysis [35, 109, 198]. More specifically, as we shall see, various types of inference problems in graphical models can be understood in terms of mappings between mean parameters and canonical parameters.

### 3.1 Exponential representations via maximum entropy

One way in which to motivate exponential family representations of graphical models is through the principle of maximum entropy [120, 259]. Here we describe a particularly simple version, applicable to a scalar random variable  $X$ , that provides helpful intuition for our subsequent development. Suppose that given  $n$  independent and identically distributed (i.i.d.) observations  $X^1, \dots, X^n$ , we compute the empirical

expectations of certain functions—namely, the quantities

$$\widehat{\mu}_\alpha := \frac{1}{n} \sum_{i=1}^n \phi_\alpha(X^i), \quad \text{for all } \alpha \in \mathcal{I}, \quad (3.1)$$

where each  $\alpha$  in some set  $\mathcal{I}$  indexes a function  $\phi_\alpha : \mathcal{X} \rightarrow \mathbb{R}$ . For example, supposing that  $d = |\mathcal{I}| = 2$  and setting  $\phi_1(x) = x$  and  $\phi_2(x) = x^2$  corresponds to observing empirical versions of the first and second moments of the random variable  $X$ . Based on the empirical expectations  $\{\widehat{\mu}_\alpha, \alpha \in \mathcal{I}\}$ , our goal is to infer a full probability distribution over the random variable  $X$ . In particular, we represent the probability distribution as a density  $p$  absolutely continuous with respect to some measure  $\nu$ . This base measure  $\nu$  might be the counting measure on  $\{0, 1, \dots, r-1\}$ , in which case  $p(\cdot)$  is a probability mass function; alternatively, for a continuous random vector, the base measure  $\nu$  could be the ordinary Lebesgue measure on  $\mathbb{R}^m$ .

We say that a distribution  $p$  is consistent with the data if

$$\mathbb{E}_p[\phi_\alpha(X)] := \int_{\mathcal{X}} \phi_\alpha(x) p(x) \nu(dx) = \widehat{\mu}_\alpha \quad \text{for all } \alpha \in \mathcal{I}.$$

In words, the expectations  $\mathbb{E}_p[\phi_\alpha(X)]$  under the distribution  $p$  are matched to the expectations under the empirical distribution. An important observation is that generically, this problem is underdetermined, in that there are many distributions  $p$  that are consistent with the observations, so that we need a principle for choosing among them.

In order to develop such a principle, we begin by defining a functional of the density  $p$ , known as the *Shannon entropy*, via

$$H(p) = - \int_{\mathcal{X}} (\log p(x)) p(x) \nu(dx), \quad (3.2)$$

The principle of maximum entropy is to choose, from among the distributions consistent with the data, the distribution  $p^*$  whose Shannon entropy is maximal. More formally, letting  $\mathcal{P}$  be the set of all probability distributions over the random variable  $X$ , the maximum entropy solution  $p^*$  is given by the solution to the following constrained opti-

mization problem:

$$p^* := \arg \max_{p \in \mathcal{P}} H(p) \quad \text{subject to } \mathbb{E}_p[\phi_\alpha(X)] = \hat{\mu}_\alpha \text{ for all } \alpha \in \mathcal{I}. \quad (3.3)$$

One interpretation of this principle is as choosing the distribution with maximal uncertainty, as measured by the entropy functional (3.2), while remaining faithful to the data. Presuming that problem (3.3) is feasible (and under certain technical conditions to be explored in the sequel), it can be shown—by calculus of variations in the general continuous case, and by ordinary calculus in the discrete case—that the optimal solution  $p^*$  takes the form

$$p_\theta(x) \propto \exp \left\{ \sum_{\alpha \in \mathcal{I}} \theta_\alpha \phi_\alpha(x) \right\}, \quad (3.4)$$

where  $\theta \in \mathbb{R}^d$  represents a parameterization of the distribution in exponential family form. From this maximum entropy perspective, the parameters  $\theta$  have a very specific interpretation as the Lagrange multipliers associated with the constraints specified by the empirical moments  $\hat{\mu}$ . We explore this connection in much more depth in the following sections.

## 3.2 Basics of exponential families

With this motivating example in mind, let us now set up the framework of exponential families in more precise terms and greater generality. At a high level, an exponential family is a parameterized family of densities, taken with respect to some underlying measure  $\nu$ .

Let  $\phi = \{\phi_\alpha \mid \alpha \in \mathcal{I}\}$  be a collection of functions  $\phi_\alpha : \mathcal{X}^m \rightarrow \mathbb{R}$ , known either as *potential functions* or *sufficient statistics*. Here  $\mathcal{I}$  is an index set with  $d = |\mathcal{I}|$  elements to be specified, so that  $\phi$  can be viewed as a vector-valued mapping from  $\mathcal{X}^m$  to  $\mathbb{R}^d$ . With the vector of sufficient statistics  $\phi$ , we associate a vector  $\theta = \{\theta_\alpha \mid \alpha \in \mathcal{I}\}$  of *canonical* or *exponential* parameters. For each fixed  $x \in \mathcal{X}^m$ , we use  $\langle \theta, \phi(x) \rangle$  to denote the Euclidean inner product in  $\mathbb{R}^d$  of the two vectors  $\theta$  and  $\phi(x)$ .

With this notation, the *exponential family* associated with  $\phi$  con-

sists of the following parameterized collection of density functions

$$p_\theta(x) = \exp\{\langle \theta, \phi(x) \rangle - A(\theta)\}, \quad (3.5)$$

taken with respect<sup>1</sup> to  $d\nu$ . The quantity  $A$ , known as the *log partition function* or *cumulant function*, is defined by the integral:

$$A(\theta) = \log \int_{\mathcal{X}^m} \exp\langle \theta, \phi(x) \rangle \nu(dx). \quad (3.6)$$

Presuming that the integral is finite, this definition ensures that  $p_\theta$  is properly normalized (i.e.,  $\int_{\mathcal{X}^m} p_\theta(x) \nu(dx) = 1$ ). With the set of potentials  $\phi$  fixed, each parameter vector  $\theta$  indexes a particular member  $p_\theta$  of the family. The canonical parameters  $\theta$  of interest belong to the set

$$\Omega := \{\theta \in \mathbb{R}^d \mid A(\theta) < +\infty\}. \quad (3.7)$$

We will see shortly that  $A$  is a convex function of  $\theta$ , which in turn implies that  $\Omega$  must be a convex set. The log partition function  $A$  plays a prominent role in this paper.

The following notions will be important in subsequent development:

**Regular families:** An exponential family for which the domain  $\Omega$  of equation (3.7) is an open set is known as a *regular* family. Although there do exist exponential families for which  $\Omega$  is closed (for instance, see Brown [40]), herein we restrict our attention to regular exponential families.

**Minimal:** It is typical to define an exponential family with a collection of functions  $\phi = \{\phi_\alpha\}$  for which there does not exist a non-zero vector  $a \in \mathbb{R}^d$  such that the linear combination

$$\langle a, \phi(x) \rangle = \sum_{\alpha \in \mathcal{I}} a_\alpha \phi_\alpha(x)$$

is equal to a constant ( $\nu$ -almost everywhere.). This condition gives rise to a so-called *minimal representation*, in which there is a unique parameter vector  $\theta$  associated with each distribution.

---

<sup>1</sup> More precisely, for any measurable set  $S$ , we have  $\mathbb{P}[X \in S] = \int_S p_\theta(x) \nu(dx)$ .

**Overcomplete:** Instead of a minimal representation, it can be convenient to use a non-minimal or *overcomplete representation*, in which there exists a linear combination  $\langle a, \phi(x) \rangle$  that is equal to a constant ( $\nu$ -almost everywhere). In this case, there exists an entire affine subset of parameter vectors  $\theta$ , each associated with the same distribution. The reader might question the utility of an overcomplete representation. Indeed, from a statistical perspective, it can be undesirable since the identifiability of the parameter vector  $\theta$  is lost. However, this notion of overcompleteness is useful in understanding the sum-product algorithm and its variants (see Section 4).

Table 3.1 provides some examples of well-known scalar exponential families. Observe that all of these families are both regular (since  $\Omega$  is open), and minimal (since the collection of sufficient statistics  $\phi$  do not satisfy any linear relations).

Family	$\mathcal{X}$	$\nu$ $h(\cdot)$	Sufficient statistics $\langle \theta, \phi(x) \rangle$	$A(\theta)$	$\Omega$
Bernoulli	$\{0, 1\}$	Counting	$\theta x$	$\log[1 + \exp(\theta)]$	$\mathbb{R}$
Gaussian Location family	$\mathbb{R}$	Lebesgue $h(x) = \frac{\exp(-x^2/2)}{\sqrt{2\pi}}$	$\theta x$	$\frac{1}{2}\theta^2$	$\mathbb{R}$
Gaussian Location-scale	$\mathbb{R}$	Lebesgue $h(x) = \frac{1}{\sqrt{2\pi}}$	$\theta_1 x + \theta_2 x^2$	$-\frac{\theta_1^2}{4\theta_2} - \frac{1}{2} \log(-2\theta_2)$	$\{(\theta_1, \theta_2) \in \mathbb{R}^2 \mid \theta_2 < 0\}$
Exponential	$(0, +\infty)$	Lebesgue	$\theta x$	$-\log(-\theta)$	$(-\infty, 0)$
Poisson	$\{0, 1, 2, \dots\}$	Counting $h(x) = 1/x!$	$\theta x$	$\exp(\theta)$	$\mathbb{R}$
Beta	$(0, 1)$	Lebesgue	$\theta_1 \log x + \theta_2 \log(1-x)$	$\sum_{i=1}^2 \log \Gamma(\theta_i + 1)$ $-\log \Gamma(\sum_{i=1}^2 (\theta_i + 1))$	$(-1, +\infty)^2$

**Table 3.1.** Several well-known classes of scalar random variables as exponential families. In all cases, the base measure  $\nu$  is either Lebesgue or counting measure, suitably restricted to the space  $\mathcal{X}$ , and (in some cases) modulated by a factor  $h(x)$ . All of these examples are both minimal and regular.

### 3.3 Examples of graphical models in exponential form

The scalar examples in Table 3.1 serve as building blocks for the construction of more complex exponential families for which graphical structure plays a role. Whereas earlier we described graphical models in terms of products of functions, as in equations (2.1) and (2.2), these products become additive decompositions in the exponential family setting. Here we discuss a few well-known examples of graphical models as exponential families. Those readers familiar with such formulations may skip directly to Section 3.4, where we continue our general discussion of exponential families.

**Example 3 (Ising model).** The *Ising model* from statistical physics [11, 116] is a classical example of a graphical model in exponential form. Consider a graph  $G = (V, E)$  and suppose that the random variable  $X_s$  associated with node  $s \in V$  is Bernoulli, say taking the “spin” values  $\{-1, +1\}$ . In the context of statistical physics, these values might represent the orientations of magnets in a field, or the presence/absence of particles in a gas. Moreover, the Ising model and variations on it have been used in image processing and spatial statistics [25, 86, 98], where  $X_s$  might correspond to pixel values in a black-and-white image.

Components  $X_s$  and  $X_t$  of the full random vector  $X$  are allowed to interact directly only if  $s$  and  $t$  are joined by an edge in the graph. This set-up leads to an exponential family consisting of the densities

$$p_\theta(x) = \exp\left\{\sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t - A(\theta)\right\}, \quad (3.8)$$

where the base measure  $\nu$  is the counting measure<sup>2</sup> restricted to  $\{0, 1\}^m$ . Here  $\theta_{st} \in \mathbb{R}$  is the strength of edge  $(s, t)$ , and  $\theta_s \in \mathbb{R}$  is a potential for node  $s$ , which models an “external field” in statistical physics, or a noisy observation in spatial statistics. Strictly speaking, the family of densities (3.8) is more general than the classical Ising model, in which  $\theta_{st}$  is constant for all edges.

<sup>2</sup>Explicitly, for each singleton set  $\{x\}$ , this counting measure is defined by  $\nu(\{x\}) = 1$  if  $x \in \{0, 1\}^m$  and  $\nu(\{x\}) = 0$  otherwise, and extended to arbitrary sets by sub-additivity.



As an exponential family, the Ising index set  $\mathcal{I}$  consists of the union  $V \cup E$ , and the dimension of the family is  $d = m + |E|$ . The log partition function is given by the sum

$$A(\theta) = \log \sum_{x \in \{0,1\}^m} \exp \left\{ \sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t \right\}. \quad (3.9)$$

Since this sum is finite for all choices of  $\theta \in \mathbb{R}^d$ , the domain  $\Omega$  is the full space  $\mathbb{R}^d$ , and the family is regular. Moreover, it is a minimal representation, since there is no non-trivial linear combination of the potentials equal to a constant  $\nu$ -a.e..

The standard Ising model can be generalized in a number of different ways. Although equation (3.8) includes only pairwise interactions, higher-order interactions among the random variables can also be included. For example, in order to include coupling within the 3-clique  $\{s, t, u\}$ , we add a monomial of the form  $x_s x_t x_u$ , with corresponding canonical parameter  $\theta_{stu}$ , to equation (3.8). More generally, to incorporate coupling in  $k$ -cliques, we can add monomials up to order  $k$ , which lead to so-called  $k$ -spin models in the statistical physics literature. At the upper extreme, taking  $k = m$  amounts to connecting all nodes in the graphical model, which allows one to represent any distribution over a binary random vector  $x \in \{0,1\}^m$ .



**Example 4 (Metric labeling and Potts model).** Here we consider another generalization of the Ising model: suppose that the random variable  $X_s$  at each node  $s \in V$  takes values in the discrete space  $\mathcal{X} := \{0, 1, \dots, r-1\}$ , for some integer  $r > 2$ . One interpretation of the states  $j \in \mathcal{X}$  is as labels, for instance defining membership in an image segmentation problem. Each pairing of a node  $s \in V$  and a state  $j \in \mathcal{X}$  yields a sufficient statistic

$$\mathbb{I}_{s;j}(x_s) = \begin{cases} 1 & \text{if } x_s = j \\ 0 & \text{otherwise,} \end{cases} \quad (3.10)$$

with an associated vector  $\theta_s = \{\theta_{s;j}, j = 0, \dots, r-1\}$  of canonical parameters. Moreover, for each edge  $(s, t)$  and pair of values  $(j, k) \in$

$\mathcal{X} \times \mathcal{X}$ , define the sufficient statistics

$$\mathbb{I}_{st;jk}(x_s, x_t) = \begin{cases} 1 & \text{if } x_s = j \text{ and } x_t = k, \\ 0 & \text{otherwise,} \end{cases} \quad (3.11)$$

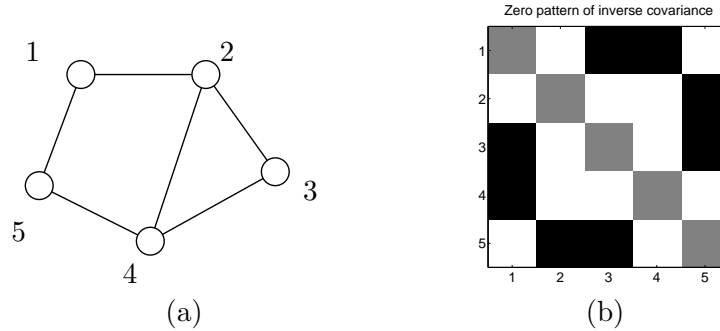
as well as the associated parameter  $\theta_{st;jk} \in \mathbb{R}$ . A special case is the *metric labeling problem*, in which a metric  $d : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$  specifies the parameterization—that is,  $\theta_{st;jk} = -d(j, k)$  for all  $(j, k) \in \mathcal{X} \times \mathcal{X}$ . Consequently, the canonical parameters satisfy the relations  $\theta_{st;kk} = 0$  for all  $k \in \mathcal{X}$ ,  $\theta_{st;jk} < 0$  for all  $j \neq k$ , and satisfy the reversed triangle inequality (i.e.,  $\theta_{st;jl} \geq \theta_{st;jk} + \theta_{st;kl}$  for all triples  $(j, k, \ell)$ ). Another special case is the *Potts model* from statistical physics, in which case  $\theta_{st;kk} = \alpha$  for all  $k \in \mathcal{X}$ , and  $\theta_{st;jk} = \beta$  for all  $j \neq k$ .

Viewed as an exponential family, the chosen collection of sufficient statistics defines an exponential family with dimension  $d = r|V| + r^2|E|$ . Like the Ising model, the log partition function is everywhere finite, so that the family is regular. In contrast to the Ising model, however, the family is overcomplete: indeed, the sufficient statistics satisfy various linear relations—for instance,  $\sum_{j \in \mathcal{X}} \mathbb{I}_{s;j}(x_s) = 1$  for all  $x_s \in \mathcal{X}$ . ♣

We now turn to an important class of graphical models based on continuous random variables:

**Example 5 (Gaussian MRF).** Given an undirected graph  $G$  with vertex set  $V = \{1, \dots, m\}$ , a Gaussian Markov random field [215] consists of a multivariate Gaussian random vector  $(X_1, \dots, X_m)$  that respects the Markov properties of  $G$  (see Section 2.2). It can be represented in exponential form using the collection of sufficient statistics  $\{x_s, x_s^2 \mid s \in V\} \cup \{x_s x_t \mid (s, t) \in E\}$ . We define a  $m$ -vector  $\theta$  of parameters associated with the vector of sufficient statistics  $x = (x_1, \dots, x_m)$ , and a symmetric matrix  $\Theta \in \mathbb{R}^{m \times m}$  associated with the matrix  $xx^T$ . Concretely, the matrix  $\Theta$  is the inverse covariance or precision matrix, and by the Hammersley-Clifford theorem, it has the property that  $\Theta_{st} = 0$  if  $(s, t) \notin E$ , as illustrated in Figure 3.1. Consequently, the dimension of the resulting exponential family is  $d = 2m + |E|$ .

With this set-up, the multivariate Gaussian is an exponential fam-



**Fig. 3.1.** (a) Undirected graphical model on five nodes. (b) For a Gaussian Markov random field, the zero pattern of the inverse covariance or precision matrix respects the graph structure: for any pair  $i \neq j$ , if  $(i, j) \notin E$ , then  $\Theta_{ij} = 0$ .

ily<sup>3</sup> of the form

$$p_{\theta}(x) = \exp \left\{ \langle \theta, x \rangle + \frac{1}{2} \langle \langle \Theta, xx^T \rangle \rangle - A(\theta, \Theta) \right\}, \quad (3.12)$$

where  $\langle \theta, x \rangle := \sum_{i=1}^m \theta_i x_i$  is the Euclidean inner product on  $\mathbb{R}^m$ , and

$$\langle \langle \Theta, xx^T \rangle \rangle := \text{trace}(\Theta xx^T) = \sum_{i=1}^m \sum_{j=1}^m \Theta_{ij} x_i x_j \quad (3.13)$$

is the Frobenius inner product on symmetric matrices. The integral defining  $A(\theta, \Theta)$  is finite only if  $\Theta \prec 0$ , so that

$$\Omega = \{(\theta, \Theta) \in \mathbb{R}^d \mid \Theta \prec 0\}. \quad (3.14)$$

♣

Graphical models are not limited to cases in which the random variables at each node belong to the same exponential family. More generally, we can consider heterogeneous combinations of exponential family members, as illustrated by the following examples.

<sup>3</sup>Our inclusion of the  $\frac{1}{2}$ -factor in the term  $\frac{1}{2} \langle \langle \Theta, xx^T \rangle \rangle$  is for later technical convenience.

**Example 6 (Mixture models).** As shown in Figure 3.2(a), a scalar mixture model has a very simple graphical interpretation. Concretely, in order to form a finite mixture of Gaussians, let  $X_s$  be a multinomial variable, taking values in  $\{0, 1, \dots, r-1\}$ . The role of  $X_s$  is to specify the choice of mixture component, so that our mixture model has  $r$  components in total. As in the Potts model described in Example 4, the distribution of this random variable is exponential family with sufficient statistics  $\{\mathbb{I}_j[x_s], j = 0, \dots, r-1\}$ , and associated canonical parameters  $\{\alpha_{s;0}, \dots, \alpha_{s;r-1}\}$ .

In order to form the finite mixture, conditioned on  $X_s = j$ , we now let  $Y_s$  be conditionally Gaussian with mean and variance  $(\mu_j, \sigma_j^2)$ . Each such conditional distribution can be written in exponential family form in terms of the sufficient statistics  $\{y_s, y_s^2\}$ , with an associated pair of canonical parameters  $(\gamma_{s;j}, \gamma'_{s;j}) := (\frac{\mu_j}{\sigma_j^2}, -\frac{1}{2\sigma_j^2})$ . Overall, we obtain the exponential family with the canonical parameters

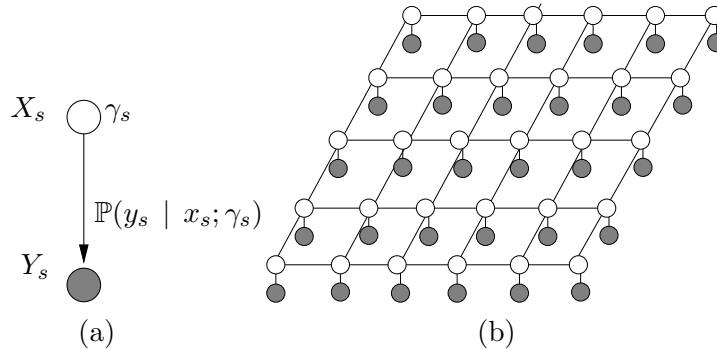
$$\theta_s := \left( \underbrace{\alpha_{s;0}, \dots, \alpha_{s;r-1}}_{\alpha_s}, \underbrace{\gamma_{s;0}, \dots, \gamma_{s;r-1}}_{\gamma_s}, \underbrace{\gamma'_{s;0}, \dots, \gamma'_{s;r-1}}_{\gamma'_s} \right),$$

and a density of the form

$$\begin{aligned} p_{\theta_s}(y_s, x_s) &= p_{\alpha}(x_s) p_{\gamma_s}(y_s | x_s) \\ &\propto \exp \left\{ \sum_{j=0}^{r-1} \alpha_{s;j} \mathbb{I}_j(x_s) + \sum_{j=0}^{r-1} \mathbb{I}_j(x_s) [\gamma_{s;j} y_s + \gamma'_{s;j} y_s^2] \right\}. \end{aligned} \quad (3.15)$$

The pair  $(X_s, Y_s)$  serves a basic block for building more sophisticated graphical models, suitable for describing multivariate data  $((X_1, Y_1), \dots, (X_m, Y_m))$ . For instance, suppose that the vector  $X = (X_1, \dots, X_m)$  of multinomial variates is modeled as a Markov random field (see Example 4 on the Potts model) with respect to some underlying graph  $G = (V, E)$ , and the variables  $Y_s$  are conditionally independent given  $\{X = x\}$ . These assumptions lead to an exponential family  $p_{\theta}(y, x)$  with density

$$p_{\alpha}(x) p_{\gamma}(y | x) \propto \exp \left\{ \sum_{s \in V} \alpha_s(x_s) + \sum_{(s,t) \in E} \alpha_{st}(x_s, x_t) \right\} \prod_{s \in V} p_{\gamma_s}(y_s | x_s), \quad (3.16)$$



**Fig. 3.2.** (a) Graphical representation of a finite mixture of Gaussians:  $X_s$  is a multinomial label for the mixture component, whereas  $Y_s$  is conditionally Gaussian given  $X_s = j$ , with Gaussian density  $p_{\gamma_s}(y_s | x_s)$  in exponential form. (b) Coupled mixture-of-Gaussian graphical model, in which the vector  $X = (X_1, \dots, X_m)$  are a Markov random field on an undirected graph, and the elements of  $(Y_1, \dots, Y_m)$  are conditionally independent given  $X$ .

corresponding to a product of the Potts-type distribution  $p_\alpha(x)$  over  $X$ , and the local conditional distributions  $p_{\gamma_s}(y_s | x_s)$ . Here we have used  $\alpha_s(x_s)$  as a shorthand for the exponential family representation  $\sum_{j=0}^{r-1} \alpha_{s;j} \mathbb{I}_j(x_s)$ , and similarly for the quantity  $\alpha_{st}(x_s, x_t)$ ; see Example 4 where this notation was used. ♣

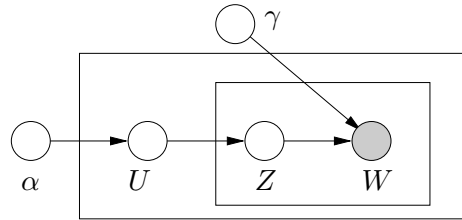
Whereas the mixture model just described is a two-level hierarchy, the following example involves three distinct levels:

**Example 7 (Latent Dirichlet allocation).** The *latent Dirichlet allocation* model [29] is a particular type of hierarchical Bayes model for capturing the statistical dependencies among words in a corpus of documents. It involves three different types of random variables: “documents”  $U$ , “topics”  $Z$  and “words”  $W$ . Words are multinomial random variables ranging over some vocabulary. Topics are also multinomial random variables. Associated to each value of the topic variable there is a distribution over words. A document is a distribution over topics. Finally, a corpus is defined by placing a distribution on the documents. For the latent Dirichlet allocation model, this latter distribution is a Dirichlet distribution.

More formally, words  $W$  are drawn from a multinomial distribution,  $\mathbb{P}(W = j | Z = i; \gamma) = \exp(\gamma_{ij})$ , for  $j = 0, 1, \dots, k - 1$ , where  $\gamma_{ij}$  is a parameter encoding the probability of the  $j$ th word under the  $i$ th topic. This conditional distribution can be expressed as an exponential family in terms of indicator functions as follows:

$$p_{\gamma}(w|z) \propto \exp \left( \sum_{i=0}^{r-1} \sum_{j=0}^{k-1} \gamma_{ij} \mathbb{I}_i(z) \mathbb{I}_j(w) \right), \quad (3.17)$$

where  $\mathbb{I}_i(z)$  is an  $\{0, 1\}$ -valued indicator for the event  $\{Z = i\}$ , and similarly for  $\mathbb{I}_j(w)$ . At the next level of the hierarchy (see Figure 3.3),



**Fig. 3.3.** Graphical illustration of the Latent Dirichlet allocation (LDA) model. The word variable  $W$  is multinomial conditioned on the underlying topic  $Z$ , where  $\gamma$  specifies the topic distributions. The topics  $Z$  are also modeled as multinomial variables, with distributions parameterized by a probability vector  $U$  that follows a Dirichlet distribution with parameter  $\alpha$ . This model is an example of a hierarchical Bayesian model. The rectangles, known as plates, denote replication of the random variables.

the topic variable  $Z$  also follows a multinomial distribution whose parameters are determined by the Dirichlet variable as follows:

$$p(z | u) \propto \exp \left\{ \sum_{i=0}^{r-1} \mathbb{I}_i[z] \log u_i \right\}. \quad (3.18)$$

Finally, at the top level of the hierarchy, the Dirichlet variable  $U$  has a density with respect to Lebesgue measure of the form  $p_{\alpha}(u) \propto \exp\{\sum_{i=0}^{r-1} \alpha_i \log u_i\}$ . Overall then, for a single triplet  $X := (U, Z, W)$ , the LDA model is an exponential family with param-

eter vector  $\theta := (\alpha, \gamma)$ , and an associated density of the form

$$p_\alpha(u)p(z | u) p_\gamma(w | z) \propto \exp \left\{ \sum_{i=0}^{r-1} \alpha_i \log u_i + \sum_{i=0}^{r-1} \mathbb{I}_i[z] \log u_i \right\} \\ \times \exp \left\{ \sum_{i=0}^{r-1} \sum_{j=0}^{k-1} \gamma_{ij} \mathbb{I}_i[z] \mathbb{I}_j[w] \right\}. \quad (3.19)$$

The sufficient statistics  $\phi$  consist of the collections of functions  $\{\log u_i\}$ ,  $\{\mathbb{I}_i[z] \log u_i\}$  and  $\{\mathbb{I}_i[z] \mathbb{I}_j[w]\}$ . As illustrated in Figure 3.3, the full LDA model entails replicating these types of local structures many times. ♣

Many graphical models include what are known as “hard-core” constraints, meaning that a subset of configurations are forbidden. Instances of such problems include decoding of binary linear codes, and other combinatorial optimization problems (e.g., graph matching, covering, packing, etc.) At first glance, it might seem that such families of distributions cannot be represented as exponential families, since the density  $p_\theta$  in an exponential family is strictly positive—that is,  $p_\theta(x) > 0$  for all  $x \in \mathcal{X}^m$ . In the following example, we show that these hard-core constraints can be incorporated within the exponential family framework by making appropriate use of the underlying base measure  $\nu$ .

**Example 8 (Models with hard constraints).** One important domain in which hard-core constraints arise is communication theory, and in particular the problem of error-control coding [84, 164, 196]. To motivate the coding problem, suppose that two people—following an old convention, let us call them Alice and Bob—wish to communicate. We assume that they can communicate by transmitting a sequence of 0s and 1s, but make the problem interesting by assuming that the communication channel linking them behaves in a random manner. Concretely, in a “bit-flipping” channel, any binary digit transmitted by Alice is received correctly by Bob only with probability  $1 - \epsilon$ . As a probabilistic model, this channel can be modeled by the conditional distribution

$$p(y | x) := \begin{cases} 1 - \epsilon & \text{if } x = y \\ \epsilon & \text{if } x \neq y, \end{cases}$$

where  $x \in \{0, 1\}$  represents the bit transmitted by Alice, and  $y \in \{0, 1\}$  represents the bit received by Bob.

In order to mitigate the “noisiness” of the channel, Alice and Bob agree on the following *block coding* scheme: instead of transmitting a single bit, they communicate using strings  $(x_1, x_2, \dots, x_m)$  of bits, and moreover, they agree to use only a subset of the total number  $2^m$  of length  $m$  binary strings. These privileged strings, which are known as *codewords*, can be defined using a particular type of graphical model. Concretely, suppose that Alice decides to choose from among strings  $x \in \{0, 1\}^m$  that satisfy a set of parity checks—say of the form  $x_1 \oplus x_2 \oplus x_3 = 0$ , where  $\oplus$  denotes addition in modulo two arithmetic. Let us consider a collection  $F$  of such parity checks; each  $a \in F$  enforces a parity check constraint on some subset  $N(a) \subset \{1, \dots, m\}$  of bits. If we define the indicator function

$$\psi_a(x_{N(a)}) := \begin{cases} 1 & \text{if } \bigoplus_{i \in N(a)} x_i = 0 \\ 0 & \text{otherwise,} \end{cases} \quad (3.20)$$

then a *binary linear code* consists of all bit strings  $(x_1, x_2, \dots, x_m) \in \{0, 1\}^m$  for which  $\prod_{a \in F} \psi_a(x_{N(a)}) = 1$ . These constraints are known as “hard-core,” since they take the value 0 for some settings of  $x$ , and hence eliminate such configurations from the support of the distribution. Letting  $(y_1, \dots, y_m) \in \{0, 1\}^m$  denote the string of bits received by Bob, his goal is to use these observations to infer which codeword was transmitted by Alice. Depending on the error metric used, this decoding problem corresponds to either computing marginals or modes of the posterior distribution

$$p(x_1, \dots, x_m \mid y_1, \dots, y_m) \propto \prod_{i=1}^m p(y_i \mid x_i) \prod_{a \in F} \psi_a(x_{N(a)}) \quad (3.21)$$

This distribution can be described by a factor graph, with the bits represented as (circular) variable nodes, and the parity check indicator functions represented at the (square) factor nodes (see Figure 2.9).

The code can also be viewed as an exponential family on this graph, if we take densities with respect to an *appropriately chosen base measure*. In particular, define the counting measure restricted to valid code-



words as follows

$$\nu(dx) := \left[ \prod_{a \in F} \psi_a(x_{N(a)}) \right] dx_1 \dots dx_m. \quad (3.22)$$

A little calculation shows that conditional distribution  $p(y_i | x_i)$  can be written in an exponential form as

$$p(y_i | x_i) \propto \exp(\theta_i x_i), \quad (3.23)$$

where the exponential parameter  $\theta_i$  is defined by the observation  $y_i$  and the channel noise parameter  $\epsilon$  via the relation

$$\theta_i = y_i \log \frac{1-\epsilon}{\epsilon} + (1-y_i) \log \frac{\epsilon}{1-\epsilon} = (1-2y_i) \log \frac{\epsilon}{1-\epsilon}. \quad (3.24)$$

Note that we are using the fact that the vector  $(y_1, y_2, \dots, y_m)$  is observed, and hence can be viewed as fixed. With this set-up, the distribution (3.21) can be cast as an exponential family, where the density is taken with respect to the restricted counting measure (3.22), and has the form  $p_\theta(x) = \exp(\sum_{i=1}^m \theta_i x_i)$ .



### 3.4 Mean parameterization and inference problems

Thus far, we have characterized any exponential family member  $p_\theta(x)$  by its vector of canonical parameters  $\theta \in \Omega$ . As we discuss in this section, it turns out that any exponential family has an alternative parameterization in terms of a vector of *mean parameters*. Moreover, various statistical computations, among them marginalization and maximum likelihood estimation, can be understood as transforming from one parameterization to the other.

We digress momentarily to make an important remark regarding the role of observed variables or conditioning for the marginalization problem. As discussed earlier, applications of graphical models frequently involve conditioning on a subset of random variables—say  $Y$ —that represent observed quantities, and computing marginals under the posterior distribution  $p_\theta(x | y)$ . The application to error-correcting coding, just discussed in Example 8, is one such example. So as to simplify notation in our development to follow, it is convenient to no longer make direct reference to observed variables, but rather discuss marginalization

and other computational problems only in reference to unconditional forms of exponential families. For such computational purposes, there is no loss of generality in doing so, since the effects of observations can always be absorbed by modifying the canonical parameters  $\theta$  and/or the sufficient statistics  $\phi$  appropriately. (For instance, in Example 8, the noisy observation  $y_i$  at node  $i$  contributes a new factor  $\exp(\theta_i x_i)$  to the exponential family factorization, as described in equations (3.23) and (3.24).)

### 3.4.1 Mean parameter spaces and marginal polytopes

Let  $p$  be a given density defined with respect to the underlying base measure  $\nu$ ; for the moment, we do not assume that  $p$  is a member of an exponential family defined with respect to  $\nu$ . The mean parameter  $\mu_\alpha$  associated with a sufficient statistic  $\phi_\alpha : \mathcal{X}^m \rightarrow \mathbb{R}$  is defined by the expectation

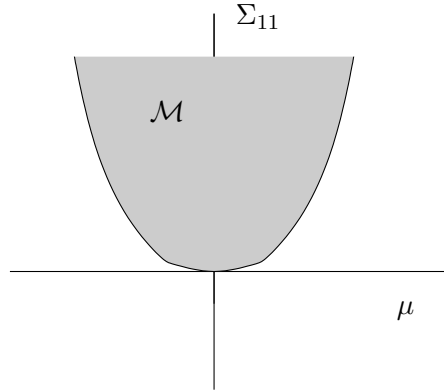
$$\mu_\alpha = \mathbb{E}_p[\phi_\alpha(X)] = \int \phi_\alpha(x) p(x) \nu(dx), \quad \text{for } \alpha \in \mathcal{I}. \quad (3.25)$$

In this way, we define a vector of *mean parameters*  $(\mu_1, \dots, \mu_d)$ , one for each of the  $|\mathcal{I}| = d$  sufficient statistics  $\phi_\alpha$ , with respect to an arbitrary density  $p$ . It turns out that under suitable technical conditions, this vector provides an alternative parameterization of the exponential family defined by  $\{\phi_\alpha\}$  and  $\nu$ . An interesting object is the set of all such vectors  $\mu \in \mathbb{R}^d$  traced out as the underlying density  $p$  is varied. More formally, we define

$$\mathcal{M} := \left\{ \mu \in \mathbb{R}^d \mid \exists p \text{ s. t. } \mathbb{E}_p[\phi_\alpha(X)] = \mu_\alpha \forall \alpha \in \mathcal{I} \right\}. \quad (3.26)$$

We illustrate this definition with a continuation of a previous example:

**Example 9 (Gaussian MRF mean parameters).** Using the canonical parameterization of the Gaussian Markov random field provided in Example 5, the mean parameters for a Gaussian Markov random field are the second-order moment matrix  $\Sigma := \mathbb{E}[XX^T] \in \mathbb{R}^{m \times m}$ , and the mean vector  $\mu = \mathbb{E}[X] \in \mathbb{R}^m$ . For this particular model, it is straightforward to characterize the set  $\mathcal{M}$  of globally realizable mean parameters



**Fig. 3.4.** Illustration of the set  $\mathcal{M}$  for a scalar Gaussian: the model has two mean parameters  $\mu = \mathbb{E}[X]$  and  $\Sigma_{11} = \mathbb{E}[X^2]$ , which must satisfy the quadratic constraint  $\Sigma_{11} - \mu^2 \geq 0$ . Note that the set  $\mathcal{M}$  is convex, which is a general property.

$(\mu, \Sigma)$ . We begin by recognizing that if  $(\mu, \Sigma)$  are realized by some distribution (not necessarily Gaussian), then  $\Sigma - \mu\mu^T$  must be a valid covariance matrix of the random vector  $X$ , implying that the positive semidefiniteness (PSD) condition  $\Sigma - \mu\mu^T \succeq 0$  must hold. Conversely, any pair  $(\mu, \Sigma)$  for which the PSD constraint holds, we may construct a multivariate Gaussian distribution with mean  $\mu$ , and (possibly degenerate) covariance  $\Sigma - \mu\mu^T$ , which by construction realizes  $(\mu, \Sigma)$ . Thus, we have established that for a Gaussian Markov random field, we have

$$\mathcal{M} = \{(\mu, \Sigma) \in \mathbb{R}^m \times \mathcal{S}_+^m \mid \Sigma - \mu\mu^T \succeq 0\}, \quad (3.27)$$

where  $\mathcal{S}_+^m$  denotes the set of  $m \times m$  symmetric positive semidefinite matrices. Figure 3.4 illustrates this set in the scalar case ( $m = 1$ ), where the mean parameters  $\mu = \mathbb{E}[X]$  and  $\Sigma_{11} = \mathbb{E}[X^2]$  must satisfy the semidefinite constraint  $\Sigma_{11} - \mu^2 \geq 0$ . ♣

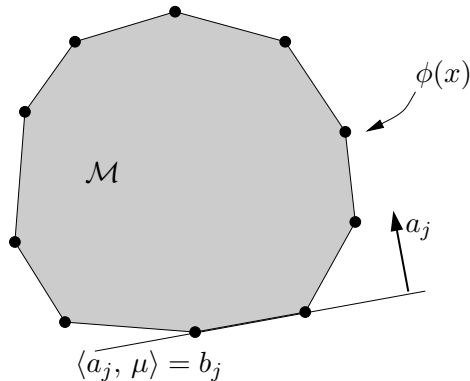
For any graphical model, the set  $\mathcal{M}$  is always a convex subset of  $\mathbb{R}^d$ . Indeed, if  $\mu$  and  $\mu'$  are both elements of  $\mathcal{M}$ , then there must exist distributions  $p$  and  $p'$  that realize them (i.e., such that  $\mu = \mathbb{E}_p[\phi(X)]$  and similarly for  $\mu'$ ). For any  $\lambda \in [0, 1]$ , the convex combination  $\mu(\lambda) = \lambda\mu + (1 - \lambda)\mu'$  is realized by the mixture distribution

$\lambda p + (1 - \lambda)p'$ , so that  $\mu(\lambda)$  also belongs to  $\mathcal{M}$ . In Appendix B.3, we summarize further properties of  $\mathcal{M}$  that hold for general exponential families.

The case of discrete random variables yields a set  $\mathcal{M}$  with some special properties. More specifically, for any random vector  $(X_1, X_2, \dots, X_m)$  such that the associated state space  $\mathcal{X}^m$  is finite, we can write the set  $\mathcal{M}$  as

$$\begin{aligned} \mathcal{M} &= \left\{ \mu \in \mathbb{R}^d \mid \exists p(x) \geq 0, \sum_{x \in \mathcal{X}^m} p(x) = 1 \text{ s.t. } \mu = \sum_{x \in \mathcal{X}^m} \phi(x)p(x) \right\} \\ &= \text{conv}\{\phi(x), x \in \mathcal{X}^m\}, \end{aligned} \quad (3.28)$$

where  $\text{conv}$  denotes the convex hull operation (see Appendix A.2). Consequently, when  $|\mathcal{X}^m|$  is finite, the set  $\mathcal{M}$  is—by definition—a *convex polytope*.



**Fig. 3.5.** Generic illustration of  $\mathcal{M}$  for a discrete random variable with  $|\mathcal{X}^m|$  finite. In this case, the set  $\mathcal{M}$  is a convex polytope, corresponding to the convex hull of  $\{\phi(x) \mid x \in \mathcal{X}^m\}$ . By the Minkowski-Weyl theorem, this polytope can also be written as the intersection of a finite number of halfspaces, each of the form  $\{\mu \in \mathbb{R}^d \mid \langle a_j, \mu \rangle \geq b_j\}$  for some pair  $(a_j, b_j) \in \mathbb{R}^d \times \mathbb{R}$ .

The Minkowski-Weyl theorem [198], stated in Appendix A.2, provides an alternative description of a convex polytope. As opposed to the convex hull of a finite collection of vectors, any polytope  $\mathcal{M}$  can also be characterized by a finite collection of linear inequality constraints. Explicitly, for any polytope  $\mathcal{M}$ , there exists a collection

$\{(a_j, b_j) \in \mathbb{R}^d \times \mathbb{R} \mid j \in \mathcal{J}\}$  with  $|\mathcal{J}|$  finite such that

$$\mathcal{M} = \{\mu \in \mathbb{R}^d \mid \langle a_j, \mu \rangle \geq b_j \ \forall j \in \mathcal{J}\}. \quad (3.29)$$

In geometric terms, this representation shows that  $\mathcal{M}$  is equal to the intersection of a finite collection of half-spaces, as illustrated in Figure 3.5. Let us show the distinction between the convex hull (3.28) and linear inequality (3.29) representations using the Ising model.

**Example 10 (Ising mean parameters).** Continuing from Example 3, the sufficient statistics for the Ising model are the singleton functions  $\{x_s, s \in V\}$  and the pairwise functions  $\{x_s x_t, (s, t) \in E\}$ . The vector of sufficient statistics takes the form

$$\phi(x) := (x_1, x_2, \dots, x_m, x_s x_t, (s, t) \in E) \in \mathbb{R}^{m+|E|}. \quad (3.30)$$

The associated mean parameters correspond to particular marginal probabilities, associated with nodes and edges of the graph  $G$  as

$$\mu_s = \mathbb{E}_p[X_s] = \mathbb{P}[X_s = 1] \quad \text{for all } s \in V, \text{ and} \quad (3.31a)$$

$$\mu_{st} = \mathbb{E}_p[X_s X_t] = \mathbb{P}[(X_s, X_t) = (1, 1)] \quad \text{for all } (s, t) \in E. \quad (3.31b)$$

Consequently, the mean parameter vector  $\mu \in \mathbb{R}^{|V|+|E|}$  consists of marginal probabilities over singletons ( $\mu_s$ ), and pairwise marginals over variable pairs on graph edges ( $\mu_{st}$ ). The set  $\mathcal{M}$  consists of the convex hull of  $\{\phi(x), x \in \{0, 1\}^m\}$ , where  $\phi$  is given in equation (3.30). In probabilistic terms, the set  $\mathcal{M}$  corresponds to the set of all singleton and pairwise marginal probabilities that can be realized by some distribution over  $(X_1, \dots, X_m) \in \{0, 1\}^m$ . In the polyhedral combinatorics literature, this set is known as the *correlation polytope*, or the *cut polytope* [66, 183].

To make these ideas more concrete, consider the simplest non-trivial case: namely, a pair of variables  $(X_1, X_2)$ , and the graph consisting of the single edge joining them. In this case,  $\mathcal{M}$  is a polytope in three dimensions (two nodes plus one edge): it is the convex hull of the vectors  $\{(x_1, x_2, x_1 x_2) \mid (x_1, x_2) \in \{0, 1\}^2\}$ , or more explicitly

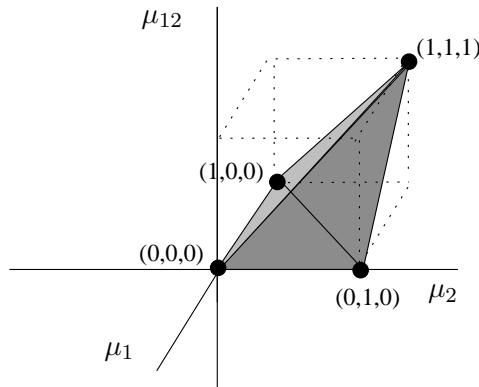
$$\text{conv}\{(0, 0, 0), (1, 0, 0), (0, 1, 0), (1, 1, 1)\},$$

as illustrated in Figure 3.6.

Let us also consider the half-space representation (3.29) for this case. Elementary probability theory and a little calculation shows that the three mean parameters  $(\mu_1, \mu_2, \mu_{12})$  must satisfy the constraints  $0 \leq \mu_{12} \leq \mu_i$  for  $i = 1, 2$  and  $1 + \mu_{12} - \mu_1 - \mu_2 \geq 0$ . We can write these constraints in matrix-vector form as

$$\begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & -1 \\ 0 & 1 & -1 \\ -1 & -1 & 1 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_{12} \end{bmatrix} \geq \begin{bmatrix} 0 \\ 0 \\ 0 \\ -1 \end{bmatrix}.$$

These four constraints provide an alternative characterization of the three-dimensional polytope illustrated in Figure 3.6.



**Fig. 3.6.** Illustration of  $\mathcal{M}$  for the special case of an Ising model with two variables  $(X_1, X_2) \in \{0, 1\}^2$ . The four mean parameters  $\mu_1 = \mathbb{E}[X_1]$ ,  $\mu_2 = \mathbb{E}[X_2]$  and  $\mu_{12} = \mathbb{E}[X_1 X_2]$  must satisfy the constraints  $0 \leq \mu_{12} \leq \mu_i$  for  $i = 1, 2$ , and  $1 + \mu_{12} - \mu_1 - \mu_2 \geq 0$ . These constraints carve out a polytope with four facets, contained within the unit hypercube.



Our next example deals with a interesting family of polytopes that arises in error-control coding and binary matroid theory:

**Example 11 (Codeword polytopes and binary matroids).** Recall the definition of a binary linear code from Example 8: it corresponds to the subset of binary strings  $x \in \{0, 1\}^m$  that satisfy a set of parity

check relations. Specifically, the base measure in the exponential family representation enforces the constraint  $\prod_{a \in F} \psi_a(x_{N(a)}) = 1$ , where each  $a$  indexes the subset  $N(a) \subseteq \{1, 2, \dots, m\}$ , and the associated parity check  $\psi_a$  imposes the constraint  $\psi_a(x_{N(a)}) = 1$  if  $\oplus_{i \in N(a)} x_i = 0$ , and  $\psi_a(x_{N(a)}) = 0$  otherwise. We define the code

$$\mathbb{C} := \{x \in \{0, 1\}^m \mid \prod_{a \in F} \psi_a(x_{N(a)}) = 1\}.$$

Viewed as an exponential family, the parity checks define the base measure  $\nu$ , and the sufficient statistics are simply  $\phi(x) = (x_1, x_2, \dots, x_m)$ . Consequently, the set  $\mathcal{M}$  for this problem corresponds to the *codeword polytope*—namely, the convex hull of all possible codewords

$$\begin{aligned} \mathcal{M} &= \text{conv} \{x \in \{0, 1\}^m \mid \prod_{a \in F} \psi_a(x_{N(a)}) = 1\} \\ &= \text{conv}\{x \in \mathbb{C}\}. \end{aligned} \tag{3.32}$$

To provide a concrete illustration, consider the code on  $m = 3$  bits, defined by the single parity check relation  $x_1 \oplus x_2 \oplus x_3 = 0$ . Figure 3.7(a) shows the factor graph representation of this toy code. This parity check eliminates half of the  $2^3 = 8$  possible binary sequences of length 3. The codeword polytope is simply the convex hull of  $\{(0, 0, 0), (0, 1, 1), (1, 0, 1), (1, 1, 0)\}$ , as illustrated in Figure 3.7(b). Equivalently, we can represent this codeword polytope in terms of half-space constraints. For this single parity check code, the codeword polytope is defined by the four inequalities

$$(1 - \mu_1) + (1 - \mu_2) + (1 - \mu_3) \geq 1, \tag{3.33a}$$

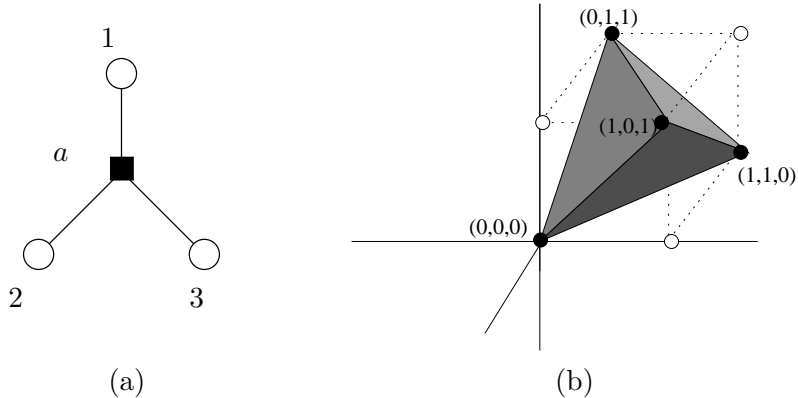
$$(1 - \mu_1) + \mu_2 + \mu_3 \geq 1, \tag{3.33b}$$

$$\mu_1 + (1 - \mu_2) + \mu_3 \geq 1, \quad \text{and} \tag{3.33c}$$

$$\mu_1 + \mu_2 + (1 - \mu_3) \geq 1. \tag{3.33d}$$

As we discuss at more length in Section 8.4.4, these inequalities can be understood as requiring that  $(\mu_1, \mu_2, \mu_3)$  is at least distance 1 from each of the forbidden odd-parity vertices of the hypercube  $\{0, 1\}^3$ .

Of course, for larger code lengths  $m$  and many parity checks, the associated codeword polytope has a much more complicated structure.



**Fig. 3.7.** Illustration of  $\mathcal{M}$  for the special case of the binary linear code  $\mathbb{C} = \{(0, 0, 0), (0, 1, 1), (1, 1, 0), (1, 1, 0)\}$ . This code is defined by a single parity check

It plays a key role in the error-control decoding problem [75], studied intensively in the coding and information theory communities. Any binary linear code can be identified with a binary matroid [182], in which context the codeword polytope is referred to as the *cycle polytope* of the binary matroid. There is a rich literature in combinatorics on the structure of these codeword or cycle polytopes [8, 210, 100, 66]. ♣

Examples 10 and 11 are specific instances of a more general object that we refer to as the *marginal polytope* for a discrete graphical model. The marginal polytope is defined for any graphical model with multinomial random variables  $X_s \in \mathcal{X}_s := \{0, 1, \dots, r_s - 1\}$  at each vertex  $s \in V$ ; note that the cardinality  $|\mathcal{X}_s| = r_s$  can differ from node to node. Consider the exponential family defined in terms of  $\{0, 1\}$ -valued indicator functions

$$\forall s \in V, j \in \mathcal{X}_s, \quad \mathbb{I}_{s;j}(x_s) := \begin{cases} 1 & \text{if } x_s = j \\ 0 & \text{otherwise.} \end{cases} \quad (3.34)$$

$$\forall (s, t) \in E, (j, k) \quad \mathbb{I}_{st;jk}(x_s, x_t) := \begin{cases} 1 & \text{if } (x_s, x_t) = (j, k) \\ 0 & \text{otherwise.} \end{cases}$$

We refer to the sufficient statistics (3.34) as the *standard overcomplete*



*representation.* Its overcompleteness was discussed previously in Example 4.

With this choice of sufficient statistics, the mean parameters take a very intuitive form: in particular, for each node  $s \in V$

$$\mu_{s;j} = \mathbb{E}_p[\mathbb{I}_j(X_s)] = \mathbb{P}[X_s = j] \quad \forall j \in \mathcal{X}_s, \quad (3.35)$$

and for each edge  $(s, t) \in E$ , we have

$$\mu_{st;jk} = \mathbb{E}_p[\mathbb{I}_{st;jk}(X_s, X_t)] = \mathbb{P}[X_s = j, X_t = k] \quad \forall (j, k) \in \mathcal{X}_s \times \mathcal{X}_t. \quad (3.36)$$

Thus, the mean parameters correspond to singleton marginal distributions  $\mu_s$  and pairwise marginal distributions  $\mu_{st}$  associated with the nodes and edges of the graph. In this case, we refer to the set  $\mathcal{M}$  as the *marginal polytope* associated with the graph, and denote it by  $\mathbb{M}(G)$ . Explicitly, it is given by

$$\mathbb{M}(G) := \{\mu \in \mathbb{R}^d \mid \exists p \text{ such that (3.35) holds } \forall (s; j), \text{ and} \\ \text{(3.36) holds } \forall (st; jk)\}. \quad (3.37)$$

Note that the correlation polytope for the Ising model presented in Example 10 is a special case of a marginal polytope, obtained for  $X_s \in \{0, 1\}$  for all nodes  $s$ . The only difference is we have defined marginal polytopes with respect to the standard overcomplete basis of indicator functions, whereas the Ising model is usually parameterized as a minimal exponential family. The codeword polytope of Example 11 is another special case of a marginal polytope. In this case, the reduction requires two steps: first, we convert the factor graph representation of the code—for instance, as shown in Figure 3.7(a)—to an equivalent pairwise Markov random field, involving binary variables at each bit node, and higher-order discrete variables at each factor node. (See Appendix E.2 for details of this procedure for converting from factor graphs to pairwise MRFs.) The marginal polytope associated with this pairwise MRF is simply a lifted version of the codeword polytope. We discuss these and other examples of marginal polytopes in more detail in later sections.

For the toy models considered explicitly in Examples 10 and 11, the number of half-space constraints  $|\mathcal{J}|$  required to characterize the

marginal polytopes was very small ( $|\mathcal{J}| = 4$  in both cases). It is natural to ask how the number of constraints required grows as a function of the graph size. Interestingly, we will see later that this so-called facet complexity depends critically on the graph structure. For trees, any marginal polytope is characterized by local constraints—involving only pairs of random variables on edges—with the total number growing only linearly in the graph size. In sharp contrast, for general graphs with cycles, the constraints are very non-local and the growth in their number is astonishingly fast. For the special case of the Ising model, the book by Deza and Laurent [66] contains a wealth of information about the correlation/cut polytope. The intractability of representing marginal polytopes in a compact manner is one underlying cause of the complexity in performing statistical computation.

### 3.4.2 Role of mean parameters in inference problems

The preceding examples suggest that mean parameters have a central role to play in the marginalization problem. For the multivariate Gaussian (Example 9), an efficient algorithm for computing the mean parameterization provides us with both the Gaussian mean vector, as well as the associated covariance matrix. For the Ising model (see Example 10), the mean parameters completely specify the singleton and pairwise marginals of the probability distribution; the same statement holds more general for the multinomial graphical model defined by the standard overcomplete parameterization (3.34). Even more broadly, the computation of the *forward mapping*, from the canonical parameters  $\theta \in \Omega$  to the mean parameters  $\mu \in \mathcal{M}$ , can be viewed as a fundamental class of inference problems in exponential family models. Although computing the mapping is straightforward for most low-dimensional exponential families, the computation of this forward mapping is extremely difficult for many high-dimensional exponential families.

The *backward mapping*, namely from mean parameters  $\mu \in \mathcal{M}$  to canonical parameters  $\theta \in \Omega$ , also has a natural statistical interpretation. In particular, suppose that we are given a set of samples  $X_1^n := \{X^1, \dots, X^n\}$ , drawn independently from an exponential family member  $p_\theta(x)$ , where the parameter  $\theta$  is unknown. If the goal is

to estimate  $\theta$ , the classical principle of maximum likelihood dictates obtaining an estimate  $\hat{\theta}$  by maximizing the likelihood of the data, or equivalently (after taking logarithms and rescaling), maximizing the quantity

$$\ell(\theta; X_1^n) := \frac{1}{n} \sum_{i=1}^n \log p_\theta(X^i) = \langle \theta, \hat{\mu} \rangle - A(\theta), \quad (3.38)$$

where  $\hat{\mu} := \widehat{\mathbb{E}}[\phi(X)] = \frac{1}{n} \sum_{i=1}^n \phi(X^i)$  is the vector of *empirical mean parameters* defined by the data  $X_1^n$ . The maximum likelihood estimate  $\hat{\theta}$  is chosen to achieve the maximum of this objective function. Note that computing  $\hat{\theta}$  is, in general, another challenging problem, since the objective function involves the log partition function  $A$ . As will be demonstrated by the development to follow, under suitable conditions, the maximum likelihood estimate is unique, and specified by the stationarity condition  $\mathbb{E}_{\hat{\theta}}[\phi(X)] = \hat{\mu}$ . Finding the unique solution to this equation is equivalent to computing the backward mapping  $\mu \mapsto \theta$ , from mean parameters to canonical parameters. In general, computing this inverse mapping is also computationally intensive.

### 3.5 Properties of $A$

With this background in mind, we now turn to a deeper exploration of the properties of the cumulant function  $A$ . Perhaps the most important property of  $A$  is its convexity, which brings us into convex analysis, and more specifically leads us to the study of the the conjugate dual function of  $A$ . Under suitable conditions, the function  $A$  and its conjugate dual  $A^*$ —or more precisely, their derivatives—turn out to define a one-to-one and surjective mapping between the canonical and mean parameters. As discussed above, the mapping between canonical and mean parameters is the core challenge in various statistical computations in high-dimensional graphical models.

### 3.5.1 Derivatives and convexity

Recall that a real-valued function  $g$  is *convex* if, for any two  $x, y$  belonging to the domain of  $g$  and any scalar  $\lambda \in (0, 1)$ , the inequality

$$g(\lambda x + (1 - \lambda)y) \leq \lambda g(x) + (1 - \lambda)g(y) \quad (3.39)$$

holds. Geometrically, the line connecting the function values  $g(x)$  and  $g(y)$  lies above the graph of the function itself. The function is *strictly convex* if the inequality (3.39) is strict for all  $x \neq y$ . (See Appendix A.2.5 for some additional properties of convex functions.) We begin by establishing that the log partition function is smooth, and convex in terms of  $\theta$ .

**Proposition 2.** *The cumulant function*

$$A(\theta) := \log \int_{\mathcal{X}^m} \exp\langle \theta, \phi(x) \rangle \nu(dx) \quad (3.40)$$

*associated with any regular exponential family has the following properties:*

(a) *It has derivatives of all orders on its domain  $\Omega$ . The first two derivatives yield the cumulants of the random vector  $\phi(X)$  as follows:*

$$\frac{\partial A}{\partial \theta_\alpha}(\theta) = \mathbb{E}_\theta[\phi_\alpha(X)] := \int \phi_\alpha(x) p_\theta(x) \nu(dx). \quad (3.41a)$$

$$\frac{\partial^2 A}{\partial \theta_\alpha \partial \theta_\beta}(\theta) = \mathbb{E}_\theta[\phi_\alpha(X)\phi_\beta(X)] - \mathbb{E}_\theta[\phi_\alpha(X)]\mathbb{E}_\theta[\phi_\beta(X)] \quad (3.41b)$$

(b) *Moreover,  $A$  is a convex function of  $\theta$  on its domain  $\Omega$ , and strictly so if the representation is minimal.*

*Proof.* Let us assume that differentiating through the integral (3.40) defining  $A$  is permitted; verifying the validity of this assumption is a standard argument using the dominated convergence theorem (e.g.,

Brown [40]). Under this assumption, we have

$$\begin{aligned}
\frac{\partial A}{\partial \theta_\alpha}(\theta) &= \frac{\partial}{\partial \theta_\alpha} \left\{ \log \int_{\mathcal{X}^m} \exp\langle \theta, \phi(x) \rangle \nu(dx) \right\} \\
&= \left[ \int_{\mathcal{X}^m} \frac{\partial}{\partial \theta_\alpha} \exp\langle \theta, \phi(x) \rangle \nu(dx) \right] / \left[ \int_{\mathcal{X}^m} \exp\langle \theta, \phi(u) \rangle \nu(du) \right] \\
&= \int_{\mathcal{X}^m} \phi_\alpha(x) \frac{\exp\langle \theta, \phi(x) \rangle \nu(dx)}{\int_{\mathcal{X}^m} \exp\langle \theta, \phi(u) \rangle \nu(du)} \\
&= \mathbb{E}_\theta[\phi_\alpha(X)],
\end{aligned}$$

which establishes equation (3.41a). The formula for the higher-order derivatives can be proven in an entirely analogous manner.

Observe from equation (3.41b) that the second order partial derivative  $\frac{\partial^2 A}{\partial \theta_\alpha \partial \theta_\beta^2}$  is equal to the covariance element  $\text{cov}\{\phi_\alpha(X), \phi_\beta(X)\}$ . Therefore, the full Hessian  $\nabla^2 A(\theta)$  is the covariance matrix of the random vector  $\phi(X)$ , and so is positive semidefinite on the open set  $\Omega$ , which ensures convexity (see Theorem 4.3.1 of Hiriart-Urruty and Lemaréchal [109]). If the representation is minimal, there is no non-zero vector  $a \in \mathbb{R}^d$  and constant  $b \in \mathbb{R}$  such that  $\langle a, \phi(x) \rangle = b$  holds  $\nu$ -a.e.. This condition implies  $\text{var}_\theta[\langle a, \phi(x) \rangle] = a^T \nabla^2 A(\theta) a > 0$  for all  $a \in \mathbb{R}^d$  and  $\theta \in \Omega$ ; this strict positive definiteness of the Hessian on the open set  $\Omega$  implies strict convexity [109].

□

### 3.5.2 Forward mapping to mean parameters

We now turn to an in-depth consideration of the forward mapping  $\theta \mapsto \mu$ , from the canonical parameters  $\theta \in \Omega$  defining a distribution  $p_\theta$  to its associated vector of mean parameters  $\mu \in \mathbb{R}^d$ . Note that the gradient  $\nabla A$  can be viewed as mapping from  $\Omega$  to  $\mathbb{R}^d$ . Indeed, Proposition 2 demonstrates that the range of this mapping is contained within the set  $\mathcal{M}$  of realizable mean parameters, defined previously as

$$\mathcal{M} := \left\{ \mu \in \mathbb{R}^d \mid \exists p \text{ s. t. } \mathbb{E}_p[\phi(X)] = \mu \right\}.$$

We will see that a great deal hinges on the answers to the following two questions:

- (a) when does  $\nabla A$  define a one-to-one mapping?
- (b) when does the image of  $\Omega$  under the mapping  $\nabla A$ —that is, the set  $\nabla A(\Omega)$ —fully cover the set  $\mathcal{M}$ ?

The answer to the first question is relatively straightforward, essentially depending on whether or not the exponential family is minimal. The second question is somewhat more delicate: to begin, note that our definition of  $\mathcal{M}$  allows for mean parameters  $\mu \in \mathbb{R}^d$  generated by *any* possible distribution, not just distributions  $p_\theta$  in the exponential family defined by the sufficient statistics  $\phi$ . It turns out that this extra freedom does not really enlarge the set  $\mathcal{M}$ ; as Theorem 1 makes precise, under suitable conditions, all mean parameters in  $\mathcal{M}$  can be realized by an exponential family distribution (or, for boundary points, by a limiting sequence of such distributions).

We begin with a result addressing the first question:

**Proposition 3.** *The gradient mapping  $\nabla A : \Omega \rightarrow \mathcal{M}$  is one-to-one if and only if the exponential representation is minimal.*

*Proof.* If the representation is not minimal, then there must exist a non-zero vector  $\gamma \in \mathbb{R}^d$  for which  $\langle \gamma, \phi(x) \rangle$  is a constant (almost surely with respect to  $\nu$ ). Given any parameter  $\theta^1 \in \Omega$ , let us define another parameter vector  $\theta^2 = \theta^1 + t\gamma$ , where  $t \in \mathbb{R}$ . Since  $\Omega$  is open, choosing  $t$  sufficiently small ensures that  $\theta^2 \in \Omega$  as well. By the condition on the vector  $\gamma$ , the densities  $p_{\theta^1}$  and  $p_{\theta^2}$  induce the same probability distribution (only their normalization constants differ). For this pair, we have  $\nabla A(\theta^1) = \nabla A(\theta^2)$ , so that  $\nabla A$  is not one-to-one.

Conversely, if the representation is minimal, then  $A$  is strictly convex by Proposition 2. For any strictly convex and differentiable function, we have  $A(\theta^2) > A(\theta^1) + \langle \nabla A(\theta^1), \theta^2 - \theta^1 \rangle$ , for all  $\theta^1 \neq \theta^2$  in the domain  $\Omega$ . The same inequality also holds with the roles of  $\theta^1$  and  $\theta^2$  reversed; adding together these inequalities yields that

$$\langle \nabla A(\theta^1) - \nabla A(\theta^2), \theta^1 - \theta^2 \rangle > 0$$

for all distinct  $\theta^1, \theta^2 \in \Omega$ , which shows that  $\nabla A$  is one-to-one.  $\square$

In general, although the gradient mapping  $\nabla A$  is not one-to-one for an overcomplete representation, there is still a one-to-one correspondence between each element of  $\nabla A(\Omega)$  and an affine subset of  $\Omega$ . In particular, this affine subset contains all those canonical parameters  $\theta$  that are mapped to the same mean parameter. For either a minimal or an overcomplete representation, we say that a pair  $(\theta, \mu)$  is *dually coupled* if  $\mu = \nabla A(\theta)$ . This notion of dual coupling plays an important role in the sequel.

We now turn to the second issue regarding the image  $\nabla A(\Omega)$  of the domain of valid canonical parameters  $\Omega$  under the gradient mapping  $\nabla A$ . Specifically, the goal is to determine for which mean parameter vectors  $\mu \in \mathcal{M}$  does there exist a vector  $\theta = \theta(\mu) \in \Omega$  such that  $\mathbb{E}_\theta[\phi(X)] = \mu$ . The solution turns out to be rather simple: the image  $\nabla A(\Omega)$  is simply the interior  $\mathcal{M}^\circ$ . This fact is remarkable: it means that (disregarding boundary points) *all mean parameters  $\mathcal{M}$  that are realizable by some distribution can be realized by a member of the exponential family*. To provide some intuition into this fact, consider the maximum entropy problem (3.3) for a given mean parameter  $\mu$  in the interior of  $\mathcal{M}$ . As discussed earlier, when a solution to this problem exists, it necessarily takes the form of an exponential family member, say  $p_{\theta^*}$ . Moreover, from the optimality conditions for the maximum entropy problem, this exponential family member must satisfy the moment-matching conditions  $\mathbb{E}_{\theta^*}[\phi(X)] = \mu$ . Note that these moment-matching conditions are identical to those defining the maximum likelihood problem (3.38)—as we discuss in the sequel, this fact is not coincidental, but rather a consequence of the primal-dual relationship between maximum entropy and maximum likelihood.

**Theorem 1.** *In a minimal exponential family, the gradient map  $\nabla A$  is onto the interior of  $\mathcal{M}$ , denoted by  $\mathcal{M}^\circ$ . Consequently, for each  $\mu \in \mathcal{M}^\circ$ , there exists some  $\theta \in \Omega$  such that  $\mathbb{E}_\theta[\phi(X)] = \mu$ .*

We provide the proof of this result in Appendix B. In conjunction with Proposition 3, Theorem 1 guarantees that for minimal exponential families, each mean parameter  $\mu \in \mathcal{M}^\circ$  is *uniquely realized* by some density  $p_{\theta(\mu)}$  in the exponential family. However, a typical exponential

family  $\{p_\theta \mid \theta \in \Omega\}$  describes only a strict subset of all possible densities (with respect to the given base measure  $\nu$ ). In this case, there must exist at least some other density  $p$ —albeit not a member of the exponential family—that also realizes  $\mu$ . The distinguishing property of the exponential distribution  $p_{\theta(\mu)}$  is that, among the set of all distributions that realize  $\mu$ , it has the maximum entropy. The connection between  $A$  and the maximum entropy principle is specified precisely in terms of the conjugate dual function  $A^*$ , to which we now turn.

### 3.6 Conjugate duality: Maximum likelihood and maximum entropy

Conjugate duality is a cornerstone of convex analysis [198, 109], and is a natural source for variational representations. In this section, we explore the relationship between the log partition function  $A$  and its conjugate dual function  $A^*$ . This conjugate relationship is defined by a variational principle that is central to the remainder of this paper, in that it underlies a wide variety of known algorithms, both of an exact nature (e.g., the junction tree algorithm and its special cases of Kalman filtering, the forward-backward algorithm, peeling algorithms) and an approximate nature (e.g., sum-product on graphs with cycles, mean field, expectation propagation, Kikuchi methods, linear programming and semidefinite relaxations).

#### 3.6.1 General form of conjugate dual

Given a function  $A$ , the *conjugate dual function* to  $A$ , which we denote by  $A^*$ , is defined as follows:

$$A^*(\mu) := \sup_{\theta \in \Omega} \{\langle \mu, \theta \rangle - A(\theta)\}. \quad (3.42)$$

Here  $\mu \in \mathbb{R}^d$  is a fixed vector of so-called dual variables of the same dimension as  $\theta$ . Our choice of notation—i.e., using  $\mu$  again—is deliberately suggestive, in that these dual variables turn out to have a natural interpretation as mean parameters. Indeed, we have already mentioned one statistical interpretation of this variational problem (3.42); in particular, the right-hand side is the optimized value of the rescaled log



likelihood (3.38). Of course, this maximum likelihood problem only makes sense when the vector  $\mu$  belongs to the set  $\mathcal{M}$ ; an example is the vector of empirical moments  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \phi(X^i)$  induced by a data sample  $X_1^n = \{X^1, \dots, X^n\}$ . In our development, we consider the optimization problem (3.42) more broadly for any vector  $\mu \in \mathbb{R}^d$ . In this context, it is necessary to view  $A^*$  as a function taking values in the extended real line  $\mathbb{R}^* = \mathbb{R} \cup \{+\infty\}$ , as is standard in convex analysis (see Appendix A.2.5 for more details).

As we have previously intimated, the conjugate dual function (3.42) is very closely connected to entropy. Recall the definition of the Shannon entropy in (3.2). The main result of the following theorem is that when  $\mu \in \mathcal{M}^\circ$ , the value of the dual function  $A^*(\mu)$  is precisely the negative entropy of the exponential family distribution  $p_{\theta(\mu)}$ , where  $\theta(\mu)$  is the unique vector of canonical parameters satisfying the relation

$$\mathbb{E}_{\theta(\mu)}[\phi(X)] = \nabla A(\theta(\mu)) = \mu. \quad (3.43)$$

We will also find it essential to consider  $\mu \notin \mathcal{M}^\circ$ , in which case it is impossible to find canonical parameters satisfying the relation (3.43). In this case, the behavior of the supremum defining  $A^*(\mu)$  requires a more delicate analysis. In fact, denoting by  $\overline{\mathcal{M}}$  the closure of  $\mathcal{M}$ , it turns out that whenever  $\mu \notin \overline{\mathcal{M}}$ , then  $A^*(\mu) = +\infty$ . This fact is essential in the use of variational methods: it guarantees that any optimization problem involving the dual function can be reduced to an optimization problem over  $\mathcal{M}$ . Accordingly, a great deal of our discussion in the sequel will be on the structure of  $\mathcal{M}$  for various graphical models, and various approximations to  $\mathcal{M}$  for models in which its structure is overly complex.

More formally, the following theorem, proved in Appendix B.2, provides a precise characterization of the relation between  $A$  and its conjugate dual  $A^*$ :

**Theorem 2.**

(a) For any  $\mu \in \mathcal{M}^\circ$ , denote by  $\theta(\mu)$  the unique canonical parameter satisfying the dual matching condition (3.43). The conjugate dual

function  $A^*$  takes the following form:

$$A^*(\mu) = \begin{cases} -H(p_{\theta(\mu)}) & \text{if } \mu \in \mathcal{M}^\circ \\ +\infty & \text{if } \mu \notin \overline{\mathcal{M}}. \end{cases} \quad (3.44)$$

For any boundary point  $\mu \in \overline{\mathcal{M}} \setminus \mathcal{M}^\circ$  we have  $A^*(\mu) = \lim_{n \rightarrow +\infty} A^*(\mu^n)$  taken over any sequence  $\{\mu^n\} \subset \mathcal{M}^\circ$  converging to  $\mu$ .

(b) In terms of this dual, the log partition function has the following variational representation:

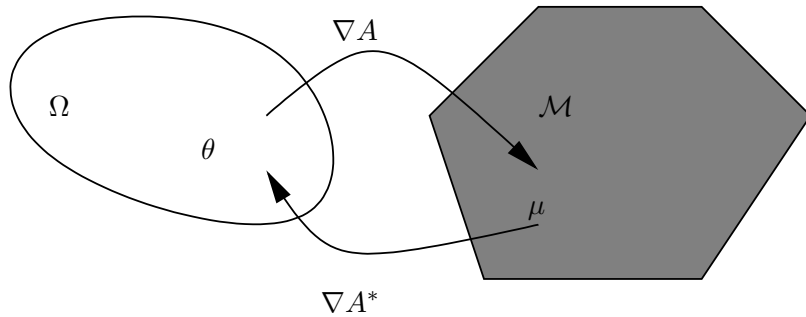
$$A(\theta) = \sup_{\mu \in \mathcal{M}} \{\langle \theta, \mu \rangle - A^*(\mu)\}. \quad (3.45)$$

(c) For all  $\theta \in \Omega$ , the supremum in equation (3.45) is attained uniquely at the vector  $\mu \in \mathcal{M}^\circ$  specified by the moment-matching conditions

$$\mu = \int_{\mathcal{X}^m} \phi(x) p_\theta(x) \nu(dx) = \mathbb{E}_\theta[\phi(X)]. \quad (3.46)$$

Theorem 2(a) provides a precise statement of the duality between the cumulant function  $A$  and entropy. It should be recognized that the dual function  $A^*$  is a slightly different object than the usual entropy (3.2): whereas the entropy maps density functions to real numbers (and so is a functional), the dual function  $A^*$  is an extended real-valued function on  $\mathbb{R}^d$ , finite only for valid mean parameters  $\mu \in \mathcal{M}$ . In particular, the value  $-A^*(\mu)$  corresponds to the optimum of the maximum entropy problem (3.3), where  $\mu \in \mathbb{R}^d$  parameterizes the constraint set. The event  $-A^*(\mu) = -\infty$  corresponds to infeasibility of the maximum entropy problem. This is an important point. Constrained optimization problems are defined both by the set being optimized over and the function being optimized. Given that the variational representation of the cumulant function in (3.45) takes the form of a maximization problem, we see that vectors  $\mu$  for which  $-A^*(\mu) = -\infty$  can certainly not be optima. Thus, as we see in (3.45), it suffices to maximize over the set  $\mathcal{M}$  instead of  $\mathbb{R}^d$ . This implies that the nature of the set  $\mathcal{M}$  plays a critical role in determining the complexity of computing the cumulant function.

Theorem 2 also clarifies the precise nature of the bijection between the sets  $\Omega$  and  $\mathcal{M}^\circ$ , which holds for any minimal exponential family. In particular, the gradient mapping  $\nabla A$  maps  $\Omega$  in a one-to-one manner onto  $\mathcal{M}^\circ$ , whereas the inverse mapping from  $\mathcal{M}^\circ$  to  $\Omega$  is given by the gradient  $\nabla A^*$  of the dual function (see Appendix B.3 for more details). Figure 3.8 provides an idealized illustration of this bijective



**Fig. 3.8.** Idealized illustration of the relation between the set  $\Omega$  of valid canonical parameters, and the set  $\mathcal{M}$  of valid mean parameters. The gradient mappings  $\nabla A$  and  $\nabla A^*$  associated with the conjugate dual pair  $(A, A^*)$  provide a bijective mapping between  $\Omega$  and the interior  $\mathcal{M}^\circ$ .

correspondence based on the gradient mappings  $(\nabla A, \nabla A^*)$ .

### 3.6.2 Some simple examples

Theorem 2 is best understood by working through some simple examples. Table 2 provides the conjugate dual pair,  $(A, A^*)$ , for several well-known exponential families of scalar random variables. For each family, the table also lists  $\Omega := \text{dom } A$ , as well as the set  $\mathcal{M}$ , which contains the effective domain of  $A^*$  (the set of values for which  $A^*$  is finite).

In the remainder of this section, we illustrate the basic ideas by working through two simple scalar examples in detail. To be clear, neither of these examples are interesting from a computational perspective—indeed, for most scalar exponential families, it is trivial to compute the mapping between canonical and mean parameters by direct methods. Nonetheless, they are useful in building intuition for

Family	$\Omega$	$A(\theta)$	$\mathcal{M}$	$A^*(\mu)$
Bernoulli	$\mathbb{R}$	$\log[1 + \exp(\theta)]$	$[0, 1]$	$\mu \log \mu + (1 - \mu) \log(1 - \mu)$
Gaussian Location family	$\mathbb{R}$	$\frac{1}{2}\theta^2$	$\mathbb{R}$	$\frac{1}{2}\mu^2$
Gaussian Location-scale	$\{(\theta_1, \theta_2) \mid \theta_2 < 0\}$	$-\frac{\theta_1^2}{4\theta_2} - \log(-2\theta_2)$	$\{(\mu_1, \mu_2) \mid \mu_2 - (\mu_1)^2 > 0\}$	$-\frac{1}{2} \log[\mu_2 - \mu_1^2]$
Exponential	$(-\infty, 0)$	$-\log(-\theta)$	$(0, +\infty)$	$-1 - \log(\mu)$
Poisson	$\mathbb{R}$	$\exp(\theta)$	$(0, +\infty)$	$\mu \log \mu - \mu$

**Table 3.2.** Conjugate dual relations of Theorem 2 for several well-known exponential families of scalar variables.

the consequences of Theorem 2. The reader interested only in the main thread may skip directly ahead to Section 4, where we use Theorem 2 in a more substantive manner—namely, to derive various approximate inference algorithms for multivariate graphical models.

**Example 12 (Conjugate duality for Bernoulli).** Consider a Bernoulli variable  $X \in \{0, 1\}$ : its distribution can be written as an exponential family with  $\phi(x) = x$ ,  $A(\theta) = \log(1 + \exp(\theta))$ , and  $\Omega = \mathbb{R}$ . In order to verify the claim in Theorem 2(a), let us compute the conjugate dual function  $A^*$  by direct methods. By the definition of conjugate duality (3.42), for any fixed  $\mu \in \mathbb{R}$ , we have

$$A^*(\mu) = \sup_{\theta \in \mathbb{R}} \{\theta \cdot \mu - \log(1 + \exp(\theta))\}. \tag{3.47}$$

Taking derivatives yields the stationary condition  $\mu = \exp(\theta)/[1 + \exp(\theta)]$ , which is simply the general matching condition (3.43) specialized to the Bernoulli model. When does this moment matching condition have a solution? If  $\mu \in (0, 1)$ , then some simple algebra shows that we may re-arrange to find the unique solution  $\theta(\mu) := \log[\mu/(1 - \mu)]$ . Since  $\mathcal{M}^\circ = (0, 1)$  for the Bernoulli model, the existence and uniqueness of this solution are particular consequences of Proposition 3 and Theorem 1. Since the objective

function (3.47) is strictly convex, the solution  $\theta(\mu)$  specifies the optimum; substituting  $\theta(\mu)$  into the objective equation (3.47) and simplifying yields that

$$\begin{aligned} A^*(\mu) &= \mu \log[\mu/(1-\mu)] - \log[1 + \frac{\mu}{1-\mu}] \\ &= \mu \log \mu + (1-\mu) \log(1-\mu), \end{aligned}$$

which is the negative entropy of Bernoulli variate with mean parameter  $\mu$ . We have thus verified Theorem 2(a) in the case that  $\mu \in (0, 1) = \mathcal{M}^\circ$ .

Now let us consider the case  $\mu \notin \overline{\mathcal{M}} = [0, 1]$ ; concretely, let us suppose that  $\mu > 1$ . In this case, there is no gradient stationary point in the optimization problem (3.47). Therefore, the supremum is specified by the limiting behavior as  $\theta \rightarrow \pm\infty$ . For  $\mu > 1$ , we claim that the objective function grows unboundedly as  $\theta \rightarrow +\infty$ . Indeed, by the convexity of  $A$ , we have  $A(0) = \log 2 \geq A(\theta) + A'(\theta)(-\theta)$ . Moreover  $A'(\theta) \leq 1$  for all  $\theta \in \mathbb{R}$ , from which we obtain the upper bound  $A(\theta) \leq \theta + \log 2$ , valid for  $\theta > 0$ . Consequently, for  $\theta > 0$ , we have

$$\mu \cdot \theta - \log[1 + \exp(\theta)] \geq (\mu - 1)\theta - \log 2,$$

showing that the objective function diverges as  $\theta \rightarrow +\infty$ , whenever  $\mu > 1$ . A similar calculation establishes the same claim for  $\mu < 0$ , showing that  $A^*(\mu) = +\infty$  for  $\mu \notin \overline{\mathcal{M}}$  and thus verifying the second part of Theorem 2(a). Finally, for the boundary points  $\mu = 0$  and  $\mu = 1$ , it can be verified by taking limits that  $A^*(0) = A^*(1) = 0$  for the Bernoulli model. (Indeed, the same statement holds more generally for any discrete model.)

Turning to the verification of Theorem 2(b), since  $A^*(\mu) = +\infty$  unless  $\mu \in [0, 1]$ , the optimization problem (3.45) reduces to

$$\max_{\mu \in [0, 1]} \{ \mu \cdot \theta - \mu \log \mu - (1-\mu) \log(1-\mu) \}.$$

Explicitly solving this concave maximization problem yields that its optimal value is  $\log[1 + \exp(\theta)]$ , which verifies the claim (3.45). Moreover, this same calculation shows that the optimum is attained uniquely at  $\mu(\theta) = \exp(\theta)/[1 + \exp(\theta)]$ , which verifies Theorem 2(c) for the Bernoulli model.



**Example 13 (Conjugate duality for exponential).** Consider the family of exponential distributions, represented as a regular exponential family with  $\phi(x) = \{x\}$ ,  $\Omega = (-\infty, 0)$ , and  $A(\theta) = -\log(-\theta)$ . From the conjugate dual definition (3.42), we have

$$A^*(\mu) = \sup_{\theta < 0} \{\mu \cdot \theta + \log(-\theta)\}. \quad (3.48)$$

In order to find the optimum, we take derivatives with respect to  $\theta$ . We thus obtain the stationary condition  $\mu = -1/\theta$ , which corresponds to the moment matching condition (3.43) specialized to this family. Hence, for all  $\mu \in \mathcal{M}^\circ = (0, \infty)$ , we have  $A^*(\mu) = -1 - \log(\mu)$ . It can be verified that the entropy of an exponential distribution with mean  $\mu > 0$  is given by  $-A^*(\mu)$ . This result confirms Theorem 2 for  $\mu \in \mathcal{M}^\circ$ . For  $\mu \notin \overline{\mathcal{M}}$ —that is, for  $\mu < 0$ —we see by inspection that the objective function  $\mu\theta + \log(-\theta)$  grows unboundedly as  $\theta \rightarrow -\infty$ , thereby demonstrating that  $A^*(\mu) = +\infty$  for all  $\mu \notin \overline{\mathcal{M}}$ . The remaining case is the boundary point  $\mu = 0$ , for which we have  $A^*(0) = +\infty$  from the definition (3.48). Note also that the negative entropy of the exponential distribution  $-1 - \log(\mu^n)$  tends to infinity for all sequence  $\mu^n \rightarrow 0^+$ , consistent with Theorem 2(a). Having computed the dual  $A^*$ , straightforward algebra shows that  $-\log(-\theta) = \sup_{\mu > 0} \{\mu \cdot \theta + 1 + \log(\mu)\}$ , with the optimum uniquely attained at  $\mu = -1/\theta$ . This calculation verifies parts (b) and (c) of Theorem 2 for the exponential variate. ♣

### 3.7 Challenges in high-dimensional models

From a computational perspective, the essential features of Theorem 2 are the representation (3.45) of the log partition function and the fact that the optimum is uniquely achieved at the mean parameters  $\mu = \mathbb{E}_\theta[\phi(X)]$ , as stated in part (c). It thus illuminates a key property of computing the log partition function  $A$ , as well as the mean parameters  $\mu$ : in principle, we can *compute both of these quantities by solving the variational problem* (3.45). Even more encouragingly, the optimization problem to be solved appears to be rather simple, at least from a superficial perspective; it is a finite-dimensional optimization problem over a convex set and the objective function is strictly concave and differentiable. Thus, the optimization lacks local optima or other un-

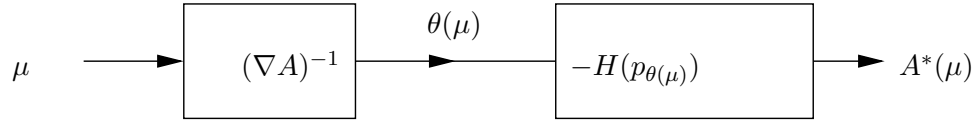
pleasant features. It is tempting, then, to assert that the problem of computing the log partition function and the associated mean parameters is now solved, since we have “reduced” it to a convex optimization problem.

In this context, the simple scalar examples of Table 3.2, for which the fundamental variational problem (3.45) had an explicit form and could be solved easily, are very misleading. For general multivariate exponential families, in contrast, there are two primary challenges associated with the variational representation:

- (a) in many cases, the constraint set  $\mathcal{M}$  of realizable mean parameters is extremely difficult to characterize in an explicit manner.
- (b) the negative entropy function  $A^*$  is defined indirectly—in a variational manner—so that it too typically lacks an explicit form.

For instance, to illustrate issue (a) concerning the nature of  $\mathcal{M}$ , for Markov random fields involving discrete random variables  $X \in \{0, 1, \dots, r - 1\}^m$ , the set  $\mathcal{M}$  is always a polytope, which we have referred to as a *marginal polytope*. In this case, at least in principle, the set  $\mathcal{M}$  can be characterized by some finite number of linear inequalities. However, for general graphs, the number of such inequalities grows rapidly with the graph size. Indeed, unless fundamental conjectures in complexity theory turn out to be false, it is not even possible to optimize a linear function over  $\mathcal{M}$  for a general discrete MRF. In addition to the complexity of the constraint set, issue (b) highlights that even evaluating the cost function at a single point  $\mu \in \mathcal{M}$ , let alone optimizing it over  $\mathcal{M}$ , is extremely difficult.

To understand the complexity inherent in evaluating the dual value  $A^*(\mu)$ , note that Theorem 2 provides only an implicit characterization of  $A^*(\mu)$  as the composition of mappings: first, the inverse mapping  $(\nabla A)^{-1} : \mathcal{M}^\circ \rightarrow \Omega$ , in which  $\mu$  maps to  $\theta(\mu)$ , corresponding to the exponential family member with mean parameters  $\mu$ ; and second, the mapping from  $\theta(\mu)$  to the negative entropy  $-H(p_{\theta(\mu)})$  of the associated exponential family density. This decomposition of the value  $A^*(\mu)$  is illustrated in Figure 3.9. Consequently, computing the dual value  $A^*(\mu)$  at some point  $\mu \in \mathcal{M}^\circ$  requires computing the inverse



**Fig. 3.9.** A block diagram decomposition of  $A^*$  as the composition of two functions. Any mean parameter  $\mu \in \mathcal{M}^\circ$  is first mapped back to an canonical parameter  $\theta(\mu)$  in the inverse image  $(\nabla A)^{-1}(\mu)$ . The value of  $A^*(\mu)$  corresponds to the negative entropy  $-H(p_{\theta(\mu)})$  of the associated exponential family density  $p_{\theta(\mu)}$ .

mapping  $(\nabla A)^{-1}(\mu)$ , in itself a non-trivial problem, and then evaluating the entropy, which requires high-dimensional integration for general graphical models. These difficulties motivate the use of approximations to  $\mathcal{M}$  and  $A^*$ . Indeed, as shown in the sections to follow, a broad class of methods for approximate marginalization are based on this strategy of finding an approximation to the exact variational principle, which is then often solved using some form of message-passing algorithm.



# 4

---

## Sum-product, Bethe-Kikuchi, and expectation-propagation

---

In this section, we begin our exploration of the variational interpretation of message-passing algorithms. We first discuss the sum-product algorithm, also known as the belief propagation algorithm. As discussed in Section 2.5, sum-product is an exact algorithm for tree-structured graphs, for which it can be derived as a divide-and-conquer algorithm. However, given the local form of the sum-product updates, there is no barrier to applying it to a graph with cycles, which yields the “loopy” form of belief propagation. In the presence of cycles, there are no general convergence or correctness guarantees associated with the sum-product algorithm, but it is nonetheless widely used to compute approximate marginals. The first part of this section describes the variational formulation of the sum-product updates in terms of the Bethe approximation. Although this approximation dates back to the work of Bethe [27], the connection to the sum-product algorithm was first elucidated by Yedidia, Freeman and Weiss [263, 264]. We then describe various natural generalizations of the Bethe approximation, including Kikuchi clustering and other hypergraph-based methods [130, 264]. Finally, we describe expectation-propagation algorithms [180, 172] and related moment-matching methods [181, 61, 112]; these are also varia-

tional methods based on Bethe-like approximations.

## 4.1 Sum-product and Bethe approximation

The simplest instantiation of the Bethe approximation applies to an undirected graphical model with potential functions involving at most pairs of variables; we refer to any such model as a *pairwise Markov random field*. In principle, by selectively introducing auxiliary variables, any undirected graphical model can be converted into an equivalent pairwise form to which the Bethe approximation can be applied; see Appendix E.2 for details of this procedure. It can also be useful to treat higher order interactions directly, which can be done using the approximations discussed in Section 4.2.

For the current section, let us assume that the given model is a pairwise Markov random field defined by a graph  $G = (V, E)$ . We also limit our discussion to discrete random vectors  $\{X_t \mid t \in V\}$ , in which each variable  $X_s$  takes values in the space  $\mathcal{X}_s = \{0, 1, \dots, r_s - 1\}$ , for which the sum-product algorithm is most widely used. The general variational principle (3.45) takes different forms, depending on the chosen form of sufficient statistics  $\phi$ . Recall from Section 3.4.1 our definition of the canonical overcomplete representation (3.34), using indicator functions for events  $\{X_s = j\}$  and  $\{X_s = j, X_t = k\}$ . Using these sufficient statistics, we define an exponential family of the form

$$p_{\theta}(x) \propto \exp \left\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \right\}, \quad (4.1)$$

where we have introduced the convenient short-hand notation

$$\theta_s(x_s) := \sum_j \theta_{s;j} \mathbb{I}_{s;j}(x_s), \text{ and} \quad (4.2a)$$

$$\theta_{st}(x_s, x_t) := \sum_{(j,k)} \theta_{st;jk} \mathbb{I}_{st;jk}(x_s, x_t). \quad (4.2b)$$

As previously discussed in Example 4, this particular parameterization is overcomplete or non-identifiable, because there exist entire affine subsets of canonical parameters that induce the same probability distribution over the random vector  $X$ .<sup>1</sup> Nonetheless, this overcom-

<sup>1</sup>A little calculation shows that the exponential family (4.1) has dimension  $d = \sum_{s \in V} r_s +$

pletteness is useful because the associated mean parameters are easily interpretable, since they correspond to singleton and pairwise marginal probabilities (see equations (3.35) and (3.36)). Using these mean parameters, it is convenient for various calculations to define the shorthand notations

$$\mu_s(x_s) := \sum_{j \in \mathcal{X}_s} \mu_{s;j} \mathbb{I}_{s;j}(x_s), \quad \text{and} \quad (4.3a)$$

$$\mu_{st}(x_s, x_t) := \sum_{(j,k) \in \mathcal{X}_s \times \mathcal{X}_t} \mu_{st;jk} \mathbb{I}_{st;jk}(x_s, x_t). \quad (4.3b)$$

Note that  $\mu_s$  is a  $|\mathcal{X}_s|$ -dimensional marginal distribution over  $X_s$ , whereas  $\mu_{st}$  is a  $|\mathcal{X}_s| \times |\mathcal{X}_t|$  matrix, representing a joint marginal over  $(X_s, X_t)$ . The *marginal polytope*  $\mathbb{M}(G)$  corresponds to the set of all singleton and pairwise marginals that are jointly realizable:

$$\mathbb{M}(G) := \{ \mu \in \mathbb{R}^d \mid \exists p \text{ with marginals } \mu_s(x_s), \mu_{st}(x_s, x_t) \} \quad (4.4)$$

#### 4.1.1 A tree-based outer bound to $\mathbb{M}(G)$

As previously discussed in Section 3.4.1, the polytope  $\mathbb{M}(G)$  can either be written as the convex hull of a finite number of vectors, one associated with each configuration  $x \in \mathcal{X}^m$ , or alternatively, as the intersection of a finite number of half-spaces. (See Figure 3.5 for an illustration.) But how to provide an explicit listing of these half-space constraints, also known as facets? In general, this problem is extremely difficult, so that we resort to listing only subsets of the constraints, thereby obtaining a polyhedral outer bound on  $\mathbb{M}(G)$ .

More specifically, consider a trial set of single node functions  $\tau_s$  and edge-based functions  $\tau_{st}$ . If these trial marginal distributions are to be globally realizable, they must of course be non-negative, and in addition, each singleton quantity  $\tau_s$  must satisfy the *normalization*

---

$\sum_{(s,t) \in E} r_s r_t$ . Given some fixed parameter vector  $\theta \in \mathbb{R}^d$ , suppose that we form a new canonical parameter  $\theta' \in \mathbb{R}^d$  by setting  $\theta'_{1;j} = \theta_{1;j} + C$  for all  $j \in \mathcal{X}$  where  $C \in \mathbb{R}$  is some fixed constant, and  $\theta'_\alpha = \theta_\alpha$  for all other indices  $\alpha$ . It is then straightforward to verify that  $p_\theta(x) = p_{\theta'}(x)$  for all  $x \in \mathcal{X}^m$ , so that  $\theta$  and  $\theta'$  describe the same probability distribution.

condition

$$\sum_{x_s} \tau_s(x_s) = 1. \quad (4.5)$$

Moreover, for each edge  $(s, t) \in E$ , the singleton  $\{\tau_s, \tau_t\}$  and pairwise quantities  $\tau_{st}$  must satisfy the *marginalization constraints*

$$\sum_{x'_t} \tau_{st}(x_s, x'_t) = \tau_s(x_s), \quad \forall x_s \in \mathcal{X}_s, \text{ and} \quad (4.6a)$$

$$\sum_{x'_s} \tau_{st}(x'_s, x_t) = \tau_t(x_t), \quad \forall x_t \in \mathcal{X}_t. \quad (4.6b)$$

These constraints define the set

$$\mathbb{L}(G) := \left\{ \tau \geq 0 \mid \begin{array}{l} \text{condition (4.5) holds for } s \in V, \text{ and} \\ \text{condition (4.6) holds for } (s, t) \in E \end{array} \right\}, \quad (4.7)$$

of *locally consistent* marginal distributions. Note that  $\mathbb{L}(G)$  is also a polytope, in fact a very simple one, since it is defined by  $\mathcal{O}(|V| + |E|)$  constraints in total.

What is the relation between the set  $\mathbb{L}(G)$  of locally consistent marginals and the set  $\mathbb{M}(G)$  of globally realizable marginals? On one hand, it is clear that  $\mathbb{M}(G)$  is always a subset of  $\mathbb{L}(G)$ , since any set of globally realizable marginals must satisfy the normalization (4.5) and marginalization (4.6) constraints. Apart from this inclusion relation, it turns out that for a graph with cycles, these two sets are rather different. However, for a tree  $T$ , the junction tree theorem, in the form of Proposition 1, guarantees that they are equivalent, as summarized in the following proposition.

**Proposition 4.** *The inclusion  $\mathbb{M}(G) \subseteq \mathbb{L}(G)$  holds for any graph. For a tree-structured graph  $T$ , the marginal polytope  $\mathbb{M}(T)$  is equal to  $\mathbb{L}(T)$ .*

*Proof.* Consider an element  $\mu$  of the full marginal polytope  $\mathbb{M}(G)$ : clearly, any such vector must satisfy the normalization and pairwise marginalization conditions defining the set  $\mathbb{L}(G)$ , from which we conclude that  $\mathbb{M}(G) \subseteq \mathbb{L}(G)$ . In order to demonstrate the reverse inclusion

for a tree-structured graph  $T$ , let  $\mu$  be an arbitrary element of  $\mathbb{L}(T)$ ; we need to show that  $\mu \in \mathbb{M}(T)$ . By definition of  $\mathbb{L}(T)$ , the vector  $\mu$  specifies a set of locally consistent singleton marginals  $\mu_s$  for vertices  $s \in V$  and pairwise marginals  $\mu_{st}$  for edges  $(s, t) \in E$ . By the junction tree theorem, we may use them to form a distribution, Markov with respect to the tree, as follows

$$p_\mu(x) := \prod_{s \in V} \mu_s(x_s) \prod_{(s,t) \in E} \frac{\mu_{st}(x_s, x_t)}{\mu_s(x_s) \mu_t(x_t)}. \quad (4.8)$$

(We take  $0/0 := 0$  in cases of zeros in the elements of  $\mu$ .) It is a consequence of the junction tree theorem or can be verified directly via an inductive “leaf-stripping” argument that with this choice of  $p_\mu$ , we have  $\mathbb{E}_{p_\mu}[\mathbb{I}_j(X_s)] = \mu_s(x_s)$  for all  $s \in V$  and  $j \in \mathcal{X}_s$ , as well as  $\mathbb{E}_{p_\mu}[\mathbb{I}_{jk}(X_s, X_t)] = \mu_{st}(x_s, x_t)$  for all  $(s, t) \in E$ , and  $(j, k) \in \mathcal{X}_s \times \mathcal{X}_t$ . Therefore, the distribution (4.8) provides a constructive certificate of the membership  $\mu \in \mathbb{M}(T)$ , which establishes that  $\mathbb{L}(T) = \mathbb{M}(T)$ .  $\square$

For a graph  $G$  with cycles, in sharp contrast to the tree case, the set  $\mathbb{L}(G)$  is a strict outer bound on  $\mathbb{M}(G)$ , in that there exist vectors  $\tau \in \mathbb{L}(G)$  that do *not* belong to  $\mathbb{M}(G)$ , for which reason we refer to members  $\tau$  of  $\mathbb{L}(G)$  as *pseudomarginals*. The following example illustrates the distinction between globally realizable marginals and pseudomarginals.

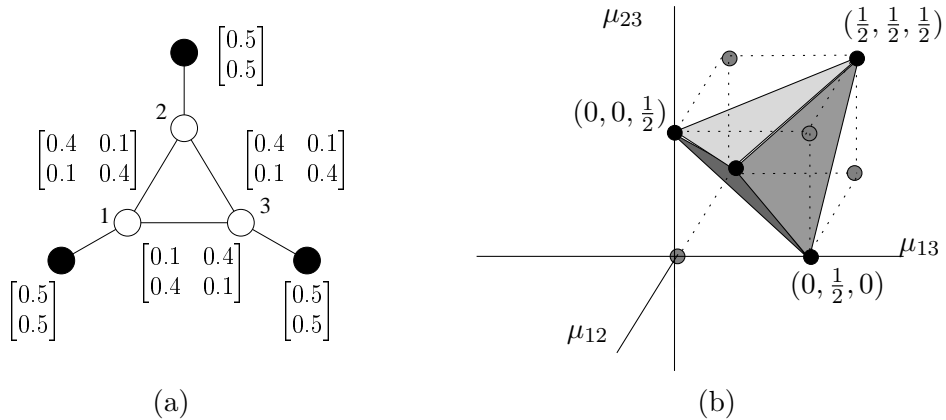
**Example 14** ( $\mathbb{L}(G)$  versus  $\mathbb{M}(G)$ ). Let us explore the relation between the two sets on the simplest graph for which they fail to be equivalent—namely, the single cycle on three vertices, denoted by  $C_3$ . Considering the binary random vector  $X \in \{0, 1\}^3$ , note that each singleton pseudomarginal  $\tau_s$ , for  $s = 1, 2, 3$ , can be viewed as a  $1 \times 2$  vector, whereas each pairwise pseudomarginal  $\tau_{st}$ , for  $(s, t) \in \{(12), (13), (23)\}$  can be viewed as a  $2 \times 2$  matrix. We define the family of pseudomarginals

$$\tau_s(x_s) := \begin{bmatrix} 0.5 & 0.5 \end{bmatrix}, \quad \text{and} \quad (4.9a)$$

$$\tau_{st}(x_s, x_t) := \begin{bmatrix} \beta_{st} & 0.5 - \beta_{st} \\ 0.5 - \beta_{st} & \beta_{st} \end{bmatrix}, \quad (4.9b)$$

where for each edge  $(s, t) \in E$ , the quantity  $\beta_{st} \in \mathbb{R}$  is a parameter to be specified.

We first observe that for any  $\beta_{st} \in [0, 0.5]$ , these pseudomarginals satisfy the normalization (4.5) and marginalization constraints (4.6), so the associated pseudomarginals (4.9) belong to  $\mathbb{L}(C_3)$ . As a particular choice, consider the collection  $\tau$  of pseudomarginals generated by setting  $\beta_{12} = \beta_{23} = 0.4$ , and  $\beta_{13} = 0.1$ , as illustrated in Figure 4.1(a). With these settings, the vector  $\tau$  is an element of  $\mathbb{L}(C_3)$ ; however, as a candidate set of global marginal distributions, certain features of the collection  $\tau$  should be suspicious. In particular, according to the putative marginals  $\tau$ , the events  $\{X_1 = X_2\}$  and  $\{X_2 = X_3\}$  should each hold with probability 0.8, whereas the event  $\{X_1 = X_3\}$  should only hold with probability 0.2. At least intuitively, this set-up appears likely to violate some type of global constraint.



**Fig. 4.1.** (a) A set of pseudomarginals associated with the nodes and edges of the graph: setting  $\beta_{12} = \beta_{23} = 0.4$  and  $\beta_{13} = 0.1$  in equation (4.9) yields a pseudomarginal vector  $\tau$  which, though locally consistent, is not globally consistent. (b) Marginal polytope  $\mathbb{M}(C_3)$  for the three node cycle; in a minimal exponential representation, it is a 6-dimensional object. Illustrated here is the slice  $\{\mu_1 = \mu_2 = \mu_3 = \frac{1}{2}\}$ , as well as and the outer bound  $\mathbb{L}(C_3)$ , also for this particular slice.

In order to prove the global invalidity of  $\tau$ , we first specify the constraints that actually define the marginal polytope  $\mathbb{M}(G)$ . For ease of illustration, let us consider a particular slice—namely,  $\mu_1 = \mu_2 =$

$\mu_3 = \frac{1}{2}$  of the marginal polytope. Viewed in this slice, the constraints defining  $\mathbb{L}(G)$  reduce to  $0 \leq \beta_{st} \leq \frac{1}{2}$  for all edges  $(s, t)$ , so that the set is simply a three-dimensional cube, as drawn with dotted lines in Figure 4.1(b). It can be shown that the sliced version of  $\mathbb{M}(G)$  is defined by these box constraints, and in addition the following cycle inequalities

$$\mu_{12} + \mu_{23} - \mu_{13} \leq \frac{1}{2}, \quad \mu_{12} - \mu_{23} + \mu_{13} \leq \frac{1}{2}, \quad (4.10a)$$

$$-\mu_{12} + \mu_{23} + \mu_{13} \leq \frac{1}{2}, \quad \text{and} \quad \mu_{12} + \mu_{23} + \mu_{13} \geq \frac{1}{2}. \quad (4.10b)$$

(See Example 42 in the sequel for the derivation of these cycle inequalities.) As illustrated by the shaded region in Figure 4.1(b), the marginal polytope  $\mathbb{M}(C_3)$  is strictly contained within the scaled cube  $[0, \frac{1}{2}]^3$ .

To conclude the discussion of our example, note that the pseudo-marginal vector  $\tau$  specified by  $\beta_{12} = \beta_{23} = 0.4$  and  $\beta_{13} = 0.1$  fails to satisfy the cycle inequalities (4.10), since in particular, we have the violation

$$\beta_{12} + \beta_{23} - \beta_{13} = 0.4 + 0.4 - 0.1 > \frac{1}{2}.$$

Therefore, the vector  $\tau$  is an instance of a set of pseudomarginals valid under the Bethe approximation, but which could never arise from a true probability distribution. ♣

#### 4.1.2 Bethe entropy approximation

We now turn to the second variational ingredient that underlies the sum-product algorithm, namely an approximation to the dual function (or negative entropy). As with the outer bound  $\mathbb{L}(G)$  in the marginal polytope, the Bethe entropy approximation is also tree-based.

For a general MRF based on a graph with cycles, the negative entropy  $A^*$ —as a function of *only* the mean parameters  $\mu$ —typically lacks a closed form expression. An important exception to this rule is the case of a tree-structured Markov random field, for which the entropy decomposes in terms of local entropies associated with the edges and nodes of the graph. In order to derive this decomposition, recall from the proof of Proposition 4 the factorization (4.8) of any tree-structured

MRF distribution in terms of marginal distributions  $\mu_s$  and  $\mu_{st}$  on the node and edges, respectively, of the tree. These marginal distributions correspond to the mean parameters under the canonical overcomplete sufficient statistics (3.34). Thus, for a tree-structured MRF, we can compute the (negative) dual value  $-A^*(\mu)$  directly, simply by computing the entropy  $H(p_\mu)$  of the distribution (4.8). Denoting by  $\mathbb{E}_\mu$  the expectation under the distribution (4.8), we obtain

$$\begin{aligned} H(p_\mu) = -A^*(\mu) &= \mathbb{E}_\mu[-\log p_\mu(X)] \\ &= \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E} I_{st}(\mu_{st}). \end{aligned} \quad (4.11)$$

The different terms in this expansion are the singleton entropy

$$H_s(\mu_s) := - \sum_{x_s \in \mathcal{X}_s} \mu_s(x_s) \log \mu_s(x_s) \quad (4.12)$$

for each node  $s \in V$ , and the mutual information

$$I_{st}(\mu_{st}) := \sum_{(x_s, x_t) \in \mathcal{X}_s \times \mathcal{X}_t} \mu_{st}(x_s, x_t) \log \frac{\mu_{st}(x_s, x_t)}{\mu_s(x_s) \mu_t(x_t)} \quad (4.13)$$

for each edge  $(s, t) \in E$ . Consequently, for a tree-structured graph, the dual function  $A^*$  can be expressed as an explicit and easily computable function of the mean parameters  $\mu$ .

With this background, the Bethe approximation to the entropy of an MRF with cycles is easily described: it simply assumes that decomposition (4.11) is approximately valid for a graph with cycles. This assumption yields the *Bethe entropy approximation*

$$-A^*(\tau) \approx H_{\text{Bethe}}(\tau) := \sum_{s \in V} H_s(\tau_s) - \sum_{(s,t) \in E} I_{st}(\tau_{st}). \quad (4.14)$$

An important fact, central in the derivation of the sum-product algorithm, is that this approximation (4.14) can be evaluated for any set of pseudomarginals  $\tau_s$  and  $\tau_{st}$  that belong to  $\mathbb{L}(G)$ . For this reason, our change in notation—from  $\mu$  for exact marginals to  $\tau$  for pseudomarginals—is deliberate.

We note in passing that Yedidia et al. [263, 264] used an alternative form of the Bethe entropy approximation (4.14), one which can be



obtained via the relation  $I_{st}(\tau_{st}) = H_s(\tau_s) + H_t(\tau_t) - H_{st}(\tau_{st})$ , where  $H_{st}$  is the joint entropy defined by the pseudomarginal  $\tau_{st}$ . Doing so and performing some algebraic manipulation yields

$$H_{Bethe}(\tau) = - \sum_{s \in V} (d_s - 1) H_s(\tau_s) + \sum_{(s,t) \in E} H_{st}(\tau_{st}), \quad (4.15)$$

where  $d_s$  corresponds to the number of neighbors of node  $s$  (i.e., the degree of node  $s$ ). However, the symmetric form (4.14) turns out to be most natural for our development in the sequel.

### 4.1.3 Bethe variational problem and sum-product

We now have the two ingredients needed to construct the Bethe approximation to the exact variational principle (3.45) from Theorem 2:

- the set  $\mathbb{L}(G)$  of locally consistent pseudomarginals (4.7) is a convex (polyhedral) outer bound on the marginal polytope  $\mathbb{M}(G)$ ; and
- the Bethe entropy (4.14) is an approximation of the exact dual function  $A^*(\tau)$ .

By combining these two ingredients, we obtain the *Bethe variational problem* (BVP):

$$\max_{\tau \in \mathbb{L}(G)} \left\{ \langle \theta, \tau \rangle + \sum_{s \in V} H_s(\tau_s) - \sum_{(s,t) \in E} I_{st}(\tau_{st}) \right\} = \max_{\tau \in \mathbb{L}(G)} \left\{ \langle \theta, \tau \rangle + H_{Bethe}(\tau) \right\}. \quad (4.16)$$

Note that this problem has a very simple structure: the cost function is given in closed form, it is concave and differentiable, and the constraint set  $\mathbb{L}(G)$  is a polytope specified by a small number of constraints. Given this special structure, one might suspect that there should exist a relatively simple algorithm for solving this optimization problem (4.16). Indeed, the sum-product algorithm turns out to be exactly such a method.

In order to develop this connection between the variational problem (4.16) and the sum-product algorithm, let  $\lambda_{ss}$  be a Lagrange multiplier associated with the normalization constraint  $C_{ss}(\tau) = 0$ , where

$$C_{ss}(\tau) := 1 - \sum_{x_s} \tau_s(x_s). \quad (4.17)$$

Moreover, for each direction  $t \rightarrow s$  of each edge and each  $x_s \in \mathcal{X}_s$ , define the constraint function

$$C_{ts}(x_s; \tau) := \tau_s(x_s) - \sum_{x_t} \tau_{st}(x_s, x_t), \quad (4.18)$$

and let  $\lambda_{st}(x_s)$  be a Lagrange multiplier associated with the constraint  $C_{ts}(x_s; \tau) = 0$ . These Lagrange multipliers turn out to be closely related to the sum-product messages, in particular via the relation  $M_{ts}(x_s) \propto \exp(\lambda_{ts}(x_s))$ .

We then consider the Lagrangian corresponding to the Bethe variational problem (4.16):

$$\begin{aligned} \mathcal{L}(\tau, \lambda; \theta) &:= \langle \theta, \tau \rangle + H_{\text{Bethe}}(\tau) \\ &+ \sum_{(s,t) \in E} \left[ \sum_{x_s} \lambda_{ts}(x_s) C_{ts}(x_s; \tau) + \sum_{x_t} \lambda_{st}(x_t) C_{st}(x_t; \tau) \right] + \sum_{s \in V} \lambda_{ss} C_{ss}(\tau). \end{aligned} \quad (4.19)$$

With these definitions, the connection between sum-product and the BVP is made precise by the following result, due to Yedidia et al. [264]:

**Theorem 3 (Sum-product and the Bethe problem).** *The sum-product updates are a Lagrangian method for attempting to solve the Bethe variational problem:*

- (a) *For any graph  $G$ , any fixed point of the sum-product updates specifies a pair  $(\tau^*, \lambda^*)$  such that*

$$\nabla_{\tau} \mathcal{L}(\tau^*, \lambda^*; \theta) = 0, \quad \text{and} \quad \nabla_{\lambda} \mathcal{L}(\tau^*, \lambda^*; \theta) = 0. \quad (4.20)$$

- (b) *For a tree-structured Markov random field (MRF), the Lagrangian equations (4.20) have a unique solution  $(\tau^*, \lambda^*)$ , where the elements of  $\tau^*$  correspond to the exact singleton and pairwise marginal distributions of the MRF. Moreover, the optimal value of the BVP is equal to the cumulant function  $A(\theta)$ .*

**Remark:** It should be noted that the Lagrangian formulation (4.19) is a partial one, because it assigns Lagrange multipliers to the normalization  $C_{ss}$  and marginalization  $C_{ts}$  constraints, but deals with the

non-negativity constraints implicitly. However, any optimum  $\tau^*$  of the Bethe variational principle with strictly positive elements must satisfy the Lagrangian conditions (4.20) from Theorem 3(a). To see this fact, note that if  $\lambda^*$  is Lagrange multiplier vector for the BVP, then any optimal solution  $\tau^*$  must maximize the Lagrangian  $\mathcal{L}(\tau, \lambda^*)$  over the positive orthant  $\tau \geq 0$  (see Bertsekas [19]). If the optimal solution  $\tau^*$  has strictly positive elements, then a necessary condition for Lagrangian optimality is the zero-gradient condition  $\nabla_{\tau} \mathcal{L}(\tau^*, \lambda^*) = 0$ , as claimed. Moreover, for graphical models where all configurations are given strictly positive mass (such as graphical models in exponential form with finite  $\theta$ ), the sum-product messages stay bounded strictly away from zero [199], so that there is always an optimum  $\tau^*$  with strictly positive elements. In this case, Theorem 3(a) guarantees that any sum-product fixed point satisfies the Lagrangian conditions necessary to be an optimum of the Bethe variational principle. For graphical models in which some configurations are assigned zero mass, further care is required in connecting fixed points to local optima; we refer the reader to Yedidia et al. [264] for details.

*Proof.* (a) Computing the partial derivative  $\nabla_{\tau} \mathcal{L}(\tau; \lambda)$  and setting it to zero yields the relations:

$$\log \tau_s(x_s) = \lambda_{ss} + \theta_s(x_s) + \sum_{t \in N(s)} \lambda_{ts}(x_s) \quad (4.21a)$$

$$\log \frac{\tau_{st}(x_s, x_t)}{\left[ \sum_{x_s} \tau_{st}(x_s, x_t) \right] \left[ \sum_{x_t} \tau_{st}(x_s, x_t) \right]} = \theta_{st}(x_s, x_t) - \lambda_{ts}(x_s) - \lambda_{st}(x_t). \quad (4.21b)$$

The condition  $\nabla_{\lambda} \mathcal{L}(\tau; \lambda)$  is equivalent to  $C_{ts}(x_s; \tau) = 0$  and  $C_{ss}(\tau) = 0$ . Using equation (4.21a) and the fact that the marginalization condition  $C_{ts}(x_s; \tau) = 0$  holds, we can re-arrange equation (4.21b) to obtain:

$$\begin{aligned} \log \tau_{st}(x_s, x_t) &= \lambda_{ss} + \lambda_{tt} + \theta_{st}(x_s, x_t) + \theta_s(x_s) + \theta_t(x_t) \\ &+ \sum_{u \in N(s) \setminus t} \lambda_{us}(x_s) + \sum_{u \in N(t) \setminus s} \lambda_{ut}(x_t). \end{aligned} \quad (4.22)$$

So as to make explicit the connection to the sum-product algorithm, let us define, for each directed edge  $t \rightarrow s$ , an  $r_s$ -vector of “messages”

$$M_{ts}(x_s) = \exp(\lambda_{ts}(x_s)), \quad \text{for all } x_s \in \{0, 1, \dots, r_s - 1\}.$$

With this notation, we can then write an equivalent form of equation (4.21a) as follows:

$$\tau_s(x_s) = \kappa \exp(\theta_s(x_s)) \prod_{t \in N(s)} M_{ts}(x_s). \quad (4.23)$$

Similarly, we have an equivalent form of equation (4.22):

$$\begin{aligned} \tau_{st}(x_s, x_t) &= \kappa' \exp(\theta_{st}(x_s, x_t) + \theta_s(x_s) + \theta_t(x_t)) \\ &\times \prod_{u \in N(s) \setminus t} M_{us}(x_s) \prod_{u \in N(t) \setminus s} M_{ut}(x_t). \end{aligned} \quad (4.24)$$

Here  $\kappa, \kappa'$  are positive constants (dependent on  $\lambda_{ss}, \lambda_{tt}$ ), chosen so that the pseudomarginals satisfy normalization conditions. Note that  $\tau_s$  and  $\tau_{st}$  so defined are non-negative.

To conclude, we need to adjust the Lagrange multipliers or messages so that the constraint  $\sum_{x_s} \tau_{st}(x_s, x_t) = \tau_s(x_s)$  is satisfied for every edge. Using equations (4.23) and (4.24) and performing some algebra, the end result is

$$M_{ts}(x_s) \propto \sum_{x_t} \exp\{\theta_{st}(x_s, x_t) + \theta_t(x_t)\} \prod_{u \in N(t) \setminus s} M_{ut}(x_t), \quad (4.25)$$

which is equivalent to the familiar sum-product update (2.8). By construction, any fixed point  $M^*$  of these updates (4.25) specifies a pair  $(\tau^*, \lambda^*)$  that satisfies the stationary conditions (4.20).

(b) See Appendix B.4 for the proof of this claim.  $\square$

The connection between the sum-product algorithm and the Bethe variational principle has a number of important consequences. First, it provides a principled basis for applying the sum-product algorithm for graphs with cycles, namely as a particular type of iterative method for attempting to satisfy Lagrangian conditions. It should be noted, however, that this connection between sum-product and the Bethe problem in itself provides no guarantees on the convergence of the sum-product updates on graphs with cycles. Indeed, whether or not the algorithm converges depends both on the potential strengths and the topology of the graph. In the standard scheduling of the messages, each

node applies equation (4.25) in parallel. Other more global schemes for message-passing are possible, and commonly used in certain applications like error-control coding [e.g., 164]; some recent work has also studied adaptive schedules for message-passing [72, 221]. Tatikonda and Jordan [225] established an elegant connection between convergence of parallel updates and Gibbs measures on the infinitely unwrapped computation tree, thereby showing that sufficient conditions for convergence can be obtained from classical conditions for uniqueness of Gibbs measures (e.g., Dobrushin’s condition or Simon’s condition [87]). In subsequent work, other researchers [115, 175, 199] have used various types of contraction arguments to obtain sharper conditions for convergence, or uniqueness of fixed points [106]. In certain For suitably weak potentials, Dobrushin-type conditions and related contraction arguments guarantee both convergence of the updates, and as a consequence, uniqueness of the associated fixed point. A parallel line of work [265, 249, 108] has explored alternatives to sum-product that are guaranteed to converge, albeit at the price of increased computational cost. However, with the exception of trees and other special cases [184, 163, 107], the Bethe variational problem is usually a non-convex problem, in that  $H_{\text{Bethe}}$  fails to be concave. As a consequence, there are frequently local optima, so that even when using a convergent algorithm, there are no guarantees that it will find the global optimum.

For each  $\theta \in \mathbb{R}^d$ , let  $A_{\text{Bethe}}(\theta)$  denote the optimal value of the Bethe variational problem (4.16). Theorem 3(b) states for any tree-structured problem, we have the equality  $A_{\text{Bethe}}(\theta) = A(\theta)$  for all  $\theta \in \mathbb{R}^d$ . Given this equivalence, it is natural to consider the relation between  $A_{\text{Bethe}}(\theta)$  and the cumulant function  $A(\theta)$  for general graphs. In general, the Bethe value  $A_{\text{Bethe}}(\theta)$  is simply an approximation to the cumulant function  $A(\theta)$ . Unlike the mean field methods to be discussed in Section 5, it is not guaranteed to provide a lower bound on the cumulant function. As will be discussed at more length in Section 7, Wainwright et al. [241] showed that “convexified” forms of the Bethe variational principle are guaranteed to yield upper bounds on the cumulant function for any graphical model. On the other hand, Sudderth et al. [219] show that  $A_{\text{Bethe}}(\theta)$  is a lower bound on the cumulant function  $A(\theta)$  for

certain classes of attractive graphical models. Such models, in which the interactions encourage random variables to agree with one another, are common in computer vision and other applications in spatial statistics. This lower-bounding property is closely related to the connection between the Bethe approximation and loop series expansions [48], discussed in Section 4.1.6.

Another important consequence of the Bethe/sum-product connection is in suggesting a number of avenues for improving upon the ordinary sum-product algorithm, via progressively better approximations to the entropy function and outer bounds on the marginal polytope. We turn to discussion of a class of such generalized sum-product algorithms beginning in Section 4.2.

#### 4.1.4 Inexactness of Bethe and sum-product

In this section, we explore some aspects of the inexactness of the sum-product algorithm. From a variational perspective, the inexactness stems from the two approximations made in setting up Bethe variational principle:

- (a) replacing the marginal polytope  $\mathbb{M}(G)$  by the polyhedral outer bound  $\mathbb{L}(G)$  and
- (b) the Bethe entropy  $H_{Bethe}$  as an approximation to the exact entropy as a function of the mean parameters.

We begin by considering the Bethe entropy approximation, and its potential inexactness:

**Example 15 (Inexactness of  $H_{Bethe}$ ).** Consider the fully connected graph  $K_4$  on four vertices, and the collection of singleton and pairwise marginal distributions given by

$$\mu_s(x_s) = \begin{bmatrix} 0.5 & 0.5 \end{bmatrix} \quad \text{for } s = 1, 2, 3, 4 \quad (4.26a)$$

$$\mu_{st}(x_s, x_t) = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \quad \forall (s, t) \in E. \quad (4.26b)$$

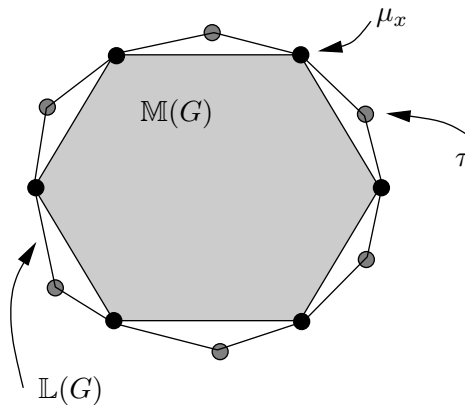
It can be verified that these marginals are globally valid, generated in particular by the distribution that places mass 0.5 on each of the configurations  $[0 \ 0 \ 0 \ 0]$  and  $[1 \ 1 \ 1 \ 1]$ . Let us calculate the Bethe

entropy approximation. Each of the four singleton entropies are given by  $H_s(\mu_s) = \log 2$ , and each of the six (one for each edge) mutual information terms are given by  $I_{st}(\mu_{st}) = \log 2$ , so that the Bethe entropy is given by

$$H_{Bethe}(\mu) = 4 \log 2 - 6 \log 2 = -2 \log 2 < 0,$$

which cannot be a true entropy. In fact, for this example, the true entropy (or value of the negative dual function) is given by  $-A^*(\mu) = \log 2 > 0$ . ♣

In addition to the inexactness of  $H_{Bethe}$  as an approximation to the negative dual function, the Bethe variational principle also involves relaxing the marginal polytope  $\mathbb{M}(G)$  to the first-order constraint set  $\mathbb{L}(G)$ . As illustrated in Example 14, the inclusion  $\mathbb{M}(C_3) \subseteq \mathbb{L}(C_3)$  holds *strictly* for the 3-node cycle  $C_3$ . The constructive procedure of Example 14 can be substantially generalized to show that the inclusion  $\mathbb{M}(G) \subset \mathbb{L}(G)$  holds strictly for any graph  $G$  with cycles. Figure 4.2



**Fig. 4.2.** Highly idealized illustration of the relation between the marginal polytope  $\mathbb{M}(G)$  and the outer bound  $\mathbb{L}(G)$ . The set  $\mathbb{L}(G)$  is always an outer bound on  $\mathbb{M}(G)$ , and the inclusion  $\mathbb{M}(G) \subset \mathbb{L}(G)$  is strict whenever  $G$  has cycles. Both sets are polytopes and so can be represented either as the convex hull of a finite number of extreme points, or as the intersection of a finite number of half-spaces, known as facets.

provides a highly idealized illustration<sup>2</sup> of the relation between  $\mathbb{M}(G)$  and  $\mathbb{L}(G)$ : both sets are polytopes, and for a graph with cycles,  $\mathbb{M}(G)$  is always strictly contained within the outer bound  $\mathbb{L}(G)$ .

Both sets are polytopes, and consequently can be represented either as the convex hull of a finite number of extreme points, or as the intersection of a finite number of half-spaces, known as facets. Letting  $\phi$  be a short-hand for the full vector of indicator functions in the standard overcomplete representation (3.34), the marginal polytope has the convex hull representation  $\mathbb{M}(G) = \text{conv}\{\phi(x) \mid x \in \mathcal{X}\}$ . Since the indicator functions are  $\{0, 1\}$ -valued, all of its extreme points consist of  $\{0, 1\}$  elements, of the form  $\mu_x := \phi(x)$  for some  $x \in \mathcal{X}^m$ ; there are a total of  $|\mathcal{X}^m|$  such extreme points. However, with the exception of tree-structured graphs, the number of facets for  $\mathbb{M}(G)$  is not known in general, even for relatively simple cases like the Ising model (see the book [66] for background on the cut or correlation polytope, which is equivalent to the marginal polytope for an Ising model.) However, the growth must be super-polynomial in the graph size, unless certain widely believed conjectures in computational complexity are false. On the other hand, the polytope  $\mathbb{L}(G)$  has a polynomial number of facets, upper bounded by any graph by  $\mathcal{O}(rm + r^2|E|)$ . It has more extreme points than  $\mathbb{M}(G)$ , since in addition to all the integral extreme points  $\{\mu_x, x \in \mathcal{X}^m\}$ , it includes other extreme points  $\tau \in \mathbb{L}(G) \setminus \mathbb{M}(G)$  that contain fractional elements (see Section 8.4 for further discussion of integral versus fractional vertices.) The number of vertices of  $\mathbb{L}(G)$  is not known in general.

The strict inclusion of  $\mathbb{M}(G)$  within  $\mathbb{L}(G)$ —and the fundamental role of the latter set in the Bethe variational problem (4.16)—leads to the following question: *do solutions to the Bethe variational problem ever fall into the gap  $\mathbb{M}(G) \setminus \mathbb{L}(G)$ ?* The optimistically inclined would hope that these points would somehow be excluded as optima of the Bethe variational problem, but such hope turns out to be misguided. In fact, for every element  $\tau$  of  $\mathbb{L}(G)$ , it is possible to construct a distribu-

---

<sup>2</sup>In particular, this picture is misleading in that it suggests that  $\mathbb{L}(G)$  has more facets and more vertices than  $\mathbb{M}(G)$ ; in fact, the polytope  $\mathbb{L}(G)$  has fewer facets and more vertices, but this is difficult to convey in a two-dimensional representation.



tion  $p_\theta$  such that if the sum-product algorithm is run on the problem, then  $\tau$  arises from the messages defined by a sum-product fixed point. In order to understand this fact, we need to describe the reparameterization interpretation of the sum-product algorithm [237], to which we now turn.

#### 4.1.5 Bethe optima and reparameterization

One view of the junction tree algorithm, as described in Section 2.5.2, is as follows: taking as input a set of potential functions on the cliques of some graph, it returns as output an *alternative factorization* of the same distribution in terms of local marginal distributions on the cliques and separator sets of a junction tree. In the special case of an ordinary tree, the alternative factorization is a product of local marginals at single nodes and edges of the tree, as in equation (4.8). Indeed, the sum-product algorithm for trees can be understood as an efficient method for computing this alternative parameterization.

It turns out that the same interpretation applies to arbitrary graphs with cycles: more precisely, any fixed point of the sum-product algorithm—and even more generally, any local optimum of the Bethe variational principle—specifies a reparameterization of the original distribution  $p_\theta$ . More formally, we have [237]:

**Proposition 5** (Reparameterization properties of Bethe approximation). *Letting  $\tau^* = \{\tau_s^*, s \in V\} \cup \{\tau_{st}^*, (s, t) \in E\}$  denote any optimum of the Bethe variational principle, consider the distribution defined as*

$$p_{\tau^*}(x) := \frac{1}{Z(\tau^*)} \prod_{s \in V} \tau_s^*(x_s) \prod_{(s,t) \in E} \frac{\tau_{st}^*(x_s, x_t)}{\tau_s^*(x_s) \tau_t^*(x_t)}, \quad (4.27)$$

*Then  $p_{\tau^*}$  is a reparameterization of  $p_\theta$ —that is,  $p_{\tau^*}(x) = p_\theta(x)$  for all  $x \in \mathcal{X}^m$ .*

Note that this type of reparameterization is possible only because the exponential family is defined by an overcomplete set of sufficient statistics, involving the indicator functions (3.34). The reparameterization (4.27) is the analog of the tree-structured factorization (4.8), but as applied to a graph with cycles. In contrast to the tree case, the

normalization constant  $Z(\tau^*)$  is not equal to one (in general). Whereas the junction tree theorem guarantees this reparameterization for any (junction) tree, it is not immediately obvious that even one reparameterization need exist for a general graph. In fact, the result establishes that every graph has at least one such reparameterization, and some graphs may have multiple reparameterizations of the form (4.27); in particular, this is the case for any problem for which the BVP has multiple optima. Moreover, the reparameterization viewpoint provides some insight into the *approximation error*: that is, the difference between the exact marginals  $\mu_s$  of  $p_\theta(x)$  and the approximations  $\tau_s^*$  computed by the sum-product algorithm. Indeed, using equation (4.27), it is possible to derive an exact expression for the error in the sum-product algorithm, as well as computable error bounds, as described in more detail in Wainwright et al. [237].

We now show how the reparameterization characterization (4.27) enables us to specify, for *any* pseudomarginal  $\tau$  in the interior of  $\mathbb{L}(G)$ , a distribution  $p_\theta$  for which  $\tau$  is a fixed point of the sum-product algorithm. The following example illustrates this construction.

**Example 16 (Fooling the sum-product algorithm).** Let us return to the simplest graph for which sum-product is not exact—namely, a single cycle with three nodes (see Example 14). Consider candidate marginal distributions  $\{\tau_s, \tau_{st}\}$  of the form illustrated in Figure 4.1(a), with  $\beta_{12} = \beta_{23} = 0.4$  and  $\beta_{13} = 0.1$ . As discussed in Example 14, this setting yields a set of pseudomarginals  $\tau$  that lie in  $\mathbb{L}(G)$  but *not* in  $\mathbb{M}(G)$ , and therefore could not possibly arise from any global probability distribution.

Let us now demonstrate how, for an appropriately chosen distribution  $p_\theta$  on the graph, the sum-product algorithm can be “fooled” into converging to this pseudomarginal vector  $\tau$ . Using the canonical overcomplete representation (3.34), consider a set of canonical parameters

of the form

$$\theta_s(x_s) := \log \tau_s(x_s) = \log [0.5 \ 0.5] \quad \forall s \in V \quad (4.28a)$$

$$\begin{aligned} \theta_{st}(x_s, x_t) &:= \log \frac{\tau_{st}(x_s, x_t)}{\tau_s(x_s)\tau_t(x_t)} && (4.28b) \\ &= \log 4 \begin{bmatrix} \beta_{st} & 0.5 - \beta_{st} \\ 0.5 - \beta_{st} & \beta_{st} \end{bmatrix} && \forall (s, t) \in E, \end{aligned}$$

where we have adopted the shorthand notation from equation (4.1). With these canonical parameters, suppose that we apply the sum-product algorithm to the Markov random field  $p_\theta$ , using the uniform message initialization  $M_{ts}(x_s) \propto [0.5 \ 0.5]$ . A little bit of algebra using the sum-product update (4.25) shows that for this parameter choice, the uniform messages  $M$  already define a fixed point of the sum-product algorithm. Moreover, if we compute the associated pseudomarginals specified by  $M$  and  $\theta$ , they are equal to the previously specified  $\tau_s, \tau_{st}$ . In summary, the sum-product algorithm—when applied to the distribution  $p_\theta$  defined by the canonical parameters (4.28)—produces as its output the pseudomarginal  $\tau$  as its estimate of the true marginals. (The reader might object to the fact that the problem construction ensured the sum-product algorithm was already at this particular fixed point, and so obviates the possibility of the updates converging to some other fixed point if initialized in a different way. However, it is known [244, 107] that for any discrete Markov random field in exponential family form with at most a single cycle, the sum-product has a unique fixed point, and always converges to it. Therefore, the sum-product fixed point that we have constructed (4.28) is the unique fixed point for this problem.) ♣

More generally, the constructive approach illustrated in the preceding example applies to an arbitrary member of the interior<sup>3</sup> of  $\mathbb{L}(G)$ . Therefore, for all pseudomarginals  $\tau \in \mathbb{L}(G)$ , including those that are not globally valid, there exists a distribution  $p_\theta$  for which  $\tau$  arises from a sum-product fixed point.

<sup>3</sup>Strictly speaking, it applies to members of the relative interior since, as described in the overcomplete representation (3.34), the set  $\mathbb{L}(G)$  is not full-dimensional and hence has an empty interior.

#### 4.1.6 Bethe and loop series expansions

In this section, we discuss the loop series expansions of Chertkov and Chernyak [48]. These expansions provide exact representation of the cumulant function as a sum of terms, with the first term corresponding to the Bethe approximation  $A_{Bethe}(\theta)$ , and higher-order terms obtained by adding in so-called loop corrections. They provided two derivations of their loop series: one applies a trigonometric identity to a Fourier representation of binary variables, while the second is based upon a saddle point approximation obtained via an auxiliary field of complex variables. In this section, we describe a more direct derivation of the loop expansion based on the reparameterization characterization of sum-product fixed points given in Proposition 5. Although the loop series expansion can be developed for general factor graphs, considering the case of a pairwise Markov random field with binary variables—that is, the Ising model from Example 3—suffices to illustrate the basic ideas. (See Sudderth et al. [219] for the derivation for more general factor graphs.)

Before stating the result, we require a few preliminary definitions. Given an undirected graph  $G = (V, E)$  and some subset  $F \subseteq E$  of the edge set  $E$ , we let  $G(F)$  denote the induced subgraph associated with  $F$ —that is, the graph with edge set  $F$ , and vertex set

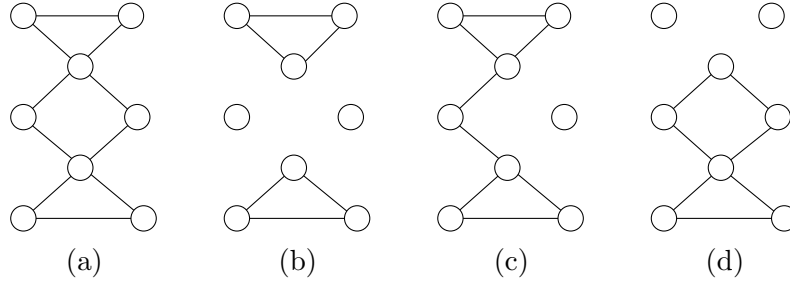
$$V(F) = \{t \in V \mid (t, u) \in F \text{ for some } u\}.$$

For any vertex  $s \in V$ , we define its degree with respect to  $F$  as

$$d_s(F) = |\{t \in V \mid (s, t) \in F\}|. \quad (4.29)$$

Following Chertkov and Chernyak [48], we define a *generalized loop* to be a subgraph  $G(F)$  for which all nodes  $s \in V$  have degree  $d_s(F) \neq 1$ . Otherwise stated, for every node  $s$ , it either does not belong to  $G(F)$  so that  $d_s(F) = 0$ , or it has degree  $d_s(F) \geq 2$ . See Figure 4.3 for an illustration of the concept of a generalized loop. Note also that a graph without cycles (i.e., a tree or forest graph) does not have any generalized loops.

Consider a BP fixed point for a pairwise MRF with binary variables—that is, an Ising model (3.8). For binary variables



**Fig. 4.3.** Illustration of generalized loops. (a) Original graph. (b)–(d) Various generalized loops associated with the graph in (a). In this particular case, the original graph is a generalized loop for itself.

$X_s \in \{0, 1\}$ , the singleton and edgewise pseudomarginals associated with a BP fixed point can be parameterized as

$$\tau_s(x_s) = \begin{bmatrix} 1 - \tau_s \\ \tau_s \end{bmatrix}, \quad \text{and} \quad \tau_{st}(x_s, x_t) = \begin{bmatrix} 1 - \tau_s - \tau_t + \tau_{st} & \tau_t - \tau_{st} \\ \tau_s - \tau_{st} & \tau_{st} \end{bmatrix} \quad (4.30)$$

Some calculation shows that membership in the set  $\mathbb{L}(G)$  is equivalent to imposing the following four inequalities

$$\tau_{st} \geq 0, \quad 1 - \tau_s - \tau_t + \tau_{st} \geq 0, \quad \tau_s - \tau_{st} \geq 0, \quad \text{and} \quad \tau_t - \tau_{st} \geq 0,$$

for each edge  $(s, t) \in E$ . See also Example 10 for some discussion of the mean parameters for an Ising model.

Using this parameterization, for each edge  $(s, t) \in E$ , we define the *edge weight*

$$\beta_{st} := \frac{\tau_{st} - \tau_s \tau_t}{\tau_s(1 - \tau_s)\tau_t(1 - \tau_t)}, \quad (4.31)$$

which extends naturally to the *subgraph weight*  $\beta_F := \prod_{(s,t) \in F} \beta_{st}$ . With these definitions, we have [219]:

**Proposition 6.** *Consider a pairwise Markov random field with binary variables (3.8), and let  $A_{\text{Bethe}}(\theta)$  be the optimized free energy (4.16) evaluated at a BP fixed point  $\tau = \{\tau_s, s \in V\} \cup \{\tau_{st}, (s, t) \in E\}$ . The cumulant function  $A(\theta)$  is equal to the loop series expansion:*

$$A(\theta) = A_{\text{Bethe}}(\theta) + \log \left\{ 1 + \sum_{\emptyset \neq F \subseteq E} \beta_F \prod_{s \in V} \mathbb{E}_{\tau_s} [(X_s - \tau_s)^{d_s(F)}] \right\}. \quad (4.32)$$

Before proving Proposition 6, we pause to make some remarks. By definition, we have

$$\begin{aligned}\mathbb{E}_{\tau_s}[(X_s - \tau_s)^d] &= (1 - \tau_s)(-\tau_s)^d + \tau_s(1 - \tau_s)^d \\ &= \tau_s(1 - \tau_s)[(1 - \tau_s)^{d-1} + (-1)^d(\tau_s)^{d-1}]\end{aligned}\quad (4.33)$$

corresponding to  $d^{\text{th}}$  central moments of a Bernoulli variable with parameter  $\tau_s \in [0, 1]$ . Consequently, for any  $F \subseteq E$  such that  $d_s(F) = 1$  for at least one  $s \in V$ , then the associated term in the expansion (4.32) vanishes. For this reason, only generalized loops  $F$  lead to non-zero terms in the expansion (4.32). The terms associated with these generalized loops effectively define corrections to the Bethe estimate  $A_{\text{Bethe}}(\theta)$  of the cumulant function. Tree-structured graphs do not contain any non-trivial generalized loops, which provides an alternative proof of the exactness of the Bethe approximation for trees.

The first term  $A_{\text{Bethe}}(\theta)$  in the loop expansion is easily computed from any BP fixed point, since it simply corresponds to the optimized value of the Bethe free energy (4.16). However, explicit computation of the full sequence of loop corrections—and hence exact calculation of the partition function—is intractable for general (non-tree) models. For instance, any fully connected graph with  $n \geq 5$  nodes has more than  $2^n$  generalized loops. In some cases, accounting for a small set of significant loop corrections may lead to improved approximations to the partition function [97], or more accurate approximations of the marginals for LDPC codes [49].

*Proof of Proposition 6:*

We begin by noting that the loop series expansion (4.32) is invariant to any reparameterization of  $\theta$ , since such a reparameterization simply shifts both  $A(\theta)$  and  $A_{\text{Bethe}}(\theta)$  by the same constant. Consequently, it suffices to prove the loop series expansion for a single parameterization; in particular, we prove it for the reparameterization given by the canonical parameters

$$\tilde{\theta}_s(x_s) = \log \tau_s(x_s), \quad \text{and} \quad \tilde{\theta}_{st}(x_s, x_t) = \log \frac{\tau_{st}(x_s, x_t)}{\tau_s(x_s)\tau_t(x_t)}, \quad (4.34)$$

as specified by any BP fixed point (see Proposition 5). For this reparam-

eterization, a little calculation shows that  $A_{Bethe}(\tilde{\theta}) = 0$ . Consequently, it suffices to show that for this particular parameterization (4.34), we have the equality

$$A(\tilde{\theta}) = \log \left\{ 1 + \sum_{\emptyset \neq F \subseteq E} \beta_F \prod_{s \in V} \mathbb{E}_{\tau_s} [(X_s - \tau_s)^{d_s(F)}] \right\}. \quad (4.35)$$

Using the representation (4.30), a little calculation shows that

$$\frac{\tau_{st}(x_s, x_t)}{\tau_s(x_s)\tau_t(x_t)} = 1 + \beta_{st}(x_s - \tau_s)(x_t - \tau_t). \quad (4.36)$$

By definition of  $\tilde{\theta}$ , we have

$$\exp(A(\tilde{\theta})) = \sum_x \prod_{s \in V} \tau_s(x_s) \prod_{(s,t) \in E} \frac{\tau_{st}(x_s, x_t)}{\tau_s(x_s)\tau_t(x_t)}.$$

Let  $\mathbb{E}$  denote expectation taken with respect to the product distribution  $\tau_{fact}(x) = \prod_s \tau_s(x_s)$ . With this notation, and applying the identity (4.36), we have

$$\begin{aligned} \exp(A(\tilde{\theta})) &= \sum_{x \in \{0,1\}^n} \prod_{s \in V} \tau_s(x_s) \prod_{(s,t) \in E} \frac{\tau_{st}(x_s, x_t)}{\tau_s(x_s)\tau_t(x_t)} \\ &= \mathbb{E} \left[ \prod_{(s,t) \in E} (1 + \beta_{st}(X_s - \tau_s)(X_t - \tau_t)) \right]. \end{aligned} \quad (4.37)$$

Expanding this polynomial and using linearity of expectation, we recover one term for each non-empty subset  $F \subseteq E$  of the graph's edges:

$$\exp(A(\tilde{\theta})) = 1 + \sum_{\emptyset \neq F \subseteq E} \mathbb{E} \left[ \prod_{(s,t) \in F} \beta_{st}(X_s - \tau_s)(X_t - \tau_t) \right]. \quad (4.38)$$

The expression in equation (4.35) then follows from the independence structure of  $\tau_{fact}(x)$ , and standard formulas for the moments of Bernoulli random variables. To evaluate these terms, note that if  $d_s(F) = 1$ , it follows that  $\mathbb{E}[X_s - \tau_s] = 0$ . There is thus one loop correction for each generalized loop  $F$ , in which all connected nodes have degree at least two. □

Proposition 6 has extensions to more general factor graphs; see the papers [48, 219] for more details. Moreover, Sudderth et al. [219] show that the loop series expansion can be exploited to show that for certain types of graphical models with attractive interactions, the Bethe value  $A_{\text{Bethe}}(\theta)$  is actually a lower bound on the true cumulant function  $A(\theta)$ .

## 4.2 Kikuchi and hypertree-based methods

From our development thus far, we have seen that there are two *distinct* ways in which the Bethe variational principle (4.16) is an approximate version of the exact variational principle (3.45). First, for general graphs, the Bethe entropy (4.14) is only an approximation to the true entropy or negative dual function. Second, the constraint set  $\mathbb{L}(G)$  outer bound on the marginal polytope  $\mathbb{M}(G)$ , as illustrated in Figure 4.2. In principle, the accuracy of the Bethe variational principle could be strengthened by improving either one, or both, of these components. This section is devoted to one natural generalization of the Bethe approximation, first proposed by Yedidia et al. [263, 264] and further explored by various researchers [184, 163, 108, 235, 265], that improves both components simultaneously. The origins underlying ideas lie in the statistical physics literature, where they were referred to as cluster variational methods [6, 224, 130].

At the high level, the approximations in the Bethe approach are based on trees, which represent a special case of the junction trees. A natural strategy, then, is to strengthen the approximations by exploiting more complex junction trees. These approximations are most easily understood in terms of hypertrees, which represent an alternative way in which to describe junction trees. Accordingly, we begin with some necessary background on hypergraphs and hypertrees.

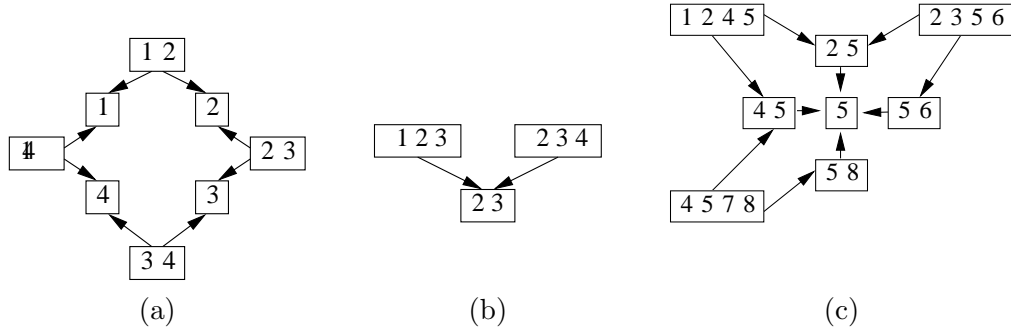
### 4.2.1 Hypergraphs and hypertrees

A hypergraph  $G = (V, E)$  is a generalization of a graph, consisting of a vertex set  $V = \{1, \dots, m\}$ , and a set of hyperedges  $E$ , where each *hyperedge*  $h$  is a particular subset of  $V$ . The hyperedges form a partially-ordered set [216], where the partial ordering is specified by inclusion. Given two distinct hyperedges  $g$  and  $h$ , one of three possibilities can



hold: (a) the hyperedge  $g$  is contained within  $h$ , in which case we write  $g \subset h$ ; (b)  $h$  is contained within  $g$ , and we write  $h \subset g$ ; and (c) neither containment relation holds, in which case  $g$  and  $h$  are incomparable. We say that a hyperedge  $h$  is *maximal* if it is not contained within any other hyperedge. With these definitions, we see that an ordinary graph is a special case of a hypergraph, in which each maximal hyperedge consists of a pair of vertices (i.e., an ordinary edge of the graph).<sup>4</sup>

A convenient graphical representation of a hypergraph is in terms of a diagram of its hyperedges, with directed edges representing the inclusion relations; such a representation is known as a *poset diagram* [163, 184, 216]. Figure 4.4 provides some simple graphical il-



**Fig. 4.4.** Graphical representations of hypergraphs. Subsets of nodes corresponding to hyperedges are shown in rectangles, whereas the arrows represent inclusion relations among hyperedges. (a) An ordinary single cycle graph represented as a hypergraph. (b) A simple hypertree of width two. (c) A more complex hypertree of width three.

lustrations of hypergraphs. Any ordinary graph, as a special case of a hypergraph, can be drawn in terms of a poset diagram; in particular, panel (a) shows the hypergraph representation of a single cycle on four nodes. Panel (b) shows a hypergraph that is not equivalent to an ordinary graph, consisting of two hyperedges of size three joined by their intersection of size two. Shown in panel (c) is a more complex hypertree, to which we will return in the sequel.

<sup>4</sup>There is a minor inconsistency in our definition of the hypertree edge set  $E$ ; for hypergraphs (unlike graphs), the set of hyperedges can include (a subset of the) individual vertices.

Given any hyperedge  $h$ , we define the sets of its *descendants* and *ancestors* in the following way:

$$\mathcal{D}(h) := \{g \in E \mid g \subset h\}, \quad \mathcal{A}(h) := \{g \in E \mid g \supset h\}. \quad (4.39)$$

For example, given the hyperedge  $h = (1245)$  in the hypergraph in Figure 4.4(c), we have  $\mathcal{A}(h) = \emptyset$  and  $\mathcal{D}(h) = \{(25), (45), (5)\}$ . We use the notation  $\mathcal{D}^+(h)$  and  $\mathcal{A}^+(h)$  as shorthand for the sets  $\mathcal{D}(h) \cup h$  and  $\mathcal{A}(h) \cup h$  respectively.

Hypertrees or acyclic hypergraphs provide an alternative way to describe the concept of junction trees, as originally described in Section 2.5.2. In particular, a hypergraph is *acyclic* if it is possible to specify a junction tree using its maximal hyperedges and their intersections. The *width* of an acyclic hypergraph is the size of the largest hyperedge minus one; we use the term *k-hypertree* to mean an acyclic hypergraph of width  $k$  consisting of a single connected component. Thus, for example, a spanning tree of an ordinary graph is a 1-hypertree, because its maximal hyperedges (i.e., ordinary edges) all have size two. As a second example, consider the hypergraph shown in Figure 4.4(c). It is clear that this hypergraph is equivalent to the junction tree with maximal cliques  $\{(1245), (4578), (2356)\}$  and separator sets  $\{(25), (45)\}$ . Since the maximal hyperedges have size four, this hypergraph is a hypertree of width three.

With this background, we now specify an alternative form of the junction tree factorization (2.10), and show how it leads to a local decomposition of the entropy. Associated with any poset is a Möbius function  $\omega : E \times E \rightarrow \mathbb{Z}$ ; see Stanley [216] and Appendix E.1 for more details. We use the Möbius function to define a bijection between the collection of marginals  $\mu := \{\mu_h\}$  associated with the hyperedges of a hypergraph, and a new set of functions  $\varphi := \{\varphi_h\}$ , as follows:

$$\log \mu_h(x_h) = \sum_{g \in \mathcal{D}^+(h)} \log \varphi_g(x_g), \quad \text{and} \quad (4.40a)$$

$$\log \varphi_h(x_h) = \sum_{g \in \mathcal{D}^+(h)} \omega(g, h) \log \mu_g(x_g). \quad (4.40b)$$

For a hypertree with an edge set containing all intersections between maximal hyperedges, the underlying distribution is guaranteed to fac-

torize as follows:

$$p_\mu(x) = \prod_{h \in E} \varphi_h(x_h; \mu). \quad (4.41)$$

Equation (4.41) is an alternative formulation of the well-known junction tree decomposition (2.10). Let us consider some examples to gain some intuition.

**Example 17 (Hypertree factorization).** (a) First suppose that the hypertree is an ordinary tree, in which case the hyperedge set consists of the union of the vertex set with the (ordinary) edge set. For any vertex  $s$ , we have  $\varphi_s(x_s) = \mu_s(x_s)$ , whereas for any edge  $(s, t)$  we have  $\varphi_{st}(x_s, x_t) = \mu_{st}(x_s, x_t) / [\mu_s(x_s) \mu_t(x_t)]$ . In this special case, equation (4.41) reduces to the tree factorization in equation (4.8).

(b) Now consider the acyclic hypergraph specified by  $\{(1245), (2356), (4578), (25), (45), (56), (58), (5)\}$ , as illustrated in Figure 4.4(c). Omitting explicit dependence on  $x$  for notational simplicity, we first calculate  $\varphi_{1245} = \frac{\mu_{1245}}{\varphi_{25}\varphi_{45}\varphi_5} = \frac{\mu_{1245}}{[\mu_{25}/\mu_5][\mu_{45}/\mu_5]\mu_5}$ , with analogous expressions for  $\varphi_{2356}$  and  $\varphi_{4578}$ . We also have  $\varphi_{25} = \mu_{25}/\mu_5$ , with analogous expressions for the other pairwise terms. Putting the pieces together yields that the density  $p_\mu$  factorizes as

$$\begin{aligned} p_\mu &= \frac{\mu_{1245}}{\frac{\mu_{25}}{\mu_5} \frac{\mu_{45}}{\mu_5} \mu_5} \frac{\mu_{2356}}{\frac{\mu_{25}}{\mu_5} \frac{\mu_{56}}{\mu_5} \mu_5} \frac{\mu_{4578}}{\frac{\mu_{45}}{\mu_5} \frac{\mu_{58}}{\mu_5} \mu_5} \frac{\mu_{25}}{\mu_5} \frac{\mu_{45}}{\mu_5} \frac{\mu_{56}}{\mu_5} \frac{\mu_{58}}{\mu_5} \mu_5 \\ &= \frac{\mu_{1245} \mu_{2356} \mu_{4578}}{\mu_{25} \mu_{45}}, \end{aligned}$$

which agrees with the expression from the junction tree formula (2.10).

♣

An immediate but important consequence of the factorization (4.41) is a local decomposition of the entropy:

$$H_{\text{hyper}}(\mu) \stackrel{(a)}{=} - \sum_{h \in E} I_h(\mu_h) \stackrel{(b)}{=} \sum_{h \in E} c(h) H_h(\mu_h), \quad (4.42)$$

where the quantities

$$I_h(\mu_h) := \sum_{x_h} \mu_h(x_h) \log \varphi_h(x_h), \quad \text{and} \quad (4.43a)$$

$$H_h(\mu_h) := - \sum_{x_h} \mu_h(x_h) \log \mu_h(x_h) \quad (4.43b)$$

are, respectively, the multi-information and entropy associated with the hyperedge  $h$ , and the quantities

$$c(f) := \sum_{e \in \mathcal{A}^+(f)} \omega(f, e) \quad (4.44)$$

are known as *overcounting numbers*. Equality (a) in the hypertree entropy (4.42) follows immediately from the hypertree factorization (4.41) and the definition of  $I_h$ . Equality (b) follows by applying the Möbius inversion relation (4.40) between  $\log \varphi_h(x)$  and  $\log \mu_h(x_h)$ , expanding, and simplifying. We illustrate the decomposition (4.42) by continuing with Example 17:

**Example 18 (Hypertree entropies).** (a) For an ordinary tree, there are two types of multi-information: for an edge  $(s, t)$ ,  $I_{st}$  is equivalent to the ordinary mutual information, whereas for any vertex  $s \in V$ , the term  $I_s$  is equal to the negative entropy  $-H_s$ . Consequently, in this special case, equation (4.42) is equivalent to the tree entropy given in equation (4.11). The overcounting numbers for a tree are  $c((s, t)) = 1$  for any edge  $(s, t)$ , and  $c(s) = 1 - d(s)$  for any vertex  $s$ , where  $d(s)$  denotes the number of neighbors of  $s$ .

(b) Consider again the hypertree in Figure 4.4(c). On the basis of our previous calculations in Example 17(c), we calculate  $I_{1245} = -[H_{1245} - H_{25} - H_{45} + H_5]$ . The expressions for the other two maximal hyperedges (i.e.,  $I_{2356}$  and  $I_{4578}$ ) are analogous. Similarly, we can compute  $I_{25} = H_5 - H_{25}$ , with analogous expressions for the other hyperedges of size two. Finally, we have  $I_5 = -H_5$ . Putting the pieces together and doing some algebra yields  $H_{\text{hypertree}} = H_{1245} + H_{2356} + H_{4578} - H_{25} - H_{45}$ .



### 4.2.2 Kikuchi and related approximations

Recall that the core of the Bethe approach of Section 4 consists of a particular tree-based (Bethe) approximation to entropy, and a tree-based outer bound on the marginal polytope. The Kikuchi method and related approximations extend these tree-based approximations to ones based on more general hypertrees, as we now describe. Consider a Markov random field (MRF) defined by some (non-acyclic) hypergraph  $G = (V, E)$ , giving rise to an exponential family distribution  $p_\theta$  of the form

$$p_\theta(x) \propto \exp \left\{ \sum_{h \in E} \theta_h(x_h) \right\}. \quad (4.45)$$

Note that this equation reduces to our earlier representation (4.1) of a pairwise MRF when the hypergraph is an ordinary graph.

Let  $\tau = \{\tau_h\}$  be a collection of local marginals associated with the hyperedges  $h \in E$ . These marginals must satisfy the obvious *normalization condition*  $\sum_{x'_h} \tau_h(x'_h) = 1$ . Similarly, these local marginals must be consistent with one another wherever they overlap; more precisely, for any pair of hyperedges  $g \subset h$ , the *marginalization condition*

$$\sum_{\{x'_h \mid x'_g = x_g\}} \tau_h(x'_h) = \tau_g(x_g)$$

must hold. Imposing these normalization and marginalization conditions leads to the following constraint set:

$$\mathbb{L}_t(G) = \left\{ \tau \geq 0 \mid \sum_{x'_h} \tau_h(x'_h) = 1 \quad \forall h \in E, \right. \\ \left. \sum_{\{x'_h \mid x'_g = x_g\}} \tau_h(x'_h) = \tau_g(x_g) \quad \forall g \subset h \right\}. \quad (4.46)$$

Note that this constraint set is a natural generalization of the tree-based constraint set defined in equation (4.7). In particular, definition (4.46) coincides with definition (4.7) when the hypergraph  $G$  is simply an ordinary graph. As before, we refer to members  $\mathbb{L}_t(G)$  as *pseudomarginals*. By the junction tree conditions in Proposition 1, the local constraints

defining  $\mathbb{L}_t(G)$  are sufficient to guarantee global consistency whenever  $G$  is a hypertree.

In analogy to the Bethe entropy approximation, the entropy decomposition (4.42) motivates the following hypertree-based approximation to the entropy:

$$H_{\text{app}}(\tau) = \sum_{g \in E} c(g) H_g(\tau_g), \quad (4.47)$$

where  $H_g$  is the hyperedge-based entropy (4.43a), and  $c(g) = \sum_{f \in \mathcal{A}^+(g)} \omega(g, f)$  is the overcounting number defined in equation (4.44). This entropy approximation and the outer bound  $\mathbb{L}_t(G)$  on the marginal polytope, in conjunction, lead to the following hypertree-based approximation to the exact variational principle:

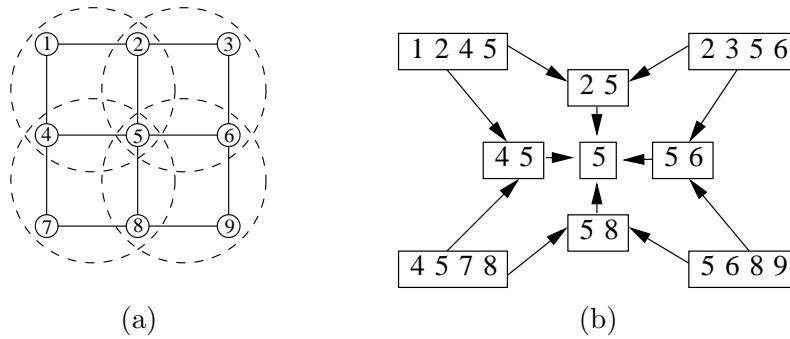
$$\max_{\tau \in \mathbb{L}_t(G)} \left\{ \langle \theta, \tau \rangle + H_{\text{app}}(\tau) \right\}. \quad (4.48)$$

This problem is the hypertree-based generalization of the Bethe variational problem (4.16).

**Example 19 (Kikuchi approximation).** To illustrate the approximate variational principle (4.48), consider the hypergraph shown in Figure 4.5(b). This hypergraph arises from applying the Kikuchi clustering method [264] to the  $3 \times 3$  lattice in panel (a); see Appendix D for more details. We determine the form of the entropy approximation  $H_{\text{app}}$  for this hypergraph by first calculating the overcounting numbers  $c(\cdot)$ . By definition,  $c(h) = 1$  for each of the four maximal hyperedges (e.g.,  $h = (1245)$ ). Since each of the 2-hyperedges has two parents, a little calculation shows that  $c(g) = -1$  for each 2-hyperedge  $g$ . A final calculation shows that  $c(5) = 1$ , so that the overall entropy approximation takes the form

$$H_{\text{app}} = [H_{1245} + H_{2356} + H_{4578} + H_{5689}] - [H_{25} + H_{45} + H_{56} + H_{58}] + H_5. \quad (4.49)$$

Since the hypergraph in panel (b) has treewidth 3, the appropriate constraint set is the polytope  $\mathbb{L}_3(G)$ , defined over the pseudomarginals  $\{\tau_h, h \in E\}$ . In particular, it imposes non-negativity constraints, nor-



**Fig. 4.5.** (a) Kikuchi clusters superimposed upon a  $3 \times 3$  lattice graph. (b) Hypergraph defined by this Kikuchi clustering.

marginalization constraints, and marginalization constraints of the form

$$\sum_{x'_1, x'_2} \tau_{1245}(x'_1, x'_2, x_4, x_5) = \tau_{45}(x_4, x_5), \quad \text{and} \quad \sum_{x'_6} \tau_{56}(x_5, x'_6) = \tau_5(x_5).$$

♣

### 4.2.3 Generalized belief propagation

In principle, the variational problem (4.48) could be solved by a number of methods. Here we describe an algorithm, referred to as *parent-to-child message-passing* by Yedidia et al. [264], that is a natural generalization of the ordinary sum-product updates for the Bethe approximation. As indicated by its name, the defining feature of this scheme is that the only messages passed are from parents to children—i.e., along directed edges in the poset representation of a hypergraph.

In the hypertree-based variational problem (4.48), the variables correspond to a pseudomarginal  $\tau_h$  for each hyperedge  $E \in E'$ . As with the earlier derivation of the sum-product algorithm, a Lagrangian formulation of this optimization problem leads to a specification of the optimizing pseudomarginals in terms of messages, which represent Lagrange multipliers associated with the constraints. There are various Lagrangian reformulations of the original problem [e.g., 263, 264, 163], which lead to different message-passing algorithms. Here we describe

the parent-to-child form of message-passing derived by Yedidia et al. [264]. Given a pair of hyperedges  $(f, g)$ , we let  $M_{f \rightarrow g}(x_g)$  denote the “message” passed from hyperedge  $f$  to hyperedge  $g$ . More precisely, this message is a function over the state space of  $x_g$  (e.g., a multi-dimensional array in the case of discrete random variables). In terms of these messages, the pseudomarginal  $\tau_h$  in parent-to-child message-passing takes the following form:

$$\tau_h(x_h) \propto \left[ \prod_{g \in \mathcal{D}^+(h)} \psi_g(x_g; \theta) \right] \left[ \prod_{g \in \mathcal{D}^+(h)} \prod_{f \in \text{Par}(g) \setminus \mathcal{D}^+(h)} M_{f \rightarrow g}(x_g) \right], \quad (4.50)$$

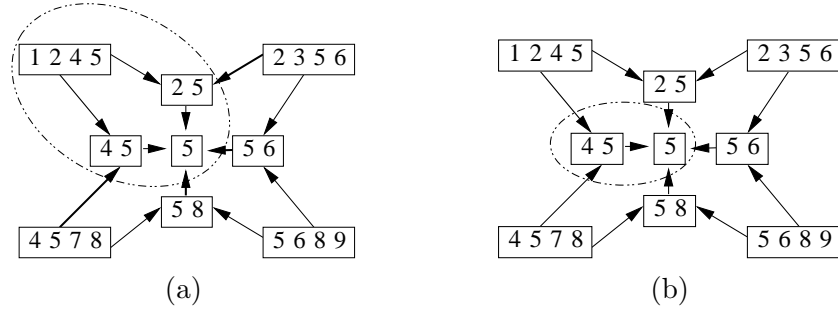
where we have introduced the convenient shorthand  $\psi_g(x_g; \theta) = \exp(\theta(x_g))$ . In this equation, the pseudomarginal  $\tau_h$  includes a compatibility function  $\psi_g$  for each hyperedge  $g$  in the set  $\mathcal{D}^+(h) := \mathcal{D}(h) \cup h$ . It also collects a message from each hyperedge  $f \notin \mathcal{D}^+(h)$  that is a parent of some hyperedge  $g \in \mathcal{D}^+(h)$ . We illustrate this construction by following up on Example 19.

**Example 20 (Parent-to-child for Kikuchi).** In order to illustrate the parent-to-child message-passing, consider the Kikuchi approximation for a  $3 \times 3$  grid, illustrated in Figure 4.5. Focusing first on the hyperedge (1245), the first term in equation (4.50) specifies a product of compatibility functions  $\psi_g$  as  $g$  ranges over  $\mathcal{D}^+(1245)$ , which in this case yields the product  $\psi_{1245}\psi_{25}\psi_{45}\psi_5$ . We then take the product over messages from hyperedges that are parents of hyperedges in  $\mathcal{D}^+\{(1245)\}$ , excluding hyperedges in  $\mathcal{D}^+\{(1245)\}$  itself. Figure 4.6(a) provides an illustration; the set  $\mathcal{D}^+\{(1245)\}$  is given by the hyperedges within the dotted ellipses. In this case, the set  $\cup_g \text{Par}(g) \setminus \mathcal{D}^+(h)$  is given by (2356) and (4578), corresponding to the parents of (25) and (45) respectively, combined with hyperedges (56) and (58), which are both parents of hyperedge (5). The overall result is an expression of the following form:

$$\begin{aligned} \tau_{1245} \propto & \psi'_{12} \psi'_{14} \psi'_{25} \psi'_{45} \psi'_1 \psi'_2 \psi'_4 \psi'_5 \\ & \times M_{(2356) \rightarrow (25)} M_{(4578) \rightarrow (45)} M_{(56) \rightarrow 5} M_{(58) \rightarrow 5}. \end{aligned}$$

By symmetry, the expressions for the pseudomarginals on the other 4-hyperedges are analogous. By similar arguments, it is straightforward





**Fig. 4.6.** Illustration of relevant regions for parent-to-child message-passing in a Kikuchi approximation. (a) Message-passing for hyperedge (1245). Set of descendants  $\mathcal{D}^+\{(1245)\}$  is shown within a dotted ellipse. Relevant parents for  $\tau_{1245}$  consists of the set  $\{(2356), (4578), (56), (58)\}$ . (b) Message-passing for hyperedge (45). Dotted ellipse shows descendant set  $\mathcal{D}^+\{(45)\}$ . In this case, relevant parent hyperedges are  $\{(1245), (4578), (25), (56), (58)\}$ .

to compute the following expression for  $\tau_{45}$  and  $\tau_5$ :

$$\begin{aligned} \tau_{45} &\propto \psi'_{45} \psi'_4 \psi'_5 M_{(1245) \rightarrow (45)} M_{(4578) \rightarrow (45)} M_{(25) \rightarrow 5} M_{(56) \rightarrow 5} M_{(58) \rightarrow 5} \\ \tau_5 &\propto \psi'_5 M_{(45) \rightarrow 5} M_{(25) \rightarrow 5} M_{(56) \rightarrow 5} M_{(58) \rightarrow 5}. \end{aligned}$$



Generalized forms of the sum-product updates follow by updating the messages so as to enforce the marginalization constraints defining membership in  $\mathbb{L}(G)$ ; as in the proof of Theorem 3, fixed points of these updates satisfy the necessary stationary conditions of the Lagrangian formulation. Further details on different variants of generalized sum-product updates can be found in various papers [263, 264, 184, 163, 125].

### 4.3 Expectation-propagation algorithms

There are a variety of other algorithms in the literature that are message-passing algorithms in the spirit of the sum-product algorithm. Examples of such algorithms include the family of expectation-propagation algorithms due to Minka [172], the related class of assumed density filtering methods [149, 161, 37], expectation-consistent inference [181], structured summary-propagation algorithms [61, 112], and

the adaptive TAP method of Opper and Winther [179, 180]. These algorithms are often defined operationally in terms of sequences of local moment-matching updates, with variational principles invoked only to characterize each individual update, not to characterize the overall approximation. In this section we show that these algorithms are in fact variational inference algorithms, involving a particular kind of approximation to the exact variational principle in Theorem 2.

The earliest forms of assumed density filtering [161] were developed for time series applications, in which the underlying graphical model is a hidden Markov model (HMM), as illustrated in Figure 2.4(a). As we have discussed, for discrete random variables, the marginal distributions for an HMM can be computed using the forward-backward algorithm, corresponding to a particular instantiation of the sum-product algorithm. Similarly, for a Gauss-Markov process, the Kalman filter also computes the means and covariances at each node of an HMM. However, given a hidden Markov model involving general continuous random variables, the message  $M_{ts}(\cdot)$  passed from node  $t$  to  $s$  is a real-valued function, and hence difficult to store and transmit.<sup>5</sup> The purpose of assumed density filtering is to circumvent the computational challenges associated with passing function-valued messages. Instead, assumed density filtering (ADF) operates by passing approximate forms of the messages, in which the true message is approximated by the closest member of some tractable class. For instance, a general continuous message might be approximated with a Gaussian message. This procedure of computing the message approximations, if closeness is measured in terms of the Kullback-Leibler divergence, can be expressed in terms moment-matching operations.

Minka [172] observed that the basic ideas underlying ADF can be generalized beyond Markov chains to arbitrary graphical models, an insight that forms the basis for the family of expectation-propagation (EP) algorithms. As with assumed density filtering, expectation-propagation [172, 171] and various related algorithms [181, 61, 112] are typically described in terms of moment-matching operations. To date,

<sup>5</sup>The Gaussian case is special, in that this functional message can always be parameterized in terms of its “mean” and “variance”, and the efficiency of Kalman filtering stems from this fact.

the close link between these algorithms and the Bethe approximation does not appear to have been widely appreciated. In this section, we show that these algorithms are methods for solving certain relaxations of the exact variational principle from Theorem 2, using “Bethe-like” entropy approximations and particular convex outer bounds on the set  $\mathcal{M}$ . More specifically, we recover the moment-matching updates of expectation-propagation as one particular type of Lagrangian method for solving the resulting optimization problem, thereby showing that expectation-propagation algorithms belong to the same class of variational methods as belief propagation. This section is a refinement of results from the thesis [235].

### 4.3.1 Entropy approximations based on term decoupling

We begin by developing a general class of entropy approximations based on decoupling an intractable collection of terms. Given a collection of random variables  $(X_1, \dots, X_m) \in \mathbb{R}^m$ , consider a collection of sufficient statistics that is partitioned as

$$\underbrace{\phi := (\phi_1, \phi_2, \dots, \phi_{d_T})}_{\text{Tractable component}}, \quad \text{and} \quad \underbrace{\Phi := (\Phi^1, \Phi^2, \dots, \Phi^{d_I})}_{\text{Intractable component}}, \quad (4.51)$$

where  $\phi_i$  are univariate statistics and the  $\Phi^i$  are generally multivariate. As will be made explicit in the examples to follow, the partitioning separates the tractable from the intractable components of the associated exponential family distribution.

Let us set up various exponential families associated with sub-collections of  $(\phi, \Phi)$ . First addressing the tractable component, the vector-valued function  $\phi : \mathcal{X}^m \rightarrow \mathbb{R}^{d_T}$  has an associated vector of canonical parameters  $\theta \in \mathbb{R}^{d_T}$ . Turning to the intractable component, for each  $i = 1, \dots, d_I$ , the function  $\Phi^i$  maps from  $\mathcal{X}^m$  to  $\mathbb{R}^b$ , and the vector  $\tilde{\theta}^i \in \mathbb{R}^b$  is the associated set of canonical parameters. Overall, the function  $\Phi = (\Phi^1, \dots, \Phi^{d_I})$  maps from  $\mathcal{X}^m$  to  $\mathbb{R}^{b \times d_I}$ , and has the associated canonical parameter vector  $\tilde{\theta} \in \mathbb{R}^{b \times d_I}$ , partitioned as  $\tilde{\theta} = (\tilde{\theta}^1, \tilde{\theta}^2, \dots, \tilde{\theta}^{d_I})$ . These families of sufficient statistics define the

exponential family

$$\begin{aligned} p(x; \theta, \tilde{\theta}) &\propto f_0(x) \exp(\langle \theta, \phi(x) \rangle) \exp(\langle \tilde{\theta}, \Phi(x) \rangle) \\ &= f_0(x) \exp(\langle \theta, \phi(x) \rangle) \prod_{i=1}^{d_I} \exp(\langle \tilde{\theta}^i, \Phi^i(x) \rangle). \end{aligned} \quad (4.52)$$

We say that any density  $p$  of the form (4.52) belongs to the  $(\phi, \Phi)$ -exponential family.

Next we define the *base model*

$$p(x; \theta, \vec{0}) \propto f_0(x) \exp(\langle \theta, \phi(x) \rangle), \quad (4.53)$$

in which the intractable sufficient statistics  $\Phi$  play no role, since we have set  $\tilde{\theta} = \vec{0}$ . We say that any distribution  $p$  of the form (4.53) belongs to the  $\phi$ -exponential family. Similarly, for each index  $i \in \{1, \dots, d_I\}$ , we define the  $\Phi^i$ -augmented distribution

$$p(x; \theta, \tilde{\theta}^i) \propto f_0(x) \exp(\langle \theta, \phi(x) \rangle) \exp(\langle \tilde{\theta}^i, \Phi^i(x) \rangle), \quad (4.54)$$

in which only a single term  $\Phi^i$  has been introduced, and we say that any  $p$  of the form (4.54) belongs to the  $(\phi, \Phi^i)$ -exponential family.

The basic premises in the tractable-intractable partitioning between  $\phi$  and  $\Phi$  are:

- First, it is possible to compute marginals exactly in polynomial time for distributions of the base form (4.53)—that is, for any member of the  $\phi$ -exponential family.
- Secondly, for each index  $i = 1, \dots, d_I$ , exact polynomial-time computation is also possible for any distribution of the  $\Phi^i$ -augmented form (4.54)—that is, for any member of the  $(\phi, \Phi^i)$ -exponential family.
- Third, it is intractable to perform exact computations in the full  $(\phi, \Phi)$ -exponential family (4.52), since it simultaneously incorporates *all* of the terms  $(\Phi^1, \dots, \Phi^{d_I})$ .

**Example 21 (Tractable/intractable partitioning for mixture models).** Let us illustrate the partitioning scheme (4.52) with the example of a Gaussian mixture model. Suppose that the random vector

$X \in \mathbb{R}^m$  has a multivariate Gaussian distribution,  $N(0, \Sigma)$ . Letting  $\varphi(y; \mu, \Lambda)$  denote the density of a random variable with a  $N(\mu, \Lambda)$  distribution, consider the two-component Gaussian mixture model

$$p(y \mid X = x) = (1 - \alpha) \varphi(y; 0, \sigma_0^2 I) + \alpha \varphi(y; x, \sigma_1^2 I), \quad (4.55)$$

where  $\alpha \in (0, 1)$  is the mixing weight,  $\sigma_0^2$  and  $\sigma_1^2$  are the variances, and  $I$  is the  $m \times m$  identity matrix.

Given  $n$  i.i.d. samples  $y^1, \dots, y^n$  from the mixture density (4.55), it is frequently of interest to compute marginals under the posterior distribution of  $X$  conditioned on  $(y^1, \dots, y^n)$ . Assuming a multivariate Gaussian prior  $X \sim N(0, \Sigma)$ , and using Bayes' theorem, we see that the posterior takes the form

$$\begin{aligned} p(x \mid y^1, \dots, y^n) &\propto \exp\left(-\frac{1}{2}x^T \Sigma^{-1}x\right) \prod_{i=1}^n p(y^i \mid X = x) \\ &= \exp\left(-\frac{1}{2}x^T \Sigma^{-1}x\right) \exp\left\{\sum_{i=1}^n \log p(y^i \mid X = x)\right\}. \end{aligned} \quad (4.56)$$

To cast this model as a special case of the partitioned exponential family (4.52), we first observe that the term  $\exp(-\frac{1}{2}x^T \Sigma^{-1}x)$  can be identified with the base term  $f_0(x) \exp(\langle \theta, \phi(x) \rangle)$ , so that we have  $d_T = m$ . On the other hand, suppose that we define  $\Phi^i(x) := \log p(y^i \mid X = x)$  for each  $i = 1, \dots, n$ . Since the observation  $y^i$  is a fixed quantity, this definition makes sense. With this definition, the term  $\exp\{\sum_{i=1}^n \log p(y^i \mid X = x)\}$  corresponds to the product  $\prod_{i=1}^{d_I} \exp(\langle \tilde{\theta}^i, \Phi^i(x) \rangle)$  in equation (4.52). Note that we have  $d_I = n$  and  $b = 1$ , with  $\tilde{\theta}^i = 1$  for all  $i = 1, \dots, n$ .

As a particular case of the  $(\phi, \Phi)$ -exponential family in this setting, the base distribution  $p(x; \theta, \vec{0}) \propto \exp(-\frac{1}{2}x^T \Sigma^{-1}x)$  corresponds to a multivariate Gaussian, for which exact calculations of marginals is possible in polynomial time (at most cubic time). Similarly, for each  $i = 1, \dots, n$ , the  $\Phi^i$ -augmented distribution (4.54) is proportional to

$$\exp\left(-\frac{1}{2}x^T \Sigma^{-1}x\right) [(1 - \alpha) \varphi(y^i; 0, \sigma_0^2 I) + \alpha \varphi(y^i; x, \sigma_1^2 I)], \quad (4.57)$$

This distribution is a Gaussian mixture with two components, so that it is also possible to compute marginals exactly in cubic time. However,

the full distribution (4.56) is a Gaussian mixture with  $2^n$  components, so that the complexity of computing exact marginals is exponential in the problem size. ♣

Returning to the main thread, let us develop a few more features of the distribution of interest (4.45). Since it is an exponential family, it has a partitioned set of mean parameters  $(\mu, \tilde{\mu}) \in \mathbb{R}^{d_T} \times \mathbb{R}^{d_I \times b}$ , where

$$\mu = \mathbb{E}[\phi(X)], \quad \text{and} \quad (\tilde{\mu}^1, \dots, \tilde{\mu}^{d_I}) = \mathbb{E}[(\Phi^1(X), \dots, \Phi^{d_I}(X))].$$

As an exponential family, the general variational principle from Theorem 2 is applicable. As usual, the relevant quantities in the variational principle (3.45) are the set

$$\mathcal{M}(\phi, \Phi) := \{(\mu, \tilde{\mu}) \mid (\mu, \tilde{\mu}) = \mathbb{E}_p[(\phi(X), \Phi(X))] \text{ for some } p\}, \quad (4.58)$$

and the entropy, or negative dual function  $H(\mu, \tilde{\mu}) = -A^*(\mu, \tilde{\mu})$ . Given our assumption that exact computation under the full distribution (4.45) is intractable, there must be challenges associated with characterizing the mean parameter space (4.58) and/or the entropy function. Accordingly, we now describe a natural approximation to these quantities, based on the partitioned structure of the distribution (4.45), that leads to the class of expectation-propagation algorithms.

Associated with the base distribution (4.53) is the set

$$\mathcal{M}(\phi) := \{\mu \in \mathbb{R}^{d_T} \mid \mu = \mathbb{E}_p[\phi(X)] \text{ for some density } p\}, \quad (4.59)$$

corresponding to the globally realizable mean parameters for the base distribution viewed as a  $d_T$ -dimensional exponential family. By Theorem 1, for any  $\mu \in \mathcal{M}^\circ(\phi)$ , there exists an exponential family member  $p_{\theta(\mu)}$  that realizes it. Moreover, by our assumption that the base distribution is tractable, we can compute the entropy  $H(\mu)$  of this distribution. Similarly, for each  $i = 1, \dots, d_I$ , the  $\Phi^i$ -augmented distribution (4.54) is associated with the mean parameter space

$$\mathcal{M}(\phi, \Phi^i) := \{(\mu, \tilde{\mu}^i) \in \mathbb{R}^{d_T} \times \mathbb{R}^b \mid (\mu, \tilde{\mu}^i) = \mathbb{E}_p[(\phi(X), \Phi^i(X))] \text{ for some density } p\}. \quad (4.60)$$

By similar reasoning as above, for any  $(\mu, \tilde{\mu}^i) \in \mathcal{M}^\circ(\phi, \Phi^i)$ , there is a member of the  $(\phi, \Phi^i)$ -exponential family with these mean parameters.

Furthermore, since we have assumed that the  $\Phi^i$ -distribution (4.54) is tractable, it is easy to determine membership in the set  $\mathcal{M}(\phi; \Phi^i)$ , and moreover the entropy  $H(\mu, \tilde{\mu}^i)$  can be computed easily.

With these basic ingredients, we can now define an outer bound on the set  $\mathcal{M}(\phi, \Phi)$ . Given a candidate set of mean parameters  $(\tau, \tilde{\tau}) \in \mathbb{R}^{d_T} \times \mathbb{R}^{d_I \times b}$ , we define for each  $i = 1, 2, \dots, d_I$  the coordinate projection operator  $\Pi^i : \mathbb{R}^{d_T} \times \mathbb{R}^{d_I \times b} \rightarrow \mathbb{R}^{d_T} \times \mathbb{R}^b$  that operates as

$$(\tau, \tilde{\tau}) \xrightarrow{\Pi^i} (\tau, \tilde{\tau}^i) \in \mathbb{R}^{d_T} \times \mathbb{R}^b.$$

We then define the set

$$\mathcal{L}(\phi; \Phi) := \{(\tau, \tilde{\tau}) \mid \tau \in \mathcal{M}(\phi), \Pi^i(\tau, \tilde{\tau}) \in \mathcal{M}(\phi, \Phi^i) \quad \forall i = 1, \dots, d_I\}. \quad (4.61)$$

Note that  $\mathcal{L}(\phi; \Phi)$  is a convex set, and moreover it is an outer bound on the true set  $\mathcal{M}(\phi; \Phi)$  of mean parameters.

We now develop an entropy approximation that is tailored to the structure of  $\mathcal{L}(\phi; \Phi)$ . We begin by observing that for any  $(\tau, \tilde{\tau}) \in \mathcal{L}(\phi; \Phi)$  and for each  $i = 1, \dots, d_I$ , there is a member of the  $(\phi, \Phi^i)$ -exponential family with mean parameters  $(\tau, \tilde{\tau}^i)$ . This assertion follows because the projected mean parameters  $(\tau, \tilde{\tau}^i)$  belong to  $\mathcal{M}(\phi; \Phi^i)$ , so that Theorem 1 can be applied. We let  $H(\tau, \tilde{\tau}^i)$  denote the entropy of this exponential family member. Similarly, since  $\tau \in \mathcal{M}(\phi)$  by definition of  $\mathcal{L}(\phi; \Phi)$ , there is a member of the  $\phi$ -exponential family with mean parameter  $\tau$ ; we denote its entropy by  $H(\tau)$ . With these ingredients, we define the following *term-by-term entropy approximation*

$$H_{\text{ep}}(\tau, \tilde{\tau}) := H(\tau) + \sum_{\ell=1}^{d_I} [H(\tau, \tilde{\tau}^\ell) - H(\tau)]. \quad (4.62)$$

Combining this entropy approximation with the convex outer bound (4.61) yields the optimization problem

$$\max_{(\tau, \tilde{\tau}) \in \mathcal{L}(\phi; \Phi)} \{ \langle \tau, \theta \rangle + \langle \tilde{\tau}, \tilde{\theta} \rangle + H_{\text{ep}}(\tau, \tilde{\tau}) \}. \quad (4.63)$$

This optimization problem is an approximation to the exact variational principle from Theorem 2 for the  $(\phi, \Phi)$ -exponential family (4.45), and it underlies the family of expectation-propagation algorithms. It is

closely related to the Bethe variational principle, as the examples to follow should clarify.

**Example 22 (Sum-product and Bethe approximation).** To provide some intuition, we begin by rederiving the Bethe approximation from the point of view of (4.63). That is, we derive Bethe entropy approximation (4.14) as a particular case of the term-by-term entropy approximation (4.62), and the tree-based outer bound  $L(G)$  from equation (4.7) as a particular case of the convex outer bound (4.61) on  $\mathcal{M}(\phi; \Phi)$ .

Consider a pairwise Markov random field based on an undirected graph  $G = (V, E)$ , involving a discrete variable  $X_s \in \{0, 1, \dots, r-1\}$  at each vertex  $s \in V$ ; as we have seen, this can be expressed as an exponential family in the form (4.1) using the standard overcomplete parameterization (3.34) with indicator functions. Taking the point of view of (4.63), we partition the sufficient statistics as follows: the sufficient statistics associated with nodes are defined to be the tractable set and those associated with edges are defined to be the intractable set). Thus, using the functional shorthands  $\theta_s(\cdot)$  and  $\theta_{st}(\cdot, \cdot)$  defined in equation (4.2), the base distribution (4.53) takes the form

$$p(x; \theta_1, \dots, \theta_m, \vec{0}) \propto \prod_{s \in V} \exp(\theta_s(x_s)). \quad (4.64)$$

In this particular case, the terms  $\Phi^i$  to be added correspond to the functions  $\theta_{st}(\cdot, \cdot)$ —that is, the index  $i$  runs over the edge set of the graph. For edge  $(u, v)$ , the  $\Phi^{uv}$ -augmented distribution (4.54) takes the form

$$p(x; \theta_1, \dots, \theta_m, \theta_{uv}) \propto \left[ \prod_{s \in V} \exp(\theta_s(x_s)) \right] \exp(\theta_{uv}(x_u, x_v)). \quad (4.65)$$

The mean parameters associated with the standard overcomplete representation are singleton and pairwise marginal distributions, which we denote by  $\tau_s(\cdot)$  and  $\tau_{st}(\cdot, \cdot)$  respectively. The entropy of the base distribution depends only on the singleton marginals, and given the product structure (4.64), takes the simple form

$$H(\tau_1, \dots, \tau_m) = \sum_{s \in V} H(\tau_s)$$



where  $H(\tau_s) = -\sum_{x_s} \tau_s(x_s) \log \tau_s(x_s)$  is the entropy of the marginal distribution. Similarly, since the augmented distribution (4.65) has only a single edge added (and factorizes over a cycle-free graph), its entropy has the explicit form

$$\begin{aligned} H(\tau_1, \dots, \tau_m, \tau_{uv}) &= \sum_{s \in V} H(\tau_s) + [H(\tau_{uv}) - H(\tau_u) - H(\tau_v)] \\ &= \sum_{s \in V} H(\tau_s) - I(\tau_{uv}), \end{aligned}$$

where  $H(\tau_{uv}) := -\sum_{x_u, x_v} \tau_{uv}(x_u, x_v) \log \tau_{uv}(x_u, x_v)$  is the joint entropy, and  $I(\tau_{uv}) := H(\tau_u) + H(\tau_v) - H(\tau_{uv})$  is the mutual information. Putting together the pieces, we find that the term-by-term entropy approximation (4.62)—for this particular problem—has the form

$$H_{\text{ep}}(\tau) = \sum_{s \in V} H(\tau_s) - \sum_{(s,t) \in E} I_{st}(\tau_{st}),$$

which is precisely the Bethe entropy approximation defined previously (4.14).

Next we show how the outer bound  $\mathcal{L}(\phi; \Phi)$  defined in equation (4.61) specializes to  $\mathbb{L}(G)$ . Consider a candidate set of local marginal distributions

$$\tau = (\{\tau_s, s \in V\}, \{\tau_{st}, (s, t) \in E\}).$$

In the current setting, the set  $\mathcal{M}(\phi)$  corresponds to the set of all globally realizable marginals  $\{\tau_s, s \in V\}$  under a factorized distribution, so that the inclusion  $\{\tau_s, s \in V\} \in \mathcal{M}(\phi)$  is equivalent to the non-negativity constraints  $\tau_s(x_s) \geq 0$ , and the local normalization constraints

$$\sum_{x_s} \tau_s(x_s) = 1 \quad \text{for all } s \in V. \quad (4.66)$$

Recall that the index  $i$  runs over edges of the graph; if, for instance  $\ell = (u, v)$ , then the projected marginals  $\Pi^{uv}(\tau)$  are given by  $(\tau_1, \dots, \tau_m, \tau_{uv})$ . The set  $\mathcal{M}(\phi; \Phi^{uv})$  is traced out by all globally consistent marginals of this form, and is equivalent to the marginal polytope  $\mathbb{M}(G_{uv})$ , where  $G_{uv}$  denotes the graph with a single edge  $(u, v)$ . Since this graph is a tree, the inclusion  $\Pi^{uv}(\tau) \in \mathcal{M}(\phi; \Phi^{uv})$  is equivalent

to having  $\Pi^{uv}(\tau)$  satisfy—in addition to the non-negativity and local normalization constraints (4.66)—the marginalization conditions

$$\sum_{x_v} \tau_{uv}(x_u, x_v) = \tau_u(x_u), \quad \text{and} \quad \sum_{x_u} \tau_{uv}(x_u, x_v) = \tau_v(x_v).$$

Therefore, the full collection of inclusions  $\Pi^{uv}(\tau) \in \mathcal{M}(\phi; \Phi^{uv})$ , as  $(u, v)$  runs over the graph edge set  $E$ , specify the same conditions defining the first-order relaxed constraint set  $\mathbb{L}(G)$  from equation (4.7). ♣

### 4.3.2 Optimality in terms of moment-matching

Returning to the main thread, we now consider a natural Lagrangian method for attempting to solve the expectation-propagation variational principle (4.63). As we will show, these Lagrangian updates reduce to moment-matching, so that the usual expectation-propagation updates are recovered.

Our Lagrangian formulation is based on the following two steps:

- First we augment the space of pseudo-mean-parameters over which we optimize, so that the original constraints defining  $\mathcal{L}(\phi; \Phi)$  are decoupled.
- Second, we add new constraints—to be penalized with Lagrange multipliers—so as to enforce the remaining constraints required for membership in  $\mathcal{L}(\phi; \Phi)$ .

Beginning the augmentation step, let us duplicate the vector  $\tau$  a total of  $d_I$  times, defining thereby defining  $d_I$  new vectors  $\eta^i$  and imposing the constraint that  $\eta^i = \tau$  for each index  $i = 1, \dots, d_I$ . This yields a large collection of pseudo-mean-parameters

$$\{\tau, (\eta^i, \tilde{\tau}^i), i = 1, \dots, d_I\} \in \mathbb{R}^{d_T} \times (\mathbb{R}^{d_T} \times \mathbb{R}^b)^{d_I},$$

and we recast the variational principle (4.63) using these pseudo-mean-parameters as follows:

$$\max_{\{\tau, (\eta^i, \tilde{\tau}^i)\}} \left\{ \langle \tau, \theta \rangle + \sum_{i=1}^{d_I} \langle \tilde{\tau}^i, \tilde{\theta}^i \rangle + \underbrace{H(\tau) + \sum_{i=1}^{d_I} [H(\eta^i, \tilde{\tau}^i) - H(\eta^i)]}_{F(\tau; (\eta^i, \tilde{\tau}^i))} \right\}, \quad (4.67)$$

subject to the constraints  $(\eta^i, \tilde{\tau}^i) \in \mathcal{M}(\phi; \Phi^i)$  for all  $i = 1, \dots, d_I$ , and

$$\tau = \eta^i \quad \forall i = 1, \dots, d_I. \quad (4.68)$$

Let us now consider a particular iterative scheme for solving the reformulated problem (4.67). In particular, for each  $i = 1, \dots, d_I$ , define a vector of Lagrange multipliers  $\lambda^i \in \mathbb{R}^{d_T}$  associated with the constraints  $\tau = \eta^i$ . We then form the Lagrangian function

$$L(\tau; \lambda) = \langle \tau, \theta \rangle + \sum_{i=1}^{d_I} \langle \tilde{\tau}^i, \tilde{\theta}^i \rangle + F(\tau; (\eta^i, \tilde{\tau}^i)) + \sum_{i=1}^{d_I} \langle \lambda^i, \tau - \eta^i \rangle. \quad (4.69)$$

This Lagrangian is a partial one, since we are enforcing the constraints  $\tau \in \mathcal{M}(\phi)$  and  $(\eta^i, \tilde{\tau}^i) \in \mathcal{M}(\phi; \Phi^i)$  explicitly.

Consider an optimum solution  $\{\tau, (\eta^i, \tilde{\tau}^i), i = 1, \dots, d_I\}$  of the optimization problem (4.67) that satisfies the following properties: (a) the vector  $\tau$  belongs to the (relative) interior  $\mathcal{M}^\circ(\phi)$ , and (b) for each  $i = 1, \dots, d_I$ , the vector  $(\eta^i, \tilde{\tau}^i)$  belongs to the relative interior  $\mathcal{M}^\circ(\phi; \Phi^i)$ . Any such solution must satisfy the zero-gradient conditions associated with the partial Lagrangian (4.69)—namely

$$\nabla_\tau L(\tau; \lambda) = 0, \quad (4.70a)$$

$$\nabla_{(\eta^i, \tilde{\tau}^i)} L(\tau; \lambda) = 0 \quad \text{for } i = 1, \dots, d_I, \text{ and} \quad (4.70b)$$

$$\nabla_\lambda L(\tau; \lambda) = 0. \quad (4.70c)$$

Since the vector  $\tau$  belongs to  $\mathcal{M}^\circ(\phi)$ , it specifies a distribution in the  $\phi$ -exponential family. By explicitly computing the Lagrangian condition (4.70a) and performing some algebra—essentially the same steps as the proof of Theorem 3—we find that this exponential family member can be written, in terms of the original parameter vector  $\theta$  and the Lagrange multipliers  $\lambda$ , as follows:

$$q(x; \theta, \lambda) \propto f_0(x) \exp \left( \left\langle \theta + \sum_{i=1}^{d_I} \lambda^i, \phi(x) \right\rangle \right). \quad (4.71)$$

Similarly, since for each  $i = 1, \dots, d_I$ , the vector  $(\tau, \tilde{\tau}^i)$  belongs to  $\mathcal{M}^\circ(\phi; \Phi^i)$ , it also specifies a distribution in the  $(\phi, \Phi^i)$ -exponential

family. Explicitly computing the condition (4.70b) and performing some algebra shows that this distribution can be expressed as

$$q^i(x; \theta, \tilde{\theta}^i, \lambda) \propto f_0(x) \exp(\langle \theta + \sum_{\ell \neq i} \lambda^\ell, \phi(x) \rangle + \langle \tilde{\theta}^i, \Phi^i(x) \rangle), \quad (4.72)$$

The final Lagrangian condition (4.70c) ensures that the constraints (4.68) are satisfied. Noting that  $\tau = \mathbb{E}_q[\phi(X)]$  and  $\eta^i = \mathbb{E}_{q^i}[\phi(X)]$ , these constraints reduce to the moment-matching conditions

$$\int q(x; \theta, \lambda) \phi(x) \nu(dx) = \int q^i(x; \theta, \tilde{\theta}^i, \lambda) \phi(x) \nu(dx), \quad (4.73)$$

for  $i = 1, \dots, d_I$ .

On the basis of equations (4.71), (4.72) and (4.73), we arrive at the expectation-propagation updates, as summarized in Box 4.7. We note

**Expectation-propagation (EP) updates:**

- (1) At iteration  $n = 0$ , initialize the Lagrange multiplier vectors  $(\lambda^1, \dots, \lambda^{d_I})$ .
- (2) At each iteration,  $n = 1, 2, \dots$ , choose some index  $\ell(n) \in \{1, \dots, d_I\}$ , and
  - (a) Using equation (4.72), form the distribution  $q^{\ell(n)}$  and compute the mean parameter

$$\eta^{\ell(n)} := \int q^{\ell(n)}(x) \phi(x) \nu(dx) = \mathbb{E}_{q^{\ell(n)}}[\phi(X)]. \quad (4.74)$$

- (b) Using equation (4.71), form the distribution  $q$  and adjust  $\lambda^{\ell(n)}$  to satisfy the moment-matching condition

$$\mathbb{E}_q[\phi(X)] = \eta^{\ell(n)}. \quad (4.75)$$

**Fig. 4.7.** Steps involved in the expectation-propagation updates. It is a Lagrangian algorithm for attempting to solve the Bethe-like approximation (4.67), or equivalently (4.63).

that the algorithm is well defined, since for any realizable mean parameter  $\tau$ , there is always a solution for  $\lambda^{\ell(n)}$  in equation (4.75). This fact

can be seen by considering the  $d_T$ -dimensional exponential family defined by the pair  $(\lambda^{\ell(n)}, \phi)$ , and then applying Theorem 1. Moreover, it follows that any fixed point of these EP updates satisfies the necessary Lagrangian conditions for optimality in the program (4.67). From this fact, we make an important conclusion—namely, fixed points of the EP updates are guaranteed to exist, assuming that the optimization problem (4.67) has at least one optimum. As with the sum-product algorithm and the Bethe variational principle, there are no guarantees that the EP updates converge in general. However, it would be relatively straightforward to develop convergent algorithms for finding at least local optimal of the variational problem (4.67), or equivalently (4.63), possibly along the lines of convergent algorithms developed for the ordinary Bethe variational problem [265, 249, 108].

At a high level, the key points to take away are the following: within the variational framework, expectation-propagation algorithms are based on a Bethe-like entropy approximation (4.62), and a particular convex outer bound on the set of mean parameters (4.61). Moreover, the moment-matching steps in the EP algorithm (4.7) arise from a Lagrangian approach for attempting to solve the relaxed variational principle (4.63).

We conclude by illustrating the specific forms taken by Algorithm 4.7 for some concrete examples:

**Example 23 (Sum-product as moment-matching).** As a continuation of Example 22, we now show how the updates (4.74) and (4.75) of Algorithm 4.7 reduce to the sum-product updates. In this particular example, the index  $\ell$  ranges over edges of the graph. Associated with  $\ell = (u, v)$  are the pair of Lagrange multiplier vectors  $(\lambda_{uv}(x_v), \lambda_{vu}(x_u))$ . In terms of these quantities, the base distribution (4.71) can be written

as

$$\begin{aligned}
p(x; \theta, \lambda) &\propto \prod_{s \in V} \exp(\theta_s(x_s)) \prod_{(u,v) \in E} \exp(\lambda_{uv}(x_v) + \lambda_{vu}(x_u)) \\
&= \prod_{s \in V} \exp\left(\theta_s(x_s) + \sum_{t \in N(s)} \lambda_{ts}(x_s)\right) \\
&\propto \prod_{s \in V} \tau_s(x_s)
\end{aligned} \tag{4.76}$$

where we have defined the pseudo-marginals

$$\tau_s(x_s) \propto \exp\left(\theta_s(x_s) + \sum_{t \in N(s)} \lambda_{ts}(x_s)\right).$$

It is worthwhile noting the similarity to the sum-product expression (4.23) for the singleton pseudomarginals  $\tau_s$ , obtained in the proof of Theorem 3. Similarly, when  $\ell = (u, v)$ , the augmented distribution (4.72) is given by

$$\begin{aligned}
q_{(u,v)}(x; \theta; \lambda) &\propto p(x; \theta, \lambda) \exp(\theta_{uv}(x_u, x_v) - \lambda_{vu}(x_u) - \lambda_{uv}(x_v)) \\
&= \left[ \prod_{s \in V} \tau_s(x_s) \right] \exp(\theta_{uv}(x_u, x_v) - \lambda_{vu}(x_u) - \lambda_{uv}(x_v)).
\end{aligned} \tag{4.77}$$

Now consider the updates (4.74) and (4.75). If index  $\ell = (u, v)$  is chosen, the update (4.74) is equivalent to computing the singleton marginals of the distribution (4.77). The update (4.75) dictates that the Lagrange multipliers  $\lambda_{uv}(x_v)$  and  $\lambda_{vu}(x_u)$  should be adjusted so that the marginals  $\{\tau_u, \tau_v\}$  of the distribution (4.76) match these marginals. Following a little bit of algebra, these fixed point conditions reduce to the sum-product updates along edge  $(u, v)$ , where the messages are given as exponentiated Lagrange multipliers by  $M_{uv}(x_v) = \exp(\lambda_{uv}(x_v))$ .

♣

From the perspective of term-by-term approximations, the sum-product algorithm uses a product base distribution (see equation (4.64)). A natural extension, then, is to consider a base distribution with more structure.

**Example 24 (Tree-structured EP).** We illustrate this idea by deriving the tree-structured EP algorithm [171], applied to a pairwise Markov random field on a graph  $G = (V, E)$ . We use the same set-up and parameterization as in Examples 22 and 23. Given some fixed spanning tree  $T = (V, E(T))$  of the graph, the base distribution (4.53) is given by

$$p(x; \theta, \vec{0}) \propto \prod_{s \in V} \exp(\theta_s(x_s)) \prod_{(s,t) \in E(T)} \exp(\theta_{st}(x_s, x_t)). \quad (4.78)$$

In this case, the index  $i$  runs over edges  $(u, v) \in E \setminus E(T)$ ; given such an edge  $(u, v)$ , the  $\Phi^{uv}$ -augmented distribution (4.54) is obtained by adding in the term  $\theta_{uv}$ :

$$p(x; \theta, \theta_{uv}) \propto p(x; \theta, \vec{0}) \exp(\theta_{uv}(x_u, x_v)). \quad (4.79)$$

Let us now describe the analogs of the sets  $\mathcal{M}(\phi; \Phi)$ ,  $\mathcal{M}(\phi)$  and  $\mathcal{M}(\phi; \Phi_\ell)$  for this setting. First, as we saw in Example 22, the set of mean parameters associated with the full model (4.45) is given by

$$\mu := \{\mu_s, s \in V\} \cup \{\mu_{st}, (s, t) \in E\}.$$

Accordingly, the analog of  $\mathcal{M}(\phi; \Phi)$  is the marginal polytope  $\mathbb{M}(G)$ . Second, if we use the tree-structured distribution (4.78) as the base distribution, the relevant subset of mean parameters are

$$\mu(T) := \{\mu_s, s \in V\} \cup \{\mu_{st}, (s, t) \in E(T)\}.$$

The analog of the set  $\mathcal{M}(\phi)$  is the marginal polytope  $\mathbb{M}(T)$ , or equivalently—using Proposition 4—the set  $\mathbb{L}(T)$ . Finally, if we add in edge  $\ell = (u, v) \notin E(T)$ , then the analog of the set  $\mathcal{M}(\phi; \Phi)$  is the marginal polytope  $\mathbb{M}(T \cup (u, v))$ .

The analog of  $\mathcal{L}(\phi; \Phi)$  is an interesting object, which we refer to as the *tree-EP outer bound* on the marginal polytope  $\mathbb{M}(G)$ . It consists of all vectors  $\tau = \{\tau_s, s \in V\} \cup \{\tau_{st}, (s, t) \in E\}$  such that

- (a) the inclusion  $\tau(T) \in \mathbb{M}(T)$  holds, and
- (b) for all  $(u, v) \notin T$ , the inclusion  $(\tau(T), \tau_{uv}) \in \mathbb{M}(T \cup (u, v))$  holds.

Observe that the set thus defined is contained within  $\mathbb{L}(G)$ : in particular, the condition  $\tau(T) \in \mathbb{M}(T)$  ensures non-negativity, normalization,

and marginalization for all mean parameters associated with the tree, and the inclusion  $(\tau(T), \tau_{uv}) \in \mathbb{M}(T \cup (u, v))$  ensures that  $\tau_{uv}$  is non-negative, and satisfies the marginalization constraints associated with  $\tau_u$  and  $\tau_v$ . For graphs with cycles, this tree-EP outer bound is strictly contained within  $\mathbb{L}(G)$ ; for instance, for the single cycle  $C_3$  on three nodes, the tree-EP outer bound is equivalent to  $\mathbb{M}(C_3)$  (e.g., see Wainwright [235]). But the set  $\mathbb{M}(C_3)$  is strictly contained within  $\mathbb{L}(C_3)$ , as shown in Example 14.

The entropy approximation  $H_{\text{ep}}$  associated with tree-EP is easy to define. Given a candidate set of pseudo-mean-parameters  $\tau$  in the tree-EP outer bound, we define the tree-structured entropy associated with the base distribution (4.78)

$$H(\tau(T)) := \sum_{s \in V} H(\tau_s) - \sum_{(s,t) \in E(T)} I_{st}(\tau_{st}). \quad (4.80)$$

Similarly, for each edge  $(u, v) \notin E(T)$ , we define the entropy  $H(\tau(T), \tau_{uv})$  associated with the augmented distribution (4.79); unlike the base case, this entropy does not have an explicit form (since the graph is not in junction tree form), but it can be computed easily. The tree-EP entropy approximation then takes the form

$$H_{\text{ep}}(\tau) = H(\tau(T)) + \sum_{(u,v) \notin E(T)} [H(\tau(T), \tau_{uv}) - H(\tau(T))].$$

Finally, let us derive the moment-matching updates (4.74) and (4.75) for tree-EP. For each edge  $(u, v) \notin E(T)$ , let  $\lambda^{uv}(T)$  denote a vector of Lagrange multipliers, of the same dimension as  $\tau(T)$ , and let  $\phi(x; T)$  denote the subset of sufficient statistics associated with  $T$ . The optimized base distribution (4.71) can be represented in terms of the original base distribution and these tree-structured Lagrange multipliers as

$$q(x; \theta, \lambda) \propto p(x; \theta, \vec{0}) \prod_{(u,v) \notin E(T)} \exp(\langle \lambda^{uv}(T), \phi(x; T) \rangle).$$

If edge  $(u, v)$  is added, the augmented distribution (4.72) is given by

$$q^{uv}(x; \theta, \lambda) \propto p(x; \theta, \vec{0}) \exp(\theta_{uv}(x_u, x_v)) \prod_{(s,t) \neq (u,v)} \exp(\langle \lambda^{st}(T), \phi(x; T) \rangle).$$



For a given edge  $(u, v) \notin E(T)$ , the moment-matching step (4.75) corresponds to computing singleton marginals and pairwise marginals for edges  $(s, t) \in E(T)$  under the distribution  $q^{uv}(x; \theta, \lambda)$ . The step (4.75) corresponds to updating the Lagrange multiplier vector  $\lambda^{uv}(T)$  so the marginals of  $q(x; \theta; \lambda)$  agree with those of  $q^{uv}(x; \theta; \lambda)$  for all nodes  $s \in V$ , and all edges  $(s, t) \in E(T)$ . See the papers [171, 247, 235] for further details on tree EP.



The previous examples dealt with discrete random variables; we now return to the case of a mixture model with continuous variables, first introduced in Example 21.

**Example 25 (EP for Gaussian mixture models).** Recall from equation (4.56) the representation of the mixture distribution as an exponential family with the base potentials  $\phi(x) = \{x, xx^T\}$  and auxiliary terms  $\Phi^i(x) = \log p(y^i | x)$ , for  $i = 1, \dots, n$ . Turning to the analogs of the various mean parameter spaces, the set  $\mathcal{M}(\phi; \Phi)$  is the space of all globally realizable mean parameters

$$(\mathbb{E}[X], \mathbb{E}[XX^T], \mathbb{E}[\log p(y^i | X)]), \quad i = 1, \dots, n).$$

Recall that  $(y^1, \dots, y^\ell)$  are observed, and hence remain fixed throughout.

The set  $\mathcal{M}(\phi)$  is simply the mean parameter space for a multivariate Gaussian, as characterized in Example 9. For each  $i = 1, \dots, n$ , the constraint set  $\mathcal{M}(\phi; \Phi^i)$  corresponds to the collection of all globally realizable mean parameters of the form  $(\mathbb{E}[X], \mathbb{E}[XX^T], \mathbb{E}[\log p(y^i | X)])$ . Turning to the various entropies, the base term is the entropy  $H(\mu)$  of a multivariate Gaussian where  $\mu$  denotes the pair  $(\mathbb{E}[X], \mathbb{E}[XX^T])$ ; similarly, the augmented term  $H(\mu, \tilde{\mu}^i)$  is the entropy of the exponential family member with sufficient statistics  $(\phi(x), \log p(y^i | x))$ , and mean parameters  $(\mu, \tilde{\mu}^i)$ , where  $\tilde{\mu}^i := \mathbb{E}[\log p(y^i | X)]$ . Using these quantities, we can define the analog of the variational principle (4.63).

Turning to the moment-matching steps (4.74) and (4.75), the Lagrangian version of Gaussian-mixture EP algorithm can be formulated in terms of a collection of  $n$  Lagrange multipliers

$$(\lambda^i, \Lambda^i) \in \mathbb{R}^m \times \mathbb{R}^{m \times m}, \quad i = 1, \dots, n,$$


one for each augmented distribution. Since  $X \sim N(0, \Sigma)$  by definition, the optimized base distribution (4.71) can be written in terms of  $\Sigma^{-1}$  and the Lagrange multipliers

$$q(x; \Sigma; (\lambda, \Lambda)) \propto \exp\left(\left\langle \sum_{i=1}^n \lambda^i, x \right\rangle + \left\langle \left\langle -\frac{1}{2}\Sigma^{-1} + \sum_{i=1}^n \Lambda^i, xx^T \right\rangle \right\rangle\right). \quad (4.81)$$

Note that this is simply a multivariate Gaussian distribution. The augmented distribution (4.72) associated with any term  $\ell \in \{1, 2, \dots, d_I\}$  is given by

$$q^\ell(x) = f_0(x) \exp\left\{ \left\langle \sum_{i \neq \ell} \lambda^i, x \right\rangle + \left\langle \left\langle -\frac{1}{2}\Sigma^{-1} + \sum_{i \neq \ell} \Lambda^i, xx^T \right\rangle \right\rangle + \left\langle \tilde{\theta}_\ell, \log p(y^\ell | x) \right\rangle \right\}. \quad (4.82)$$

With this set-up, the step (4.74) corresponds to computing the mean parameters  $(\mathbb{E}_{q^\ell}[X], \mathbb{E}_{q^\ell}[XX^T])$  under the distribution (4.82), and step (4.75) corresponds to adjusting the multivariate Gaussian (4.81) so that it has these mean parameters. Fixed points of these moment-matching updates satisfy the necessary Lagrangian conditions to be optima of the associated Bethe-like approximation of the exact variational principle.

We refer the reader to Minka [169, 172] and Seeger [208] for further details of the Gaussian-mixture EP algorithm and some of its properties. 

# 5

---

## Mean field methods

---

This section is devoted to a discussion of mean field methods, which originated in the statistical physics literature [e.g., 186, 11, 44]. From the perspective of this paper, the mean field approach is based on a specific type of approximation to the exact variational principle (3.45). More specifically, as discussed in Section 3.7, there are two fundamental difficulties associated with the variational principle (3.45): the nature of the constraint set  $\mathcal{M}$ , and the lack of an explicit form for the dual function  $A^*$ . The core idea of mean field approaches is simple: let us limit the optimization to a subset of distributions for which both  $\mathcal{M}$  and  $A^*$  are relatively easy to characterize; e.g., perhaps they correspond to a graph with small treewidth. Throughout this section, we refer to any such distribution as “tractable.” The simplest choice is the family of product distributions, which gives rise to the naive mean field method. Higher-order mean field methods are obtained by choosing tractable distributions with more structure.

## 5.1 Tractable families

We base our description of mean field methods on the notion of a *tractable subgraph*, by which we mean a subgraph  $F$  of  $G = (V, E)$  over which it is feasible to perform exact calculations. The simplest example of a tractable subgraph is the fully disconnected subgraph  $F_0 = (V, \emptyset)$ , which contains all the vertices of  $G$  but *none* of the edges. Any distribution that is Markov with respect to  $F$  is then a product distribution, for which exact computations are trivial.

A bit more generally, consider an exponential family with a collection  $\phi = \{\phi_\alpha \mid \alpha \in \mathcal{I}\}$  of sufficient statistics associated with the cliques of  $G = (V, E)$ . Given a subgraph  $F$ , let  $\mathcal{I}(F) \subseteq \mathcal{I}$  be the subset of sufficient statistics associated with cliques of  $F$ . The set of all distributions that are Markov with respect to  $F$  is a sub-family of the full  $\phi$ -exponential family; it is parameterized by the subspace of canonical parameters

$$\Omega(F) := \{\theta \in \Omega \mid \theta_\alpha = 0 \quad \forall \alpha \in \mathcal{I} \setminus \mathcal{I}(F)\}. \quad (5.1)$$

We consider some examples to illustrate:

**Example 26 (Tractable subgraphs).** Suppose that  $\theta \in \Omega$  parameterizes a pairwise Markov random field, with potential functions associated with the vertices and edges of an undirected graph  $G = (V, E)$ . For each edge  $(s, t) \in E$ , let  $\theta_{(s,t)}$  denote the sub-vector of parameters associated with sufficient statistics that depend only on  $(X_s, X_t)$ . Consider the completely disconnected subgraph  $F_0 = (V, \emptyset)$ . With respect to this subgraph, permissible parameters must belong to the subspace

$$\Omega(F_0) := \{\theta \in \Omega \mid \theta_{(s,t)} = 0 \quad \forall (s, t) \in E\}, \quad (5.2)$$

The densities in this sub-family are all of the *fully factorized or product form*

$$p_\theta(x) = \prod_{s \in V} p(x_s; \theta_s), \quad (5.3)$$

where  $\theta_s$  refers to the collection of canonical parameters associated with vertex  $s$ .

To obtain a more structured approximation, one could choose a spanning tree  $T = (V, E(T))$ . In this case, we are free to choose the canonical parameters corresponding to vertices and edges in the tree  $T$ , but we must set to zero any canonical parameters corresponding to edges not in the tree. Accordingly, the subspace of tree-structured distributions is specified by the subset of canonical parameters

$$\Omega(T) := \{\theta \in \Omega \mid \theta_{(s,t)} = 0 \ \forall (s,t) \notin E(T)\}. \quad (5.4)$$



Associated with the exponential family defined by  $\phi$  and  $G$  is the set  $\mathcal{M}(G; \phi)$  of all mean parameters realizable by any distribution, as previously defined in equation (3.26). (Whereas our previous notation did not make explicit reference to  $G$  and  $\phi$ , the definition of  $\mathcal{M}$  does depend on these quantities and it is now useful to make this dependence explicit.) For a given tractable subgraph  $F$ , mean field methods are based on optimizing over the subset of mean parameters that can be obtained by the subset of exponential family densities  $\{p_\theta, \theta \in \Omega(F)\}$ —namely

$$\mathcal{M}_F(G; \phi) := \{\mu \in \mathbb{R}^d \mid \mu = \mathbb{E}_\theta[\phi(x)] \text{ for some } \theta \in \Omega(F)\}. \quad (5.5)$$

In terms of the moment mapping from Theorem 1, a more compact definition of the set  $\mathcal{M}_F(G; \phi)$  is as the image  $\nabla A(\Omega(F))$ . By Theorem 1, we have  $\mathcal{M}^\circ(G; \phi) = \nabla A(\Omega)$  so that the inclusion

$$\mathcal{M}_F^\circ(G; \phi) \subseteq \mathcal{M}^\circ(G; \phi)$$

holds for any subgraph  $F$ . For this reason, we say that  $\mathcal{M}_F$  is an *inner approximation* to the set  $\mathcal{M}$  of realizable mean parameters.

To lighten the notation in the remainder of this section, we generally drop the  $\phi$  term from  $\mathcal{M}(G; \phi)$  and  $\mathcal{M}_F(G; \phi)$ , writing  $\mathcal{M}(G)$  and  $\mathcal{M}_F(G)$  respectively. It is important to keep in mind, though, that these sets do depend on the choice of sufficient statistics.

## 5.2 Optimization and lower bounds

We now have the necessary ingredients to develop the mean field approach to approximate inference. Suppose that we are interested in

approximating some target distribution  $p_\theta$ , where  $\theta \in \Omega$ . Mean field methods generate lower bounds on the value  $A(\theta)$  of the cumulant function, as well as approximations to the mean parameters  $\mu = \mathbb{E}_\theta[\phi(X)]$  of this target distribution  $p_\theta$ .

### 5.2.1 Generic mean field procedure

The key property of any mean field method is the following fact: any valid mean parameter specifies a lower bound on the log partition function.

**Proposition 7 (Mean field lower bound).** *Any mean parameter  $\mu \in \mathcal{M}^\circ$  yields a lower bound on the cumulant function:*

$$A(\theta) \geq \langle \theta, \mu \rangle - A^*(\mu). \quad (5.6)$$

Moreover, equality holds if and only if  $\theta$  and  $\mu$  are dually coupled.

*Proof.* In convex analysis, this lower bound is known as *Fenchel's inequality*, and is an immediate consequence of the general variational principle (3.45) (the supremum over  $\mu$  is always greater than the value at any particular  $\mu$ ). Moreover, the supremum is attained when  $\mu = \nabla A(\theta)$  which means that  $\theta$  and  $\mu$  are dually coupled (by definition).

An alternative proof via Jensen's inequality is also possible, and we give a version of this proof here given its popularity in the literature on mean field methods. By definition of  $\mathcal{M}$ , for any mean parameter  $\mu \in \mathcal{M}$ , there must exist some density  $q$ , defined with respect to the base measure  $\nu$  underlying the exponential family, for which  $\mathbb{E}_q[\phi(X)] = \mu$ . We then have

$$\begin{aligned} A(\theta) &= \log \int_{\mathcal{X}^m} q(x) \frac{\exp\{\langle \theta, \phi(x) \rangle\}}{q(x)} \nu(dx) \\ &\stackrel{(a)}{\geq} \int_{\mathcal{X}^m} q(x) [\langle \theta, \phi(x) \rangle - \log q(x)] \nu(dx) \\ &\stackrel{(b)}{=} \langle \theta, \mu \rangle + H(q), \end{aligned} \quad (5.7)$$

where  $H(q) = -\mathbb{E}_q[\log q(X)]$  is the entropy of  $q$ . In this argument, step (a) follows from Jensen's inequality (see Appendix A.2.5) applied

to the negative logarithm, whereas step (b) follows from the moment-matching condition  $\mathbb{E}_q[\phi(X)] = \mu$ . Since the lower bound (5.7) holds for any density  $q$  satisfying this moment-matching condition, we may optimize over the choice of  $q$ : by Theorem 2, doing so yields the exponential family density  $q^*(x) = p_{\theta(\mu)}(x)$ , for which  $H(q^*) = -A^*(\mu)$  by construction.  $\square$

Since the dual function  $A^*$  typically lacks an explicit form, it is not possible, at least in general, to compute the lower bound (5.6). The mean field approach circumvents this difficulty by restricting the choice of  $\mu$  to the tractable subset  $\mathcal{M}_F(G)$ , for which the dual function has an explicit form. For compactness in notation, we define  $A_F^* = A^*|_{\mathcal{M}_F(G)}$ , corresponding to the dual function restricted to the set  $\mathcal{M}_F(G)$ . As long as  $\mu$  belongs to  $\mathcal{M}_F(G)$ , then the lower bound (5.6) involves  $A_F^*$ , and hence can be computed easily.

The next step of the mean field method is the natural one: find the best approximation, as measured in terms of the tightness of the lower bound (5.6). More precisely, the best lower bound from within  $\mathcal{M}_F(G)$  is given by

$$\max_{\mu \in \mathcal{M}_F(G)} \{ \langle \mu, \theta \rangle - A_F^*(\mu) \}. \quad (5.8)$$

The corresponding value of  $\mu$  is defined to be the mean field approximation to the true mean parameters.

In Section 5.3, we illustrate the use of this generic procedure in obtaining lower bounds and approximate mean parameters for various types of graphical models.

### 5.2.2 Mean field and Kullback-Leibler divergence

An important alternative interpretation of the mean field optimization problem (5.8) is as minimizing the Kullback-Leibler (KL) divergence between the approximating (tractable) distribution and the target distribution. In order to make this connection clear, we first digress to discuss various forms of the KL divergence for exponential family models.

The conjugate duality between  $A$  and  $A^*$ , as characterized in Theorem 2, leads to several alternative forms of the Kullback-Leibler (KL) divergence for exponential family members. The standard definition of

the KL divergence [54] between two distributions with densities  $q$  and  $p$  with respect to a base measure  $\nu$  is

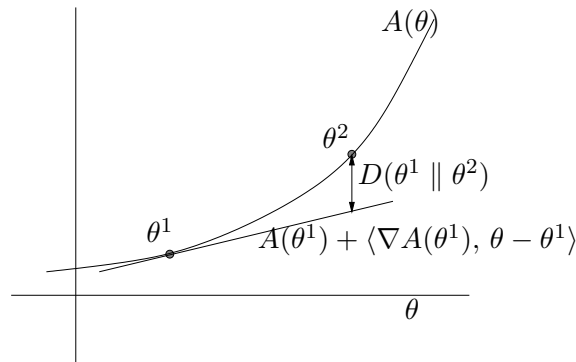
$$D(q \parallel p) := \int_{\mathcal{X}^m} \left[ \log \frac{q(x)}{p(x)} \right] q(x) \nu(dx). \quad (5.9)$$

The key result that underlies alternative representations for exponential families is Proposition 7.

Consider two canonical parameter vectors  $\theta^1, \theta^2 \in \Omega$ ; with a slight abuse of notation, we use  $D(\theta^1 \parallel \theta^2)$  to refer to the KL divergence between  $p_{\theta^1}$  and  $p_{\theta^2}$ . We use  $\mu^1$  and  $\mu^2$  to denote the respective dually-coupled mean parameters. A first alternative form of the KL divergence is obtained by substituting the exponential representations of  $p_{\theta^i}$  into equation (5.9), and then expanding and simplifying as follows:

$$\begin{aligned} D(\theta^1 \parallel \theta^2) &= \mathbb{E}_{\theta^1} \left[ \log \frac{p_{\theta^1}(x)}{p_{\theta^2}(x)} \right] \\ &= A(\theta^2) - A(\theta^1) - \langle \mu^1, \theta^2 - \theta^1 \rangle. \end{aligned} \quad (5.10)$$

We refer to this representation as the *primal form* of the KL divergence. As illustrated in Figure 5.1, this form of the KL divergence can be



**Fig. 5.1.** The hyperplane  $A(\theta^1) + \langle \nabla A(\theta^1), \theta - \theta^1 \rangle$  supports the epigraph of  $A$  at  $\theta^1$ . The Kullback-Leibler divergence  $D(\theta^1 \parallel \theta^2)$  is equal to the difference between  $A(\theta^2)$  and this hyperplane.

interpreted as the difference between  $A(\theta^2)$  and the hyperplane tangent to  $A$  at  $\theta^1$  with normal  $\nabla A(\theta^1) = \mu^1$ . This interpretation shows that the KL divergence is a particular example of a Bregman distance [39, 42].



A second form of the KL divergence can be obtained by using the fact that equality holds in Proposition 7 for dually-coupled parameters. Thus, considering the dually-coupled pair  $(\theta^1, \mu^1)$ , we can transform equation (5.10) into the following *mixed form* of the KL divergence:

$$D(\theta^1 \parallel \theta^2) \equiv D(\mu^1 \parallel \theta^2) = A(\theta^2) + A^*(\mu^1) - \langle \mu^1, \theta^2 \rangle. \quad (5.11)$$

Note that this mixed form of the divergence corresponds to the slack in the inequality (5.6). It also provides an alternative view of the variational representation given in Theorem 2(b). In particular, equation (3.45) can be rewritten as follows:

$$\min_{\mu \in \mathcal{M}} \{A(\theta) + A^*(\mu) - \langle \theta, \mu \rangle\} = 0.$$

Using equation (5.11), the variational representation in Theorem 2(b) is seen to be equivalent to the assertion that  $\min_{\mu \in \mathcal{M}} D(\mu \parallel \theta) = 0$ .

Finally, by applying equation (5.6) as an equality once again, this time for the coupled pair  $(\theta^2, \mu^2)$ , the mixed form (5.11) can be transformed into a purely *dual form* of the KL divergence:

$$\begin{aligned} D(\theta^1 \parallel \theta^2) &\equiv D(\mu^1 \parallel \mu^2) \\ &= A^*(\mu^1) - A^*(\mu^2) - \langle \theta^2, \mu^1 - \mu^2 \rangle. \end{aligned} \quad (5.12)$$

Note the symmetry between representations (5.10) and (5.12). This form of the KL divergence has an interpretation analogous to that of Figure 5.1, but with  $A$  replaced by the dual  $A^*$ .

With this background on the Kullback-Leibler divergence, let us now return to the consequences for mean field methods. For a given mean parameter  $\mu \in \mathcal{M}_F(G)$ , the difference between the log partition function  $A(\theta)$  and the quantity  $\langle \mu, \theta \rangle - A_F^*(\mu)$  to be maximized is equivalent to

$$D(\mu \parallel \theta) = A(\theta) + A_F^*(\mu) - \langle \mu, \theta \rangle,$$

corresponding to the mixed form of the Kullback-Leibler divergence defined in equation (5.11). On the basis of this relation, it can be seen that solving the variational problem (5.8) is equivalent to minimizing the KL divergence  $D(\mu \parallel \theta)$ , subject to the constraint that  $\mu \in \mathcal{M}_F(G)$ . Consequently, any mean field method can be understood as obtaining the “best” approximation to  $p_\theta$  from a family of tractable models, where approximation quality is measured by the Kullback-Leibler divergence.

### 5.3 Examples of mean field algorithms

In this section, we illustrate various instances of mean field algorithms for particular graphical models.

#### 5.3.1 Naive mean field updates

The simplest type of approximation, referred to as the naive mean field approach, is based on choosing a *product distribution*

$$p_{\theta}(x_1, x_2, \dots, x_m) := \prod_{s \in V} p(x_s; \theta_s) \quad (5.13)$$

as the tractable approximation. The naive mean field updates are a particular set of recursions for finding a stationary point of the resulting optimization problem. In the case of the Ising model, the naive mean field updates are a classical set of recursions from statistical physics, typically justified in a heuristic manner in terms of “self-consistency.” For the case of a Gaussian Markov random field, the naive mean field updates turn out to be equivalent to the Gauss-Jacobi or Gauss-Seidel algorithm for solving a linear system of equations.

**Example 27 (Naive mean field for Ising model).** As an illustration, we derive the naive mean field updates for the Ising model, first introduced in Example 3. Recall that the Ising model is characterized by the sufficient statistics  $\{x_s, s \in V\}$  and  $\{x_s x_t \mid (s, t) \in E\}$ . The associated mean parameters are the singleton and pairwise marginal probabilities

$$\mu_s = \mathbb{E}[X_s] = \mathbb{P}[X_s = 1], \text{ and } \mu_{st} = \mathbb{E}[X_s X_t] = \mathbb{P}[X_s = 1, X_t = 1],$$

respectively. The full vector  $\mu$  of mean parameters is an element of  $\mathbb{R}^{|V|+|E|}$ .

Letting  $F_0$  denote the fully disconnected graph—that is, without any edges—the tractable set  $\mathcal{M}_{F_0}(G)$  consists of all mean parameters  $\{\mu_s, \mu_{st}\}$  that arise from the product distribution (5.13). Explicitly, in

this binary case, we have

$$\mathcal{M}_{F_0}(G) := \left\{ \mu \in \mathbb{R}^{|V|+|E|} \mid 0 \leq \mu_s \leq 1 \quad \forall s \in V, \text{ and} \right. \\ \left. \mu_{st} = \mu_s \mu_t \quad \forall (s, t) \in E. \right\}, \quad (5.14)$$

where the constraints  $\mu_{st} = \mu_s \mu_t$  arise from the product nature of any distribution that is Markov with respect to  $F_0$ .

For any  $\mu \in \mathcal{M}_{F_0}(G)$ , the value of the dual function—that is, the negative entropy of a product distribution—has an explicit form in terms of  $\{\mu_s, s \in V\}$ . In particular, a straightforward computation shows that this entropy takes the form

$$\begin{aligned} -A_{F_0}^*(\mu) &= -\sum_{s \in V} [\mu_s \log \mu_s + (1 - \mu_s) \log(1 - \mu_s)] \\ &= \sum_{s \in V} H_s(\mu_s). \end{aligned} \quad (5.15)$$

With these two ingredients, we can derive the specific form of the mean field optimization problem (5.8) for the product distribution and the Ising model. Given the product structure of the tractable family, any mean parameter  $\mu \in \mathcal{M}_{F_0}(G)$  satisfies the equality  $\mu_{st} = \mu_s \mu_t$  for all  $(s, t) \in E$ . Combining the expression for the entropy in 5.15 with the characterization of  $\mathcal{M}_{F_0}(G)$  in 5.14 yields the naive mean field problem

$$A(\theta) \geq \max_{(\mu_1, \dots, \mu_m) \in [0, 1]^m} \left\{ \sum_{s \in V} \theta_s \mu_s + \sum_{(s, t) \in E} \theta_{st} \mu_s \mu_t + \sum_{s \in V} H_s(\mu_s) \right\}. \quad (5.16)$$

For any  $s \in V$ , this objective function is strictly concave as a scalar function of  $\mu_s$  with all other coordinates held fixed. Moreover, the maximum over  $\mu_s$  with the other variables  $\mu_t, t \neq s$  held fixed is attained in the open interval  $(0, 1)$ . Indeed, by taking the derivative with respect to  $\mu_s$ , setting it to zero and performing some algebra, we obtain the update equation

$$\mu_s \leftarrow \sigma\left(\theta_s + \sum_{t \in N(s)} \theta_{st} \mu_t\right), \quad (5.17)$$

where  $\sigma(z) := [1 + \exp(-z)]^{-1}$  is the logistic function. Thus, we have derived—from a variational perspective—the naive mean field updates presented earlier (2.13).

What about convergence properties and the nature of the fixed points? Applying equation (5.17) iteratively to each node in succession amounts to performing coordinate ascent of the mean field variational problem (5.16). Since the maximum is uniquely attained for every coordinate update, known results on coordinate ascent methods [19] imply that any sequence  $\{\mu^0, \mu^1, \dots\}$  generated by the updates (5.17) is guaranteed to converge to a local optimum of the naive mean field problem (5.16).

Unfortunately, the mean field problem is *non-convex* in general, so that there may be multiple local optima, and the limit point of the sequence  $\{\mu^0, \mu^1, \dots\}$  can depend strongly on the initialization  $\mu^0$ . We discuss this non-convexity and its consequences at greater depth in Section 5.4. Despite these issues, the naive mean field approximation becomes asymptotically exact for certain types of models as the number of nodes  $m$  grows to infinity [11, 262]. An example is the ferromagnetic Ising model defined on the complete graph  $K_m$  with suitably rescaled parameters  $\theta_{st} > 0$  for all  $(s, t) \in E$ ; see Baxter [11] for further discussion of such exact cases. ♣

Similarly, it is straightforward to apply the naive mean field approximation to other types of graphical models, as we illustrate for a multivariate Gaussian.

**Example 28 (Gaussian mean field).** Recall the Gaussian Markov random field, first discussed as an exponential family in Example 5. Its mean parameters consist of the vector  $\mu = \mathbb{E}[X] \in \mathbb{R}^m$ , and the symmetric<sup>1</sup> matrix  $\Sigma = \mathbb{E}[XX^T] \in \mathcal{S}_+^m$ . Suppose that we again use the completely disconnected graph  $F_0 = (V, \emptyset)$  as the tractable class. Any Gaussian distribution that is Markov with respect to  $F_0$  must have a diagonal covariance matrix, meaning that the set of tractable mean

<sup>1</sup>Strictly speaking, only elements  $[\Sigma]_{st}$  with  $(s, t) \in E$  are included in the mean parameterization, since the associated canonical parameter is zero for any pair  $(u, v) \notin E$ .

parameters takes the form

$$\mathcal{M}_{F_0}(G) = \{(\mu, \Sigma) \in \mathbb{R}^m \times \mathcal{S}_+^m \mid \Sigma - \mu\mu^T = \text{diag}(\Sigma - \mu\mu^T) \succeq 0\}. \quad (5.18)$$

For any such product distribution, the entropy (negative dual function) has the form

$$\begin{aligned} -A_{F_0}^*(\mu, \Sigma) &= \frac{m}{2} \log 2\pi e + \frac{1}{2} \log \det [\Sigma - \mu\mu^T] \\ &= \frac{m}{2} \log 2\pi e + \frac{1}{2} \sum_{s=1}^m \log(\Sigma_{ss} - \mu_s^2). \end{aligned}$$

Combining this form of the dual function with the constraints (5.18) yields that for a multivariate Gaussian, the value  $A(\theta)$  of the cumulant function is lower bounded by

$$\max_{\{(\mu, \Sigma) \mid \Sigma_{ss} - \mu_s^2 > 0, \Sigma_{st} = \mu_s \mu_t\}} \left\{ \langle \theta, \mu \rangle + \frac{1}{2} \langle \langle \Theta, \Sigma \rangle \rangle + \frac{1}{2} \sum_{s=1}^m (\Sigma_{ss} - \mu_s^2) + \frac{m}{2} \log 2\pi e \right\}, \quad (5.19)$$

where  $(\theta, \Theta)$  are the canonical parameters associated with the multivariate Gaussian (see Example 5). The optimization problem (5.19) yields the naive mean field lower bound for the multivariate Gaussian.

This optimization problem can be further simplified by substituting the constraints  $\Sigma_{st} = \mu_s \mu_t$  directly into the term  $\langle \langle \Theta, \Sigma \rangle \rangle = \sum_{s,t} \Theta_{st} \Sigma_{st}$  that appears in the objective function. Doing so and then taking derivatives with respect to the remaining optimization variables (namely,  $\mu_s$  and  $\mu_{ss}$ ) yields the stationary conditions

$$\frac{1}{2(\mu_{ss} - \mu_s^2)} = -\Theta_{ss}, \quad \text{and} \quad \frac{\mu_s}{2(\mu_{ss} - \mu_s^2)} = \theta_s + \sum_{t \in N(s)} \Theta_{st} \mu_t. \quad (5.20)$$

One particular set of updates for solving these fixed point equations is the iteration

$$\mu_s \leftarrow -\frac{1}{\Theta_{ss}} \left\{ \theta_s + \sum_{t \in N(s)} \Theta_{st} \mu_t \right\}. \quad (5.21)$$

Are these updates convergent, and what is the nature of the fixed points? Interestingly, the updates (5.21) are equivalent, depending on

the particular ordering used, to either the *Gauss-Jacobi* or the *Gauss-Seidel* methods [64] for solving the normal equations  $\mu = -\Theta^{-1}\theta$ . Therefore, if the Gaussian mean field updates (5.21) converge, they compute the correct mean vector. Moreover, the convergence behavior of such updates is well understood: for instance, the updates (5.21) are guaranteed to converge whenever  $-\Theta$  is strictly diagonally dominant; see Demmel [64] for further details on such Gauss-Jacobi and Gauss-Seidel iterations for solving matrix-vector equations.



### 5.3.2 Structured mean field

Of course, the essential principles underlying the mean field approach are not limited to fully factorized distributions. More generally, we can consider classes of tractable distributions that incorporate additional structure. This *structured mean field approach* was first proposed by Saul and Jordan [204], and further developed by various researchers [9, 254, 117].

Here we capture the structured mean field idea by discussing a general form of the updates for an approximation based on an arbitrary subgraph  $F$  of the original graph  $G$ . We do not claim that these specific updates are the best for practical purposes; the main goal here is the conceptual one of understanding the structure of the solution. Depending on the particular context, various techniques from nonlinear programming might be suitable for solving the mean field problem (5.8).

Let  $\mathcal{I}(F)$  be the subset of indices corresponding to sufficient statistics associated with  $F$ , and let  $\mu(F) := \{\mu_\alpha \mid \alpha \in \mathcal{I}(F)\}$  be the associated set of mean parameters. We let  $\mathcal{M}(F)$  denote the set of realizable mean parameters defined by the subgraph  $F$  and by the subset of sufficient statistics picked out by  $\mathcal{I}(F)$ ; it is a subset of  $\mathbb{R}^{|\mathcal{I}(F)|}$ . It is important to note that  $\mathcal{M}(F)$  differs from the previously defined (5.5) set  $\mathcal{M}_F(G)$ , which is based on the entire set of sufficient statistics  $\phi$ , and so is a subset of  $\mathbb{R}^{|\mathcal{I}|}$ .

We then observe that the mean field problem (5.8) has the following key properties:

- (a) the subvector  $\mu(F)$  can be an arbitrary member of  $\mathcal{M}(F)$ .

- (b) the dual function  $A_F^*$  actually depends only on  $\mu(F)$ , and *not* on mean parameters  $\mu_\beta$  for indices  $\beta$  in the complement  $\mathcal{I}^c(F) := \mathcal{I}(G) \setminus \mathcal{I}(F)$ .

Of course, the mean parameters  $\mu_\beta$  for indices  $\beta \in \mathcal{I}^c(F)$  do play a role in the problem; in particular, they arise within the linear term  $\langle \mu, \theta \rangle$ . Moreover, each mean parameter  $\mu_\beta$  is constrained in a nonlinear way by the choice of  $\mu(F)$ . Accordingly, for each  $\beta \in \mathcal{I}^c(F)$ , we write  $\mu_\beta = g_\beta(\mu(F))$  for some nonlinear function  $g_\beta$ , of which particular examples are given below. Based on these observations, the optimization problem (5.8) can be rewritten in the form

$$\begin{aligned} & \max_{\mu \in \mathcal{M}_F(G)} \{ \langle \theta, \mu \rangle - A_F^*(\mu) \} \\ &= \max_{\mu(F) \in \mathcal{M}(F)} \left\{ \underbrace{\sum_{\alpha \in \mathcal{I}(F)} \theta_\alpha \mu_\alpha + \sum_{\alpha \in \mathcal{I}^c(F)} \theta_\alpha g_\alpha(\mu(F)) - A_F^*(\mu(F))}_{G(\mu(F))} \right\}. \end{aligned} \quad (5.22)$$

On the left-hand side, the optimization takes place over the vector  $\mu \in \mathcal{M}_F(G)$ , which is of the same dimension as  $\theta \in \Omega \subseteq \mathbb{R}^{|\mathcal{I}|}$ . The objective function  $G$  for the optimization on the right-hand side, in contrast, is a function only of the lower-dimensional vector  $\mu(F) \in \mathcal{M}(F) \subseteq \mathbb{R}^{|\mathcal{I}(F)|}$ .

To illustrate this transformation, consider the case of naive mean field for the Ising model, where  $F \equiv F_0$  is the completely disconnected graph. In this case, each edge  $(s, t) \in E$  corresponds to an index in the set  $\mathcal{I}^c(F_0)$ ; moreover, for any such edge, we have  $g_{st}(\mu(F_0)) = \mu_s \mu_t$ . Since  $F_0$  is the completely disconnected graph,  $\mathcal{M}(F_0)$  is simply the hypercube  $[0, 1]^m$ . Therefore, for this particular example, the right-hand side of equation (5.22) is equivalent to equation (5.16).

Taking the partial derivatives of the objective function  $G$  with respect to some  $\mu_\beta$  with  $\beta \in \mathcal{I}(F)$  yields:

$$\frac{\partial G}{\partial \mu_\beta}(\mu(F)) = \theta_\beta + \sum_{\alpha \in \mathcal{I}(G) \setminus \mathcal{I}(F)} \theta_\alpha \frac{\partial g_\alpha}{\partial \mu_\beta}(\mu(F)) - \frac{\partial A_F^*}{\partial \mu_\beta}(\mu(F)). \quad (5.23)$$

Setting these partial derivatives to zero and re-arranging yields the

fixed point condition

$$\nabla A_F^*(\mu(F)) = \theta + \sum_{\alpha \in \mathcal{I}(G) \setminus \mathcal{I}(F)} \theta_\alpha \nabla g_\alpha(\mu(F)).$$

To obtain a more intuitive form of this fixed point equation, recall from Proposition 2 that the gradient  $\nabla A$  defines the forward mapping, from canonical parameters  $\theta$  to mean parameters  $\mu$ . Similarly, as we have shown in Section 3.6, the gradient  $\nabla A^*$  defines the backward mapping from mean parameters to canonical parameters. Letting  $\gamma(F)$  denote the canonical parameters that are dually coupled with  $\mu(F)$ , we can rewrite the fixed point update for component  $\beta$  as

$$\gamma_\beta(F) \leftarrow \theta_\beta + \sum_{\alpha \in \mathcal{I}(G) \setminus \mathcal{I}(F)} \theta_\alpha \frac{\partial g_\alpha}{\partial \mu_\beta}(\mu(F)). \quad (5.24)$$

After any such update, it is then necessary to adjust all mean parameters  $\mu_\delta(F)$  that depend on  $\gamma_\beta(F)$ —for instance, via junction tree updates—so that global consistency is maintained for the tractable approximation.

Alternatively, letting  $A_F$  be the conjugate dual of the function  $A_F^*$  previously defined, we have  $\nabla A_F(\gamma(F)) = \mu(F)$ . This fact follows from the Legendre duality between  $(A_F, A_F^*)$  (see Appendix B.3 for background), and the usual cumulant properties summarized in Proposition 2. By exploiting this relation, we obtain an alternative form of the update (5.24), one which involves only the mean parameters  $\mu(F)$ :

$$\mu_\beta(F) \leftarrow \frac{\partial A_F}{\partial \gamma_\beta} \left( \theta + \sum_{\alpha \in \mathcal{I}(G) \setminus \mathcal{I}(F)} \theta_\alpha \nabla g_\alpha(\mu(F)) \right). \quad (5.25)$$

We illustrate the updates (5.24) and (5.25) with some examples.

**Example 29 (Re-derivation of naive MF updates).** We first check that equation (5.25) reduces to the naive mean field updates (5.17), when  $F = F_0$  is the completely disconnected graph. For any product distribution on the Ising model, the subset of relevant mean parameters is given by  $\mu(F_0) = (\mu_1, \dots, \mu_m)$ . By the properties of the product



distribution, we have  $g_{st}(\mu(F_0)) = \mu_s \mu_t$  for all edges  $(s, t)$ , so that

$$\frac{\partial g_{st}}{\partial \mu_\alpha} = \begin{cases} \mu_t & \text{if } \alpha = s \\ \mu_s & \text{if } \alpha = t \\ 0 & \text{otherwise.} \end{cases}$$

Thus, for a given node index  $s \in V$ , the right-hand side of (5.24) is given by the sum

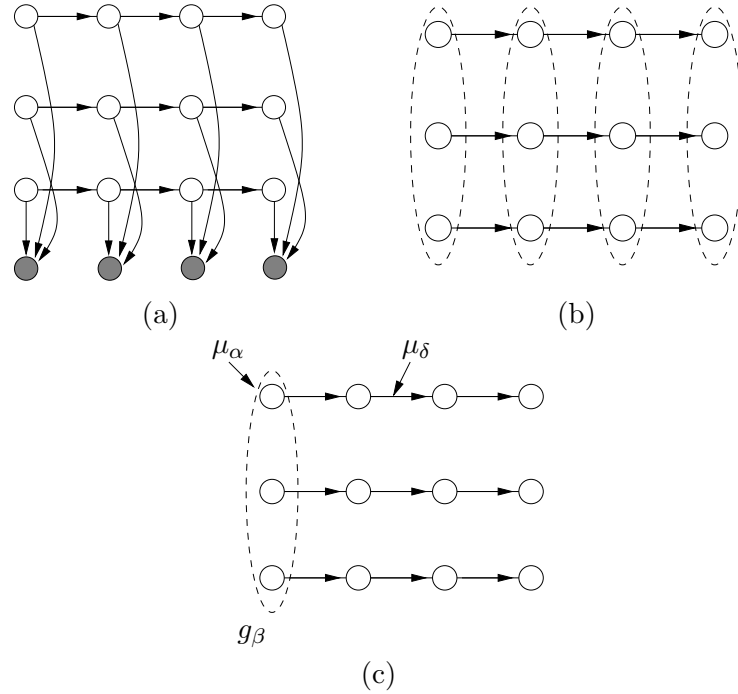
$$\theta_s + \sum_{t \in N(s)} \theta_{st} \mu_t. \quad (5.26)$$

Next, we observe that for a fully disconnected graph, the function  $A_{F_0}^*$  corresponds to a sum of negative entropies—that is,  $A_{F_0}^*(\mu(F_0)) = -\sum_{s \in V} H(\mu_s)$ —and so is a separable function. Consequently, the cumulant function  $A_{F_0}$  is also separable, and of the form  $A_{F_0}(\gamma(F_0)) = \sum_{s \in V} \log(1 + \exp(\gamma_s))$ . The partial derivatives are given by

$$\frac{\partial A_{F_0}}{\partial \gamma_\beta}(\gamma(F_0)) = \frac{\exp(\gamma_\beta)}{1 + \exp(\gamma_\beta)}, \quad (5.27)$$

which is the logistic function. Combining the expression (5.26) with the logistic function form (5.27) of the partial derivative shows that equation (5.25), when specialized to product distributions and the Ising model, reduces to the naive mean field update (5.17). ♣

**Example 30 (Structured MF for factorial HMMs).** To provide a more interesting example of the updates (5.24), consider a factorial hidden Markov model, as described in Ghahramani and Jordan [89]. Figure 5.2(a) shows the original model, which consists of a set of  $M$  Markov chains ( $M = 3$  in this diagram), which share a common observation at each time step (shaded nodes). Although the separate chains are independent a priori, the common observation induces an effective coupling among the observations. (Note that the observation nodes are linked by the moralization process that converts the directed graph into an undirected representation.) Thus, an equivalent model is shown in panel (b), where the dotted ellipses represent the induced coupling of each observation. A natural choice of approximating distribution in this



**Fig. 5.2.** Structured mean field approximation for a factorial HMM. (a) Original model consists of a set of hidden Markov models (defined on chains), coupled at each time by a common observation. (b) An equivalent model, where the ellipses represent interactions among all nodes at a fixed time, induced by the common observation. (c) Approximating distribution formed by a product of chain-structured models. Here  $\mu_\alpha$  and  $\mu_\delta$  are the sets of mean parameters associated with the indicated vertex and edge respectively.

case is based on the subgraph  $F$  consisting of the decoupled set of  $M$  chains, as illustrated in panel (c).

Now consider the nature of the quantities  $g_\beta(\mu(F))$ , which arise in the cost function (5.22). In this case, any function  $g_\beta$  will be defined on some subset of  $M$  nodes that are coupled at a given time slice (e.g., see ellipse in panel (c)). Note that this subset of nodes is independent with respect to the approximating distribution. Therefore, the function  $g_\beta(\mu(F))$  will decouple into a product of terms of the form  $f_i(\{\mu_i(F)\})$ , where each  $f_i$  is some function of the mean parameters  $\{\mu_i\} \equiv \{\mu_i(F)\}$

associated with node  $i = 1, \dots, M$  in the relevant cluster. For instance, if the factorial HMM involved binary variables and  $M = 3$  and  $\beta = (stu)$ , then  $g_{stu}(\mu) = \mu_s \mu_t \mu_u$ .

The decoupled nature of the approximation yields valuable savings on the computational side. In particular, the junction tree updates necessary to maintain consistency of the approximation can be performed by applying the forward-backward algorithm (i.e., the sum-product updates as an exact method) to each chain separately. This decoupling also has important consequences for the structure of any mean field fixed point. In particular, it can be seen that no term  $g_\beta(\mu(F))$  will ever depend on mean parameters associated with edges within any of the chains (e.g.,  $\mu_\delta$  in panel (c)). Otherwise stated, the partial derivative  $\frac{\partial g_\beta}{\partial \mu_\delta}$  is equal to 0 for all  $\beta \in \mathcal{I}(G) \setminus \mathcal{I}(F)$ . As an immediate consequence of these derivatives vanishing, the mean field canonical parameter  $\gamma_\delta(F)$  remains equal to  $\theta_\delta$  for all iterations of the updates (5.24). Any intermediate junction tree steps to maintain consistency will not affect  $\gamma_\delta(F)$  either. We conclude that it is, in fact, optimal to simply copy the edge potentials  $\theta_\delta$  from the original distribution onto each of the edges in the structured mean field approximation. In this particular form of structured mean field, only the single node potentials will be altered from their original setting. This conclusion is sensible, since the structured approximation (c) is a factorized approximation on a set of  $M$  chains, the internal structure of which is fully preserved in the approximation.

♣

In addition to structured mean field, there are various other extensions to naive mean field, which we mention only in passing here. Jaakkola and Jordan [118] explored the use of mixture distributions in improving the mean field approximation. A large class of techniques, including linear response theory and the TAP method [e.g., 191, 124, 152, 178, 250, 262], seek to improve the mean field approximation by introducing higher-order correction terms. Typically, the lower bound on the log partition function is not preserved by these higher-order methods. Leisnick and Kappen [153] proposed a class of higher-order expansions that generate tighter lower bounds.

## 5.4 Non-convexity of mean field

An important fact about the mean field approach is that the variational problem (5.22) may be non-convex, so that there may be local minima, and the mean field updates can have multiple solutions. The source of this non-convexity can be understood in two distinct ways, based on either the left-hand side or right-hand side of the formulation (5.22).

Consider first the representation of the mean field problem on the right-hand side of equation (5.22). The constraint set in this formulation—namely,  $\mathcal{M}(F)$ —is certainly convex. The cost function consists of a (concave) entropy term  $-A_F^*(\mu(F))$  and a set of terms  $\sum_{\alpha \in \mathcal{I}(F)} \theta_\alpha \mu_\alpha$  that are linear in  $\mu(F)$ . In contrast, the terms  $\sum_{\alpha \notin \mathcal{I}(F)} \theta_\alpha g_\alpha(\mu(F))$  involve the *nonlinear* functions  $g_\alpha$ , so that they may introduce non-convexity. An alternative and more geometric understanding of the non-convexity is provided by the left-hand side of equation (5.22). In this formulation, observe that the function to be maximized—namely,  $\langle \mu, \theta \rangle - A_F^*(\mu)$ —is always a concave function of  $\mu$ . Consequently, the source of any non-convexity must lie in the nature of the constraint set  $\mathcal{M}_F(G)$ .

To illustrate these two perspectives on non-convexity, let us return again to naive mean field for the Ising model.

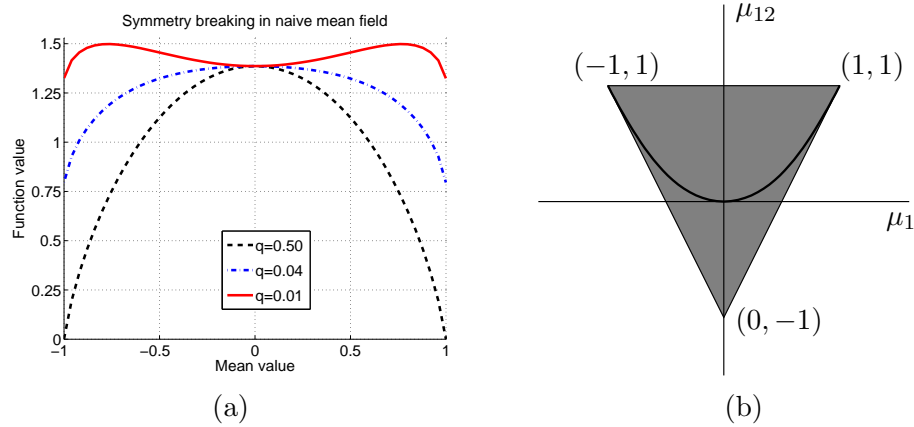
**Example 31 (Non-convexity for naive mean field).** We now consider an example, drawn from Jaakkola [117], that illustrates the non-convexity of naive mean field for a simple model. Consider a pair  $(X_1, X_2)$  of binary variates, taking values in the space<sup>2</sup>  $\{-1, +1\}^2$ , thereby defining a three-dimensional exponential family of the form  $p_\theta(x) \propto \exp(\theta_1 x_1 + \theta_2 x_2 + \theta_{12} x_1 x_2)$ , with associated mean parameters  $\mu_i = \mathbb{E}[X_i]$  and  $\mu_{12} = \mathbb{E}[X_1 X_2]$ . If the constraint  $\mu_{12} = \mu_1 \mu_2$  is imposed directly (as on the right-hand side on equation (5.22)), then the naive mean field objective function for this very special model takes the form

$$G(\mu_1, \mu_2; \theta) = \theta_{12} \mu_1 \mu_2 + \theta_1 \mu_1 + \theta_2 \mu_2 + H(\mu_1) + H(\mu_2), \quad (5.28)$$

where  $H(\mu_i) = -\frac{1}{2}(1 + \mu_i) \log \frac{1}{2}(1 + \mu_i) - \frac{1}{2}(1 - \mu_i) \log \frac{1}{2}(1 - \mu_i)$  are the

<sup>2</sup>This model, known as a “spin” representation, is a simple transformation of the  $\{0, 1\}^2$  state space considered earlier.

singleton entropies for the  $\{-1, +1\}$ -spin representation.



**Fig. 5.3.** Two different perspectives on the non-convexity of naive mean field for the Ising model. (a) Illustration of the naive mean field objective function (5.28) for three different parameter values:  $q \in \{0.50, 0.04, 0.01\}$ . For  $q = 0.50$  and  $q = 0.04$ , the global maximum is achieved at  $(\mu_1, \mu_2) = (0, 0)$ , whereas for  $q = 0.01$ , the point  $(0, 0)$  is no longer a maximum. Instead the global maximum is achieved at two non-symmetric points,  $(+\mu, -\mu)$  and  $(-\mu, +\mu)$ . (b) Non-convexity can also be seen by examining the shape of the set of fully factorized marginals for a pair of binary variables. The gray area shows the polytope defined by equations (5.29), corresponding to the intersection of  $\mathcal{M}(G)$  with the hyperplane  $\mu_1 = \mu_2$ . The (non-convex) quadratic curve  $\mu_{12} = \mu_1^2$  corresponds to the intersection of this projected polytope with the set  $\mathcal{M}_{F_0}(G)$  of fully factorized marginals.

Now let us consider a subfamily of such models, given by canonical parameters of the form

$$(\theta_1, \theta_2, \theta_{12}) = \left( 0, 0, \frac{1}{4} \log \frac{q}{1-q} \right) =: \theta(q),$$

where  $q \in (0, 1)$  is a parameter. By construction, this model is symmetric in  $X_1$  and  $X_2$ , so that for any value of  $q \in (0, 1)$ , we have  $\mathbb{E}[X_1] = \mathbb{E}[X_2] = 0$ . Moreover, some calculation shows that  $q = \mathbb{P}[X_1 = X_2]$ .

For  $q = 0.50$ , the objective function  $G(\mu_1, \mu_2; \theta(0.50))$  achieves its global maximum at  $(\mu_1, \mu_2) = (0, 0)$ , so that the mean field approximation is exact. (This exactness is to be expected since  $\theta(0.50) = (0, 0, 0)$ ,

corresponding a completely decoupled model). As  $q$  decreases away from 0.50, the objective function  $G$  starts to change, until for suitably small  $q$ , the point  $(\mu_1, \mu_2) = (0, 0)$  is no longer the global maximum—in fact, it is not even a local maximum.

To illustrate this behavior explicitly, we consider the cross-section of  $G$  obtained by setting  $\mu_1 = \tau$  and  $\mu_2 = -\tau$ , and then plot the one-dimensional function  $G(\tau, -\tau; \theta(q))$  for different values of  $q$ . As shown in Figure 5.3(a), for  $q = 0.50$ , this one-dimensional objective function has a unique global maximum at  $\tau = 0$ . As  $q$  decreases away from 0.50, the objective function gradually flattens out, as shown in the change between  $q = 0.50$  and  $q = 0.04$ . For  $q$  sufficiently close to zero, the point  $\tau = 0$  is no longer a maximum; instead, as shown in the curve for  $q = 0.01$ , there are two global maxima  $\pm\tau^*$  on either side of  $\tau = 0$ . Thus, for sufficiently small  $q$ , the maximum of the objective function (5.28) occurs at a pair  $\mu_1^* \neq \mu_2^*$ , even though the original model is always symmetric. This phenomenon, known in the physics literature as “spontaneous symmetry-breaking,” is a manifestation of non-convexity, since the optimum of any convex function will always respect symmetries in the underlying problem. Symmetry-breaking is not limited to this toy example, but also occurs with mean field methods applied to larger and more realistic graphical models, for which there may be a large number of competing modes in the objective function.

Alternatively, non-convexity in naive mean field can be understood in terms of the shape of the constraint set as an inner approximation to  $\mathcal{M}$ . For a pair of binary variates  $\{X_1, X_2\} \in \{-1, 1\}^2$ , the set  $\mathcal{M}$  is easily characterized: the mean parameters  $\mu_i = \mathbb{E}[X_i]$  and  $\mu_{12} = \mathbb{E}[X_1 X_2]$  are completely characterized by the four inequalities  $1 + ab\mu_{12} + a\mu_1 + b\mu_2 \geq 0$ , where  $\{a, b\} \in \{-1, 1\}^2$ . So as to facilitate visualization, consider a particular projection of this polytope—namely, that corresponding to intersection with the hyperplane  $\mu_1 = \mu_2$ . In this case, the four inequalities reduce to three simpler ones—namely:

$$\mu_{12} \leq 1, \quad \mu_{12} \geq 2\mu_1 - 1, \quad \mu_{12} \geq -2\mu_1 - 1. \quad (5.29)$$

Figure 5.3(b) shows the resulting two-dimensional polytope, shaded in gray. Now consider the intersection between this projected polytope and the set of factorized marginals  $\mathcal{M}_{F_0}(G)$ . The factorization condi-

tion imposes an additional constraint  $\mu_{12} = \mu_1^2$ , yielding a quadratic curve lying within the two-dimensional polytope described by the equations (5.29), as illustrated in Figure 5.3(b). Since this quadratic set is not convex, this establishes that  $\mathcal{M}_{F_0}(G)$  is not convex either. Indeed, if it were convex, then its intersection with any hyperplane would also be convex.



The geometric perspective on the set  $\mathcal{M}(G)$  and its inner approximation  $\mathcal{M}_F(G)$  reveals that more generally, mean field optimization is always non-convex for any exponential family in which the state space  $\mathcal{X}^m$  is finite. Indeed, for any such exponential family, the set  $\mathcal{M}(G)$  is a *finite convex hull*<sup>3</sup>

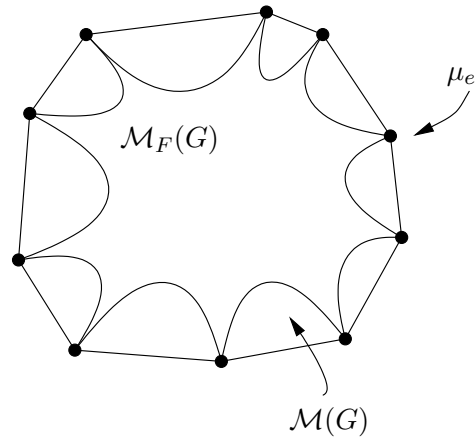
$$\mathcal{M}(G) = \text{conv} \{ \phi(e), e \in \mathcal{X}^m \} \quad (5.30)$$

in  $d$ -dimensional space, with extreme points of the form  $\mu_e := \phi(e)$  for some  $e \in \mathcal{X}^m$ . Figure 5.4 provides a highly idealized illustration of this polytope, and its relation to the mean field inner bound  $\mathcal{M}_F(G)$ .

We now claim that  $\mathcal{M}_F(G)$ —assuming that it is a strict subset of  $\mathcal{M}(G)$ —must be a non-convex set. To establish this claim, we first observe that  $\mathcal{M}_F(G)$  contains all of the extreme points  $\mu_x = \phi(x)$  of the polytope  $\mathcal{M}(G)$ . Indeed, the extreme point  $\mu_x$  is realized by the distribution that places all its mass on  $x$ , and such a distribution is Markov with respect to any graph. Therefore, if  $\mathcal{M}_F(G)$  were a convex set, then it would have to contain any convex combination of such extreme points. But from the representation (5.30), taking convex combinations of all such extreme points generates the full polytope  $\mathcal{M}(G)$ . Therefore, whenever  $\mathcal{M}_F(G)$  is a proper subset of  $\mathcal{M}(G)$ , it cannot be a convex set.

Consequently, non-convexity is an intrinsic property of mean field approximations. As suggested by Example 31, this non-convexity can have significant operational consequences, including symmetry-breaking, multiple local optima, and sensitivity to initialization.

<sup>3</sup>For instance, in the discrete case when the sufficient statistics  $\phi$  are defined by indicator functions in the standard overcomplete basis (3.34), we referred to  $\mathcal{M}(G)$  as a marginal polytope.



**Fig. 5.4.** Cartoon illustration of the set  $\mathcal{M}_F(G)$  of mean parameters that arise from tractable distributions is a non-convex inner bound on  $\mathcal{M}(G)$ . Illustrated here is the case of discrete random variables where  $\mathcal{M}(G)$  is a polytope. The circles correspond to mean parameters that arise from delta distributions, and belong to both  $\mathcal{M}(G)$  and  $\mathcal{M}_F(G)$ .

Nonetheless, mean-field methods have been used successfully in a variety of applications, and the lower bounding property of mean field methods is attractive, for instance in the context of parameter estimation, as we discuss at more length in the following section.



# 6

---

## Variational methods in parameter estimation

---

Our focus in the previous sections has been on the problems of computing or approximating the cumulant function  $A(\theta)$  and the mean parameters  $\mu = \mathbb{E}_\theta[\phi(X)]$ , assuming that the canonical parameter  $\theta$  specifying the density  $p_\theta$  was known. In this section, we turn to the inverse problem of estimating  $\theta$  on the basis of observed data. We consider the case in which the parameter  $\theta$  is assumed to be fixed but unknown (Sections 6.1 and 6.2), and the Bayesian perspective, in which the parameter  $\theta$  is viewed as a random variable (Section 6.3).

### 6.1 Estimation in fully observed models

The simplest case of parameter estimation corresponds to the case of *fully observed data*: a collection  $X_1^n := \{X^1, \dots, X^n\}$  of  $n$  independent and identically distributed (i.i.d.)  $m$ -vectors, each sampled according to  $p_\theta$ . Suppose that our goal is to estimate the unknown parameter  $\theta$ , which we view as a deterministic but non-random quantity for the moment. A classical approach to this estimation problem, dating back to Fisher, is via the principle of *maximum likelihood*: estimate  $\theta$  by maximizing the log likelihood of the data, given by  $\ell(\theta; X_1^n) := \frac{1}{n} \sum_{i=1}^n \log p_\theta(X^i)$ .

It is convenient and has no effect on the maximization to rescale the log likelihood by  $1/n$ , as we have done here.

For exponential families, the rescaled log likelihood takes a particularly simple and intuitive form. In particular, we define *empirical mean parameters*  $\hat{\mu} := \widehat{\mathbb{E}}[\phi(X)] = \frac{1}{n} \sum_{i=1}^n \phi(X^i)$  associated with the sample  $X_1^n$ . In terms of these quantities, the log likelihood for an exponential family can be written as

$$\ell(\theta; X_1^n) = \langle \theta, \hat{\mu} \rangle - A(\theta), \quad (6.1)$$

This expression highlights the close connection to maximum entropy and conjugate duality, as discussed in Section 3.6. The maximum likelihood estimate  $\hat{\theta}$  is a quantity maximizing this (random) objective function (6.1). As a corollary of Theorem 2, whenever  $\hat{\mu} \in \mathcal{M}^\circ$ , there exists a unique maximum likelihood solution. Indeed, by taking derivatives with respect to  $\theta$  of the objective function (6.1), the maximum likelihood estimate (MLE) is specified by the *moment matching conditions*  $\mathbb{E}_{\hat{\theta}}[\phi(X)] = \hat{\mu}$ . Moreover, by standard results on asymptotics of M-estimators [228], the MLE  $\hat{\theta}$  is consistent, in that it converges in probability to  $\theta$  as the sample size  $n$  tends to infinity.

### 6.1.1 Maximum likelihood for triangulated graphs

Our discussion up to this point has ignored the computational issue of solving the moment matching equation so as to obtain the MLE  $\hat{\theta}$ . As with other statistical computations in exponential families, the difficulty of solving this problem turns out to depend strongly on the graph structure. For triangulated<sup>1</sup> graphs, the MLE can be written as a closed-form function of the empirical marginals  $\hat{\mu}$ , which we describe here.

To illustrate the basic idea with a minimum of notational overhead, let us consider the simplest instance of a triangulated graph: namely, a tree  $T = (V, E)$ , with discrete random variables  $X_s$  taking values in  $\mathcal{X}_s := \{0, 1, \dots, r_s - 1\}$ . As previously discussed in Section 4, this collection of distributions can be represented as an exponential family (4.1), using indicator functions  $\mathbb{I}_{s;j}[x_s]$  for the event

<sup>1</sup>See Section 2.5.2 for further discussion of triangulated graphs.

$\{X_s = j\}$ , and  $\mathbb{I}_{st;jk}[x_s, x_t]$  for the event  $\{X_s = j, X_t = k\}$ . Moreover, the mean parameters in this exponential family are marginal probabilities—explicitly,  $\mu_{s;j} = \mathbb{P}_\theta[X_s = j]$  for  $s \in V$  and  $j \in \mathcal{X}_s$ , and  $\mu_{st;jk} = \mathbb{P}_\theta[X_s = j, X_t = k]$  for  $(s, t) \in E$  and  $(j, k) \in \mathcal{X}_s \times \mathcal{X}_t$ . Given an i.i.d. sample  $X_1^n := \{X^1, \dots, X^n\}$ , the empirical mean parameters correspond to the singleton and pairwise marginal probabilities

$$\hat{\mu}_{s;j} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{s;j}[X_s^i] \quad \text{and} \quad \hat{\mu}_{st;jk} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{st;jk}[X_s^i, X_t^i], \quad (6.2)$$

induced by the data.

With this set-up, we can exhibit the closed-form expression for the MLE in terms of the empirical marginals. For this particular exponential family, our assumption that  $\hat{\mu} \in \mathcal{M}^\circ$  means that the empirical marginals are all *strictly* positive. Consequently, the vector  $\hat{\theta}$  with elements

$$\hat{\theta}_{s;j} := \log \hat{\mu}_{s;j} \quad \forall \quad s \in V, j \in \mathcal{X}_s \quad \text{and} \quad (6.3a)$$

$$\hat{\theta}_{st;jk} := \log \frac{\hat{\mu}_{st;jk}}{\hat{\mu}_{s;j} \hat{\mu}_{t;k}} \quad \forall \quad (s, t) \in E, (j, k) \in \mathcal{X}_s \times \mathcal{X}_t \quad (6.3b)$$

is well defined.

We claim that the vector  $\hat{\theta}$  is the maximum likelihood estimate for the problem. In order to establish this claim, it suffices to show that the moment matching conditions  $\mathbb{E}_{\hat{p}_{\hat{\theta}}}[\phi(X)] = \hat{\mu}$  hold for the distribution  $p_{\hat{\theta}}$ . From the definition (6.3), this exponential family member has the form

$$\begin{aligned} p_{\hat{\theta}}(x) &= \exp \left\{ \sum_{s \in V} \hat{\theta}_s(x_s) + \sum_{(s,t) \in E} \hat{\theta}_{st}(x_s, x_t) \right\} \\ &= \prod_{s \in V} \hat{\mu}(x_s) \prod_{(s,t) \in E} \frac{\hat{\mu}_{st}(x_s, x_t)}{\hat{\mu}_s(x_s) \hat{\mu}_t(x_t)}, \end{aligned} \quad (6.4)$$

where we have used the fact (verifiable by some calculation) that  $A(\hat{\theta}) = 0$  in this parameterization. Moreover, the distribution  $p_{\hat{\theta}}$  has as its marginal distributions the empirical quantities  $\hat{\mu}_s$  and  $\hat{\mu}_{st}$ . This claim follows as a consequence of the junction tree framework (see Proposition 1); alternatively, they can be proved directly by an inductive “leaf-stripping” argument, based on a directed version of the

factorization (6.4). Consequently, the MLE has an explicit closed-form solution in terms of the empirical marginals for a tree. The same closed-form property holds more generally for any triangulated graph.

### 6.1.2 Iterative methods for computing MLEs

For non-triangulated graphical models, there is no longer a closed-form solution to the MLE problem, and iterative algorithms are needed. As an illustration, here we describe the iterative proportional fitting (IPF) algorithm [58, 57], a type of coordinate ascent method with particularly intuitive updates. To describe it, let us consider a generalization of the sufficient statistics  $\{\mathbb{I}_{s;j}(x_s)\}$  and  $\{\mathbb{I}_{st;jk}(x_s, x_t)\}$  used for trees, in which we define similar indicator functions over cliques of higher order. More specifically, for any clique  $C$  of size  $k$  in a given graph  $G$ , let us define a set of clique indicator functions

$$\mathbb{I}_J(x_C) = \prod_{s=1}^k \mathbb{I}_{s;j_s}[x_s],$$

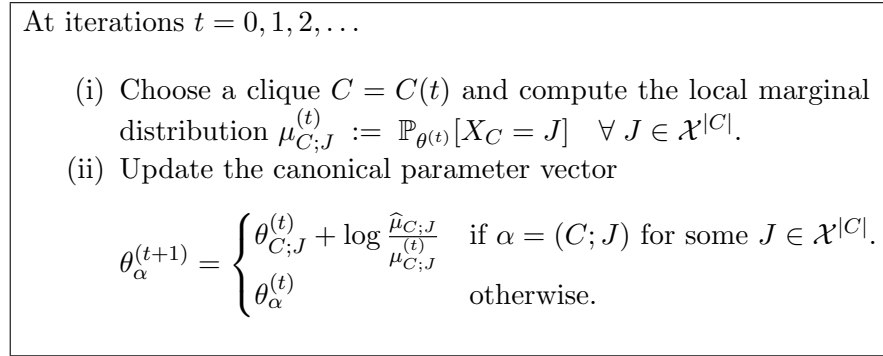
one for each configuration  $J = (j_1, \dots, j_k)$  over  $x_C$ . Each such sufficient statistic is associated with a canonical parameter  $\theta_{C;J}$ . As before, the associated mean parameters  $\mu_{C;J} = \mathbb{E}[\mathbb{I}_{C;J}(X_C)] = \mathbb{P}_\theta[X_C = J]$  are simply marginal probabilities, but now defined over cliques. Similarly, the data defines a set  $\hat{\mu}_{C;J} = \hat{\mathbb{P}}[X_C = J]$  of empirical marginal probabilities.

Since the MLE problem is differentiable and jointly concave in the vector  $\theta$ , coordinate ascent algorithms are guaranteed to converge to the global optimum. Using the cumulant generating properties of  $A$  (see Proposition 2), we can take derivatives with respect to some coordinate  $\theta_{C;J}$ , thereby obtaining

$$\frac{\partial \ell}{\partial \theta_{C;J}} = \hat{\mu}_{C;J} - \frac{\partial A}{\partial \theta_{C;J}}(\theta) = \hat{\mu}_{C;J} - \mu_{C;J}, \quad (6.5)$$

where  $\mu_{C;J} = \mathbb{E}_\theta[\mathbb{I}_{C;J}(X_C)]$  are the mean parameters defined by the current model. As a consequence, zero gradient points are specified by matching the model marginals  $\mathbb{P}_\theta[X_C = J]$  to the empirical marginals  $\hat{\mathbb{P}}[X_C = J]$ .

With this set-up, we may define a sequence of parameter estimates  $\{\theta^{(t)}\}$  via the recursion in Figure 6.1.



**Fig. 6.1:** Steps in the iterative proportional fitting (IPF) procedure.

These updates have two key properties: first, the log partition function is never changed by the updates. Indeed, let us compute the difference  $\Delta A := A(\theta^{(t+1)}) - A(\theta^{(t)})$ . We have

$$\begin{aligned} \Delta A &= \log \sum_x \exp \left\{ \langle \theta^{(t)}, \phi(x) \rangle - A(\theta^{(t)}) \right\} \exp \left\{ \sum_{J \in \mathcal{X}^{|C|}} \log \frac{\hat{\mu}_{C;J}}{\mu_{C;J}^{(t)}} \mathbb{I}_{C;J}[x_C] \right\} \\ &= \log \mathbb{E}_{\theta^{(t)}} \left[ \prod_J \left( \frac{\hat{\mu}_{C;J}}{\mu_{C;J}^{(t)}} \right)^{\mathbb{I}_{C;J}[x_C]} \right] \\ &= \log \sum_{J \in \mathcal{X}^{|C|}} \hat{\mu}_{C;J} = 0. \end{aligned}$$

Note that this invariance in  $A$  arises due to the overcompleteness of this particular exponential family—in particular, since  $\sum_J \mathbb{I}_{C;J}(x_C) = 1$  for all cliques  $C$ . Second, the update in step (ii) enforces exactly the zero-gradient condition (6.5). Indeed, after the parameter update  $\theta^{(t)} \rightarrow \theta^{(t+1)}$ , we have

$$\mathbb{E}_{\theta^{(t+1)}}[\mathbb{I}_{C;J}(X_C)] = \frac{\hat{\mu}_{C;J}}{\mu_{C;J}^{(t)}} \mathbb{E}_{\theta^{(t)}}[\mathbb{I}_{C;J}(X_C)] = \hat{\mu}_{C;J}.$$

Consequently, this IPF algorithm corresponds to as a *blockwise coordinate ascent* method for maximizing the objective function (6.1). At each iteration, the algorithm picks a particular block of coordinates—that is,  $\theta_{C;J}$ , for all configurations  $J \in \mathcal{X}^{|C|}$  defined on the block—and maximizes the objective function over this subset of coordinates. Consequently, standard theory on blockwise coordinate ascent methods [19] can be used to establish convergence of the IPF procedure. More generally, this iterative proportional fitting procedure is a special case of a broader class of iterative scaling, or successive projection algorithms; see papers [58, 56, 57] or the book [42] for further details on such algorithms and their properties.

From a computational perspective, however, there are still some open issues associated with the IPF algorithm in application to graphical models. Note that step (i) in the updates assume a “black box” routine that, given the current canonical parameter vector  $\theta^{(t)}$ , returns the associated vector of mean parameters  $\mu^{(t)}$ . For graphs of low treewidth, this computation can be carried out with the junction tree algorithm. For more general graphs not amenable to the junction tree method, one could imagine using approximate sampling methods or variational methods. The use of such approximate methods and their impact on parameter estimation is still an active area of research [220, 226, 238, 236, 248].

## 6.2 Partially observed models and expectation-maximization

A more challenging version of parameter estimation arises in the *partially observed* setting, in which the random vector  $X \sim p_\theta$  is not observed directly, but indirectly via a “noisy” version  $Y$  of  $X$ . This formulation includes the special case in which some subset where the remaining variates are unobserved—that is, “latent” or “hidden”. The expectation-maximization (EM) algorithm of Dempster et al. [65] provides a general approach to computing MLEs in this partially observed setting. Although the EM algorithm is often presented as an alternation between an expectation step (E step) and a maximization step (M step), it is also possible to take a variational perspective on EM, and view both steps as maximization steps [56, 176]. Such a perspective

illustrates how variational inference algorithms can be used in place of exact inference algorithms in the E step within the EM framework, and in particular, how the mean field approach is especially appropriate for this task.

A brief outline of our presentation in this section is as follows. We begin by deriving the EM algorithm in the exponential family setting, showing how the E step reduces to the computation of expected sufficient statistics—i.e., mean parameters. As we have seen, the variational framework provides a general class of methods for computing approximations of mean parameters. This observation suggests a general class of *variational EM algorithms*, in which the approximation provided by a variational inference algorithm is substituted for the mean parameters in the E step. In general, as a consequence of making such a substitution, one loses the exactness guarantees that are associated with the EM algorithm. In the specific case of mean field algorithms, however, one is still guaranteed that the method performs coordinate ascent on a *lower bound* of the likelihood function.

### 6.2.1 Exact EM algorithm in exponential families

We begin by deriving the exact EM algorithm for exponential families. Suppose that the set of random variables is partitioned into a vector  $Y$  of *observed* variables, and a vector  $X$  of *unobserved* variables, and the probability model is a joint exponential family distribution for  $(X, Y)$ :

$$p_{\theta}(x, y) = \exp\{\langle \theta, \phi(x, y) \rangle - A(\theta)\}. \quad (6.6)$$

Given an observation  $Y = y$ , we can also form the conditional distribution

$$\begin{aligned} p_{\theta}(x \mid y) &= \frac{\exp\{\langle \theta, \phi(x, y) \rangle\}}{\int_{\mathcal{X}^m} \exp\{\langle \theta, \phi(x, y) \rangle\} \nu(dx)} \\ &:= \exp\{\langle \theta, \phi(x, y) \rangle - A_y(\theta)\}. \end{aligned} \quad (6.7)$$

where for each *fixed*  $y$ , the log partition function  $A_y$  associated with this conditional distribution is given by the integral

$$A_y(\theta) := \log \int_{\mathcal{X}^m} \exp\{\langle \theta, \phi(x, y) \rangle\} \nu(dx). \quad (6.8)$$

Thus, we see that the conditional distribution  $p_\theta(x \mid y)$  belongs to an exponential family, with vector  $\phi(y, \cdot)$  of sufficient statistics. For each fixed  $y$ , the set  $\mathcal{M}_y$  of valid mean parameters in this exponential family takes the form

$$\mathcal{M}_y = \{\mu \mid \mu = \mathbb{E}_p[\phi(X, y)] \text{ for some } p\}, \quad (6.9)$$

where  $p$  is any density function over  $X$ , taken with respect to underlying base measure  $\nu$ . Consequently, applying Theorem 2 to this exponential family provides the variational representation

$$A_y(\theta) = \sup_{\mu_y \in \mathcal{M}_y} \{\langle \theta, \mu_y \rangle - A_y^*(\mu_y)\}, \quad (6.10)$$

where the conjugate dual is also defined variationally:

$$A_y^*(\mu_y) := \sup_{\theta \in \text{dom}(A_y)} \{\langle \mu_y, \theta \rangle - A_y(\theta)\}. \quad (6.11)$$

The maximum likelihood estimate  $\hat{\theta}$  is obtained by maximizing log probability of the observed data  $y$ , which is referred to as the *incomplete log likelihood* in the setting of EM. From equation (6.6) and the definition of  $A_y$ , it is easy to verify that this log likelihood can be written as a difference of log partition functions:

$$\log p_\theta(y) = A_y(\theta) - A(\theta). \quad (6.12)$$

From the variational representation (6.10), we obtain the lower bound  $A_y(\theta) \geq \langle \mu_y, \theta \rangle - A_y^*(\mu_y)$ , valid for any  $\mu_y \in \mathcal{M}_y$ , and hence a lower bound for the incomplete log likelihood:

$$\log p_\theta(y) \geq \langle \mu_y, \theta \rangle - A_y^*(\mu_y) - A(\theta) =: \mathcal{L}(\mu_y, \theta). \quad (6.13)$$

With this set-up, the EM algorithm is coordinate ascent on this function  $\mathcal{L}$  defining the lower bound (6.13):

$$\text{(E step)} \quad \mu_y^{(t+1)} = \arg \max_{\mu_y \in \mathcal{M}_y} \mathcal{L}(\mu_y, \theta^{(t)}) \quad (6.14a)$$

$$\text{(M step)} \quad \theta^{(t+1)} = \arg \max_{\theta \in \Omega} \mathcal{L}(\mu_y^{(t+1)}, \theta). \quad (6.14b)$$

To see the correspondence with the traditional presentation of the EM algorithm, note first that the maximization underlying the E step reduces to

$$\max_{\mu_y \in \mathcal{M}_y} \{\langle \mu_y, \theta^{(t)} \rangle - A_y^*(\mu_y)\}, \quad (6.15)$$



which by the variational representation (6.10) is equal to  $A_y(\theta^{(t)})$ , with the maximizing argument equal to the mean parameter that is dually coupled with  $\theta^{(t)}$ . Otherwise stated, the vector  $\mu_y^{(t+1)}$  that is computed by maximization in the first argument of  $\mathcal{L}(\mu_y, \theta)$  is exactly the expectation  $\mu_y^{(t+1)} = \mathbb{E}_{\theta^{(t)}}[\phi(X, y)]$ , a computation that is traditionally referred to as the E step, for obvious reasons. Moreover, the maximization underlying the M step reduces to

$$\max_{\theta \in \Omega} \{\langle \mu_y^{(t+1)}, \theta \rangle - A(\theta)\}, \quad (6.16)$$

which is simply a maximum likelihood problem based on the expected sufficient statistics  $\mu_y^{(t+1)}$ —traditionally referred to as the M step.

Moreover, given that the value achieved by the E step on the right-hand-side of (6.15) is equal to  $A_y(\theta^{(t)})$ , the inequality in (6.13) becomes an equality by (6.12). Thus, after the E step, the lower bound  $\mathcal{L}(\mu_y, \theta^{(t)})$  is actually equal to the incomplete log likelihood  $\log p(y; \theta^{(t)})$ , and the subsequent maximization of  $\mathcal{L}$  with respect to  $\theta$  in the M step is guaranteed to increase the log likelihood as well.

**Example 32 (EM for Gaussian mixtures).** We illustrate the EM algorithm with the classical example of the Gaussian mixture model, discussed earlier in Example 6. Let us focus on a scalar mixture model for simplicity, for which the observed vector  $Y$  is Gaussian mixture vector with  $r$  components, where the unobserved vector  $X \in \{1, \dots, r\}$  indexes the components. The complete likelihood of this mixture model can be written as

$$p_\theta(x, y) = \exp \left\{ \sum_{j=0}^{r-1} \mathbb{I}_j(x) [\alpha_j + \gamma_j y + \tilde{\gamma}_j y^2 - A_j(\gamma_j, \tilde{\gamma}_j)] - A(\alpha) \right\}, \quad (6.17)$$

where  $\theta = (\alpha, \gamma, \tilde{\gamma})$ , with the vector  $\alpha \in \mathbb{R}^r$  parameterizing the multinomial distribution over the hidden vector  $X$ , and the pair  $(\gamma_j, \tilde{\gamma}_j) \in \mathbb{R}^2$  parameterizing the Gaussian distribution of the  $j^{\text{th}}$  mixture component. The log normalization function  $A_j(\gamma_j, \tilde{\gamma}_j)$  is for the (conditionally) Gaussian distribution of  $Y$  given  $X = j$ , whereas the function  $A(\alpha) = \log[\sum_{j=0}^{r-1} \exp(\alpha_j)]$  normalizes the multinomial distribution. When the complete likelihood is viewed as an exponential family, the

sufficient statistics are the collection of triplets

$$\Psi_j(x, y) := \{\mathbb{I}_j(x), \mathbb{I}_j(x)y, \mathbb{I}_j(x)y^2\} \quad \text{for } j = 0, \dots, r-1.$$

Consider a collection of i.i.d. observations  $y^1, \dots, y^n$ . To each observation  $y^i$ , we associate a triplet  $(\mu^i, \eta^i, \tilde{\eta}^i) \in \mathbb{R}^r \times \mathbb{R}^r \times \mathbb{R}^r$ , corresponding to expectations of the triplet of sufficient statistics  $\Psi_j(X, y^i), j = 0, \dots, r-1$ . Since the conditional distribution has the form

$$p(x | y^i, \theta) \propto \exp \left\{ \sum_{j=0}^{r-1} \mathbb{I}_j(x) [\alpha_j + \gamma_j y^i + \tilde{\gamma}_j (y^i)^2 - A_j(\gamma_j, \tilde{\gamma}_j)] \right\},$$

some straightforward computation yields that the probability of the  $j^{\text{th}}$  mixture component—or equivalently, the mean parameter  $\mu_j^i = \mathbb{E}[\mathbb{I}_j(X)]$ —is given by

$$\begin{aligned} \mu_j^i &= \mathbb{E}[\mathbb{I}_j(X)] \\ &= \frac{\exp \{ \alpha_j + \gamma_j y^i + \tilde{\gamma}_j (y^i)^2 - A_j(\gamma_j, \tilde{\gamma}_j) \}}{\sum_{k=0}^{r-1} \exp \{ \alpha_k + \gamma_k y^i + \tilde{\gamma}_k (y^i)^2 - A_k(\gamma_k, \tilde{\gamma}_k) \}}. \end{aligned} \quad (6.18)$$

Similarly, the remaining two mean parameters are

$$\eta_j^i = \mathbb{E}[\mathbb{I}_j(X)y^i] = \mu_j^i y^i, \quad \text{and} \quad (6.19a)$$

$$\tilde{\eta}_j^i = \mathbb{E}[\mathbb{I}_j(X)y^i{}^2] = \mu_j^i (y^i)^2. \quad (6.19b)$$

The computation of these expectations (6.18) and (6.19) correspond to the E-step for this particular model.

The M-step requires solving the MLE optimization problem (6.16) over the triplet  $\theta = (\alpha, \gamma, \tilde{\gamma})$ , using the expected sufficient statistics computed in the E-step. Some computation shows that this maximization takes the form

$$\begin{aligned} \arg \max_{(\alpha, \gamma, \tilde{\gamma}) \in \Omega} & \left\{ \langle \alpha, \sum_{i=1}^n \mu^i \rangle + \langle \gamma, \sum_{i=1}^n \mu^i y^i \rangle \right. \\ & \left. + \langle \tilde{\gamma}, \sum_{i=1}^n \mu^i (y^i)^2 \rangle - \sum_{j=0}^{r-1} \sum_{i=1}^n \mu_j^i A_j(\gamma_j, \tilde{\gamma}_j) - nA(\alpha) \right\}. \end{aligned}$$

Consequently, the optimization decouples into separate maximization problems: one for the vector  $\alpha$  parameterizing the mixture component distribution, and one for each of the pairs  $(\gamma_j, \tilde{\gamma}_j)$  specifying the Gaussian mixture components. By the moment-matching property of maximum likelihood estimates, the optimum solution is to update  $\alpha$  such that

$$\mathbb{E}_\alpha[\mathbb{I}_j(X)] = \frac{1}{n} \sum_{i=1}^n (\mu_{y^i})_j \quad \text{for each } j = 0, \dots, r-1,$$

and to update the canonical parameters  $(\gamma_j, \tilde{\gamma}_j)$  specifying the  $j^{\text{th}}$  Gaussian mixture component  $(Y|X = j)$  such that the associated mean parameters—corresponding to the Gaussian mean and second moment—are matched to the data as follows:

$$\begin{aligned} \mathbb{E}_{\gamma_j, \tilde{\gamma}_j}[Y|X = j] &= \frac{\sum_{i=1}^n \mu_j^i y_j^i}{\sum_{i=1}^n \mu_j^i}, \quad \text{and} \\ \mathbb{E}_{\gamma_j, \tilde{\gamma}_j}[Y^2|X = j] &= \frac{\sum_{i=1}^n \mu_j^i (y_j^i)^2}{\sum_{i=1}^n \mu_j^i} \end{aligned}$$

where the expectations are taken under the exponential family distribution specified by  $(\gamma_j, \tilde{\gamma}_j)$ .



### 6.2.2 Variational EM

What if it is infeasible to compute the expected sufficient statistics? One possible response to this problem is to make use of a mean field approximation for the E step. In particular, given some class of tractable subgraphs  $F$ , recall the set  $\mathcal{M}_F(G)$ , as defined in equation (5.5), corresponding to those mean parameters that can be obtained by distributions that factor according to  $F$ . By using  $\mathcal{M}_F(G)$  as an inner approximation to the set  $\mathcal{M}_y$ , we can compute an approximate version of E step (6.15), of the form

$$\text{(Mean field E step)} \quad \max_{\mu \in \mathcal{M}_F(G)} \{ \langle \mu, \theta^{(t)} \rangle - A_y^*(\mu) \}, \quad (6.20)$$

This variational E step thus involves replacing the exact mean parameter  $\mathbb{E}_{\theta^{(t)}}[\phi(X, y)]$ , under the current model  $\theta^{(t)}$ , with the approximate

set of mean parameters computed by a mean field algorithm. Despite this (possibly crude) approximation, the resulting variational EM algorithm is a still coordinate ascent algorithm for  $\mathcal{L}$ . However, given that the E step no longer closes the gap between the auxiliary function  $\mathcal{L}$  and the incomplete log likelihood, it is no longer the case that the algorithm necessarily goes uphill in the latter quantity. Nonetheless, the variational mean field EM algorithm still has the attractive interpretation of maximizing a *lower bound* on the incomplete log likelihood.

In Section 4, we described a broad class of variational methods, including belief propagation and expectation propagation, that also generate approximations to mean parameters. It is tempting, and common in practice, to substitute the approximate mean parameters obtained from these relaxations in the place of the expected sufficient statistics in defining a “variational EM algorithm.” Such a substitution is particularly tempting given that these methods can yield better approximations to mean parameters than the mean field approach. Care must be taken in working with these algorithms, however, because the underlying relaxations do not (in general) guarantee lower bounds on the log partition function.<sup>2</sup> In the absence of guaranteed lower bounds, the connection to EM is thus less clear than in the mean field case; in particular, the algorithm is not guaranteed to maximize a lower bound on the incomplete log likelihood.

### 6.3 Variational Bayes

All of the variational approximations that we have discussed thus far in the paper, as well as those to be discussed in later sections, are applicable in principle to Bayesian inference problems, where the parameters are endowed with probability distributions and viewed as random variables. In the literature on the topic, however, the term “variational Bayes” has been reserved thus far for the application of the mean-field variational method to Bayesian inference [14, 88]. In this section we review this application and use the terminology of “variational Bayes” to refer to the application, but it should be borne in mind that this is

---

<sup>2</sup>However, see Sudderth et al. [219] for results on certain classes of attractive graphical models for which the Bethe approximation does provide such a lower bound.

but one potential application of variational methodology to Bayesian inference.

We consider a general set-up akin to that used in the preceding development of the EM algorithm. In particular, let the data be partitioned into an observed component  $Y$  and an unobserved component  $X$ , and assume that the complete data likelihood lies in some exponential family

$$p(x, y | \theta) = \exp \{ \langle \eta(\theta), \phi(x, y) \rangle - A(\eta(\theta)) \}. \quad (6.21)$$

The function  $\eta : \mathbb{R}^d \rightarrow \mathbb{R}^d$  provides some additional flexibility in the parameterization of the exponential family; this is convenient in the Bayesian setting.<sup>3</sup> We assume moreover that the prior distribution over  $\Theta$  also lies in some exponential family, of the *conjugate prior form*

$$p_{\xi, \lambda}(\theta) = \exp \{ \langle \xi, \eta(\theta) \rangle - \lambda A(\eta(\theta)) - B(\xi, \lambda) \}. \quad (6.22)$$

Note that this exponential family is specified by the sufficient statistics  $\{ \eta(\theta), -A(\eta(\theta)) \} \in \mathbb{R}^d \times \mathbb{R}$ , with associated canonical parameters  $(\xi, \lambda) \in \mathbb{R}^d \times \mathbb{R}$ . Its cumulant function  $B$  is defined in the usual way

$$B(\xi, \lambda) := \log \int \exp \{ \langle \xi, \eta(\theta) \rangle - \lambda A(\eta(\theta)) \} d\theta,$$

and expectations of the sufficient statistics  $\eta(\theta)$  and  $-A(\eta(\theta))$  can be obtained in the usual way by taking derivatives of  $B(\xi, \lambda)$  with respect to  $\xi$  and  $\lambda$ .

The class of models specified by the pair of equations (6.21) and (6.22) is broad. It includes as special cases Gaussian mixture models with Dirichlet priors, and the latent Dirichlet allocation model from Example 7.

We now consider a central problem in Bayesian inference, that of computing the *marginal likelihood*  $p_{\xi^*, \lambda^*}(y)$ , where  $y$  is an observed data point and where  $(\xi^*, \lambda^*)$  are fixed values of the hyperparameters.

<sup>3</sup>Up to this point, we have considered only exponential families in canonical form, for which  $[\eta(\theta)]_i = \theta_i$  for all  $i = 1, \dots, d$ . This is without loss of generality, because any regular, minimal exponential family can be reparameterized in such a canonical form [40]. However, it can be convenient to express exponential families with non-identity  $\eta$ .

This computation entails averaging over *both* the unobserved variates  $X$  and the random parameters  $\Theta$ . Working in the log domain, we have:

$$\begin{aligned} \log p_{\xi^*, \lambda^*}(y) &= \log \int \left[ \int p(x, y | \theta) dx \right] p_{\xi^*, \lambda^*}(\theta) d\theta \\ &= \log \int p_{\xi^*, \lambda^*}(\theta) p(y | \theta) d\theta. \end{aligned} \quad (6.23)$$

We now multiply and divide by  $p_{\xi, \lambda}(\theta)$  and use Jensen's inequality (see Proposition 7), to obtain the following lower bound, valid for any choice of  $(\xi, \lambda)$  in the domain of  $B$ :

$$\log p_{\xi^*, \lambda^*}(y) \geq \mathbb{E}_{\xi, \lambda}[\log p(y | \Theta)] + \mathbb{E}_{\xi, \lambda} \left[ \log \frac{p_{\xi^*, \lambda^*}(\Theta)}{p_{\xi, \lambda}(\Theta)} \right], \quad (6.24)$$

with equality for  $(\xi, \lambda) = (\xi^*, \lambda^*)$ . Here  $\mathbb{E}_{\xi, \lambda}$  denotes averaging over  $\Theta$  with respect to the distribution  $p_{\xi, \lambda}(\theta)$ .

From equation (6.12), we have  $\log p(y | \Theta) = A_y(\eta(\Theta)) - A(\eta(\Theta))$ , which when substituted into equation (6.24) yields

$$\log p_{\xi^*, \lambda^*}(y) \geq \mathbb{E}_{\xi, \lambda} [A_y(\eta(\Theta)) - A(\eta(\Theta))] + \mathbb{E}_{\xi, \lambda} \left[ \log \frac{p_{\xi^*, \lambda^*}(\Theta)}{p_{\xi, \lambda}(\Theta)} \right],$$

where  $A_y$  is the cumulant function of the conditional density  $p(x | y, \theta)$ .

For each fixed  $y$ , recall the set  $\mathcal{M}_y$  of mean parameters of the form  $\mu = \mathbb{E}[\phi(X, y)]$ . For any realization of  $\Theta$ , we could in principle apply the mean field lower bound (6.13) on the log likelihood  $A_y(\Theta)$  using a value  $\mu(\Theta) \in \mathcal{M}_y$  that varies with  $\Theta$ . We would thereby obtain that the marginal log likelihood  $\log p_{\xi^*, \lambda^*}(y)$  is lower bounded by

$$\mathbb{E}_{\xi, \lambda} [\langle \mu(\Theta), \eta(\Theta) \rangle - A_y^*(\mu(\Theta)) - A(\eta(\Theta))] + \mathbb{E}_{\xi, \lambda} \left[ \log \frac{p_{\xi^*, \lambda^*}(\Theta)}{p_{\xi, \lambda}(\Theta)} \right]. \quad (6.25)$$

At this stage, even after this sequence of lower bounds, if we were to optimize over the set of joint distributions on  $(X, \Theta)$ —or equivalently, over the choice  $\mu(\Theta)$  and the hyperparameters  $(\xi, \lambda)$  specifying the distribution of  $\Theta$ —the optimized lower bound (6.25) *would be equal* to the original marginal log likelihood  $\log p_{\xi^*, \lambda^*}(y)$ .

The variational Bayes algorithm is based on optimizing this lower bound using only product distributions over the pair  $(X, \Theta)$ . This is of-

ten described as “free-form optimization,” in that beyond the assumption of a product distribution the factors composing this product distribution are allowed to be arbitrary. But as we have seen in Section 3.1, the maximum entropy principle underlying the variational framework yields solutions that necessarily have exponential form. Thus we can exploit conjugate duality to obtain compact forms for these solutions, working with dually-coupled sets of parameters for the distributions over both  $X$  and  $\Theta$ .

We now derive the variational Bayes algorithm as a coordinate-ascent method for solving the mean field variational problem over product distributions. We begin by reformulating the objective function into a form that allows us to exploit conjugate duality. Since  $\mu$  is independent of  $\Theta$ , the optimization problem (6.25) can be simplified to

$$\langle \mu, \bar{\eta} \rangle - A_y^*(\mu) - \bar{A} + \mathbb{E}_{\xi, \lambda} \left[ \log \frac{p_{\xi^*, \lambda^*}(\Theta)}{p_{\xi, \lambda}(\Theta)} \right], \quad (6.26)$$

where  $\bar{\eta} := \mathbb{E}_{\xi, \lambda}[\eta(\Theta)]$  and  $\bar{A} := \mathbb{E}_{\xi, \lambda}[A(\Theta)]$ . Using the exponential form (6.22) of  $p_{\xi, \lambda}$ , we have

$$\mathbb{E}_{\xi, \lambda} \left[ \log \frac{p_{\xi^*, \lambda^*}(\Theta)}{p_{\xi, \lambda}(\Theta)} \right] = \langle \bar{\eta}, \xi^* - \xi \rangle + \langle -\bar{A}, \lambda^* - \lambda \rangle - B(\xi^*, \lambda^*) + B(\xi, \lambda). \quad (6.27)$$

By the conjugate duality between  $(B, B^*)$ , and since  $(\bar{\eta}, \bar{A})$  are dually coupled with  $(\xi, \lambda)$  by construction, we have

$$B^*(\bar{\eta}, \bar{A}) = \langle \bar{\eta}, \xi \rangle + \langle -\bar{A}, \lambda \rangle - B(\xi, \lambda). \quad (6.28)$$

By combining equations (6.26) through (6.28) and performing some algebra we find that the decoupled optimization problem is equivalent to maximizing the function

$$\langle \mu + \xi^*, \bar{\eta} \rangle - A_y^*(\mu) + \langle \lambda^* + 1, -\bar{A} \rangle - B^*(\bar{\eta}, \bar{A}) \quad (6.29)$$

over  $\mu \in \mathcal{M}_y$  and  $(\bar{\eta}, \bar{A}) \in \text{dom}(B)$ . Performing coordinate ascent on this function amounts to first maximizing first over  $\mu$ , and then maximizing over the mean parameters<sup>4</sup>  $(\bar{\eta}, \bar{A})$ . Doing so generates a sequence

<sup>4</sup>Equivalently, by the one-to-one correspondence between exponential and mean parameters (Theorem 1), this step corresponds to maximization over the associated canonical parameters  $(\xi, \lambda)$ .

of iterates  $(\mu^{(t)}, \bar{\eta}^{(t)}, \bar{A}^{(t)})$ , which are updated as follows

$$\mu^{(t+1)} = \arg \max_{\mu \in \mathcal{M}_y} \{ \langle \mu, \bar{\eta}^{(t)} \rangle - A_y^*(\mu) \}, \quad \text{and} \quad (6.30)$$

$$(\bar{\eta}^{(t+1)}, \bar{A}^{(t+1)}) = \arg \max_{(\bar{\eta}, \bar{A})} \{ \langle \mu^{(t+1)} + \xi^*, \bar{\eta} \rangle - (1 + \lambda^*) \bar{A} - B^*(\bar{\eta}, \bar{A}) \}. \quad (6.31)$$

In fact, both of these coordinate-wise optimizations have explicit solutions. The explicit form of equation (6.30) is given by

$$\mu^{(t+1)} = \mathbb{E}_{\bar{\eta}^{(t)}}[\phi(X, y)], \quad (6.32)$$

as can be verified by taking derivatives in equation (6.30), and using the cumulant properties of  $A_y^*$  (see Proposition 2). By analogy to the EM algorithm, this step is referred to as the “VB-E step.” Similarly, in the “VB-M step” we first update the hyperparameters  $(\xi, \lambda)$  as

$$(\xi^{(t+1)}, \lambda^{(t+1)}) = (\xi^* + \mu^{(t+1)}, \lambda^* + 1),$$

and then use these updated hyperparameters to compute the new mean parameters  $\bar{\eta}$ :

$$\bar{\eta}^{(t+1)} = \mathbb{E}_{(\xi^{(t+1)}, \lambda^{(t+1)})}[\eta(\Theta)]. \quad (6.33)$$

In summary, the variational Bayes algorithm is an iterative algorithm that alternates between two steps. The VB-E step (6.32) entails computing the conditional expectation of the sufficient statistics given the data

$$\text{(VB-E step)} \quad \mu^{(t+1)} = \int \phi(x, y) p(x | y, \bar{\eta}^{(t)}) dx,$$

with the conditional distribution specified by the current parameter estimate  $\bar{\eta}^{(t)}$ . The VB-M step (6.33) entails updating the hyperparameters to incorporate the data, and then recomputing the averaged parameter

$$\text{(VB-M step)} \quad \bar{\eta}^{(t+1)} = \int \eta(\theta) p_{\xi^{(t+1)}, \lambda^{(t+1)}}(\theta) d\theta.$$

The ordinary variational EM algorithm can be viewed as a “degenerate” form of these updates, in which the prior distribution over  $\Theta$  is always a delta function at a single point estimate of  $\theta$ .



# 7

---

## Convex relaxations and upper bounds

---

Up to this point, we have considered two broad classes of approximate variational methods: Bethe approximations and their extensions (Section 4), and mean field methods (Section 5). Mean field methods provide not only approximate mean parameters but also lower bounds on the log partition function. In contrast, the Bethe method and its extensions lead to neither upper or lower bounds on the log partition function. We have seen that both classes of methods are, at least in general, based on *non-convex* variational problems. For mean field methods, the non-convexity stems from the nature of the inner approximation to the set of mean parameters, as illustrated in Figure 5.4. For the Bethe-type approaches, the lack of convexity in the objective function—in particular, the entropy approximation—is the source. Regardless of the underlying cause, such non-convexity has a number of undesirable consequences, including multiple optima and sensitivity to the problem parameters (for the variational problem itself), and convergence issues and dependence on initialization (for iterative algorithms used to solve the variational problem).

Given that the underlying exact variational principle (3.45) is certainly convex, it is natural to consider variational approximations that

retain this convexity. In this section, we describe a broad class of such convex variational approximations, based on approximating the set  $\mathcal{M}$  with a *convex set*, and replacing the dual function  $A^*$  with a *convex function*. These two steps lead to a new variational principle, which represents a computationally tractable alternative to the original problem. Since this approximate variational principle is based on maximizing a concave and continuous function over a convex and compact set, the global optimum is unique, and is achieved. If in addition, the negative dual function  $-A^*$ , corresponding to the entropy function, is replaced by an upper bound, and the set  $\mathcal{M}$  is approximated by a convex outer bound, then the global optimum of the approximate variational principle also provides an upper bound on the log partition function. These upper bounds are complementary to the lower bounds provided by mean field procedures (see Section 5).

We begin by describing a general class of methods based on approximating the entropy by a convex combination of easily computable entropies, thereby obtaining a tractable relaxation that is guaranteed to be convex. As our first example illustrates, following this procedure using tree-structured distributions as the tractable class leads to the family of “convexified” Bethe variational problems, and associated reweighted sum-product algorithms [238, 241]. However, this basic procedure for “convexification” is quite broadly applicable; as we describe, it yields convex analogs of other known variational methods, including Kikuchi and region-graph methods [241, 255, 246], as well as expectation-propagation approximation. It has also suggested novel variational methods, also based on the notion of convex combinations, including those based on planar graph decomposition [91]. Finally, there are other convex variational relaxations that are not directly based on convex combinations of tractable distributions, including the method of conditional entropy decompositions [92], and methods based on semidefinite constraints and log-determinant programming [243].

We conclude the section by discussing some benefits of convexity in variational methods—in addition to the obvious one of unique optima—which include guarantees of algorithmic stability, as well as their potential use as surrogate likelihoods in parameter estimation.

### 7.1 Generic convex combinations and convex surrogates

We begin with a generic description of approximations based on convex combinations of tractable distributions. Consider a  $d$ -dimensional exponential family, with associated canonical parameters  $\theta = \{\theta_\alpha \mid \alpha \in \mathcal{I}\}$ , and associated sufficient statistics  $\phi = \{\phi_\alpha \mid \alpha \in \mathcal{I}\}$ . Although computing mean parameters for an arbitrary  $\theta \in \Omega$  might be intractable, it is frequently the case that such computations are tractable for certain special choices of the canonical parameters. If the exponential family is viewed as a Markov random field defined by some underlying graph  $G$ , many such special choices can be indexed by particular subgraphs  $F$ , say belonging to some class  $\mathfrak{D}$ . Examples that we consider below include  $\mathfrak{D}$  corresponding to the set of all spanning trees of  $G$ , or planar subgraphs of  $G$ . In terms of the exponential family, we can think of the subgraph  $F$  as extracting a subset of indices  $\mathcal{I}(F)$  from the full index set  $\mathcal{I}$  of potential functions, thereby defining a lower-dimensional exponential family based on  $d(F) < d$  sufficient statistics. We let  $\mathcal{M}(F)$  denote the lower-dimensional set of mean parameters associated with this exponential family; concretely, this set has the form

$$\mathcal{M}(F) := \{\mu \in \mathbb{R}^{|\mathcal{I}(F)|} \mid \exists p \text{ s.t. } \mu_\alpha = \mathbb{E}_p[\phi(X)] \ \forall \alpha \in \mathcal{I}(F)\}. \quad (7.1)$$

For any mean parameter  $\mu \in \mathcal{M}$ , let  $\mu \mapsto \mu(F)$  represent the coordinate projection, mapping from the full space  $\mathcal{I}$  to the subset  $\mathcal{I}(F)$  of indices associated with  $F$ . By definition of the sets  $\mathcal{M} \subseteq \mathbb{R}^d$  and  $\mathcal{M}(F) \subseteq \mathbb{R}^{d(F)}$ , we are guaranteed that the projected vector  $\mu(F)$  is an element of  $\mathcal{M}(F)$ . Let  $-A^*(\mu(F))$  be the negative dual function (entropy) defined by the projected mean parameter  $\mu(F)$ . A simple consequence of Theorem 2 is that this entropy is always an upper bound on the negative dual function  $A^*(\mu)$  defined in the original  $d$ -dimensional exponential family:

**Proposition 8 (Maximum entropy bounds).** *Given any mean parameter  $\mu \in \mathcal{M}$  and its projection  $\mu(F)$  onto any subgraph  $F$ , we have the bound*

$$A^*(\mu(F)) \leq A^*(\mu). \quad (7.2)$$

From Theorem 2, the negative dual value  $-A^*(\mu)$  corresponds to the entropy of the exponential family member with mean parameters  $\mu$ , with a similar interpretation for  $A^*(\mu(F))$ . Alternatively phrased, Proposition 8 states the the entropy of the exponential family member  $p_{\mu(F)}$  is always larger than the entropy of the distribution  $p_\mu$ . Intuitively, the the distribution  $p_\mu$  is obtained from a maximum entropy problem with more constraints—corresponding to the indices  $\alpha \in \mathcal{I} \setminus \mathcal{I}(F)$ —and so has lower entropy. Our proof makes this intuition precise.

*Proof.* From Theorem 2, the dual function is realized as the solution of the optimization problem

$$A^*(\mu) = \sup_{\theta \in \mathbb{R}^d} \{\langle \mu, \theta \rangle - A(\theta)\}. \quad (7.3)$$

Since the exponential family defined by  $F$  involves only a subset  $\mathcal{I}(F)$  of the potential functions, the associated dual function  $A^*(\mu(F))$  is realized by the lower-dimensional optimization problem

$$A^*(\mu(F)) = \sup_{\theta(F) \in \mathbb{R}^{d(F)}} \{\langle \mu(F), \theta(F) \rangle - A(\theta(F))\}.$$

But this optimization problem can be recast as a restricted version of the original problem (7.3):

$$A^*(\mu(F)) = \sup_{\substack{\theta \in \mathbb{R}^d, \\ \theta_\alpha = 0 \quad \forall \alpha \notin \mathcal{I}(F)}} \{\langle \mu, \theta \rangle - A(\theta)\},$$

from which the claim follows.  $\square$

From here onwards, with a slight abuse of notation, we use  $H(\mu)$  to denote the negative dual function  $-A^*(\mu)$ , and similarly  $H(\mu(F))$  to denote the subgraph-structured entropy  $-A^*(\mu(F))$ . With this notation, the upper bound (7.2) can be written as

$$H(\mu) \leq H(\mu(F)). \quad (7.4)$$

Given this bound for each subgraph  $F$ , any convex combination of subgraph-structured entropies is also an upper bound on the original

entropy. In particular, if we consider a probability distribution over the set of tractable graphs, meaning a probability vector  $\rho \in \mathbb{R}^{|\mathfrak{D}|}$  such that

$$\rho(F) \geq 0 \text{ for all tractable } F \in \mathfrak{D}, \text{ and } \sum_F \rho(F) = 1. \quad (7.5)$$

Any such distribution generates the upper bound

$$H(\mu) \leq \mathbb{E}_\rho[H(\mu(F))] := \sum_{F \in \mathfrak{D}} \rho(F) H(\mu(F)). \quad (7.6)$$

This upper bound on the entropy may be combined with any outer approximation on the set of mean parameters  $\mathcal{M}$  for which each of the subgraph-structured entropies  $H(\mu(F)) = -A^*(\mu(F))$  are well defined. The simplest outer approximation with this property is defined by requiring that each projected mean parameter  $\mu(F)$  belong to the projected set  $\mathcal{M}(F)$ . In particular, let us define the constraint set

$$\mathcal{L}(G; \mathfrak{D}) := \{ \mu \in \mathbb{R}^d \mid \mu(F) \in \mathcal{M}(F) \ \forall F \in \mathfrak{D} \}, \quad (7.7)$$

We observe that  $\mathcal{L}(G; \mathfrak{D})$  is an outer bound on  $\mathcal{M}(G)$ , since by definition of the sets  $\mathcal{M}(F)$ , any mean parameter  $\mu \in \mathcal{M}(G)$ , when projected down to  $\mu(F)$ , is also globally realizable. Moreover,  $\mathcal{L}(G; \mathfrak{D})$  is a convex set, since each of the subsets  $\mathcal{M}(F)$  in its definition are convex.

Overall, these two ingredients—the convex upper bound (7.6) and the convex outer bound (7.7)—yield the following approximate variational principle

$$B_{\mathfrak{D}}(\theta; \rho) := \sup_{\mu \in \mathcal{L}(G; \mathfrak{D})} \left\{ \langle \mu, \theta \rangle + \sum_{F \in \mathfrak{D}} \rho(F) H(\mu(F)) \right\}. \quad (7.8)$$

Note that the objective function defining  $B_{\mathfrak{D}}$  is concave; moreover, the constraint set  $\cap_F \mathcal{M}_F$  is convex, since it is the intersection of the sets  $\mathcal{M}_F$ , each of which are individually convex sets. Overall, then, the function  $B_{\mathfrak{D}}$  is a new convex function that approximates the original cumulant function  $A$ . We refer to such a function as a *convex surrogate* to  $A$ .

Of course, in order for variational principles of the form (7.8) to be useful, it must be possible to evaluate and optimize the objective function. Accordingly, in the following sections, we provide some specific choices of the tractable class  $\mathfrak{D}$  that lead to interesting and practical algorithms.

## 7.2 Variational methods from convex relaxations

In this section, we describe a number of versions of the convex surrogate (7.8), including convex versions of the Bethe/Kikuchi problems, as well as convex versions of expectation-propagation and other moment-matching methods.

### 7.2.1 Tree-reweighted sum-product and Bethe

We begin by deriving the tree-reweighted sum-product algorithm and the tree-reweighted Bethe variational principle [238, 241], corresponding to a “convexified” analog of the ordinary Bethe variational principle described in Section 4. Given an undirected graph  $G = (V, E)$ , consider a pairwise Markov random field

$$p_\theta(x) \propto \exp \left\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \right\}, \quad (7.9)$$

where we are using the standard overcomplete representation based on indicator functions at single nodes and edges (see equation (3.34)). Let the tractable class  $\mathfrak{D}$  be the set  $\mathfrak{T}$  of all spanning trees  $T = (V, E(T))$  of the graph  $G$ . (A spanning tree of a graph is a tree-structured subgraph whose vertex set covers the original graph.) Letting  $\rho$  denote a probability distribution over the set of spanning trees, equation (7.6) yields the upper bound  $H(\mu) \leq \sum_T \rho(T) H(\mu(T))$ , based on a convex combination of tree-structured entropies.

An attractive property of tree-structured entropies is that they decompose additively in terms of entropies associated with the vertices and edges of the tree. More specifically, these entropies are defined by the mean parameters, which (in the canonical overcomplete representation (3.34)) correspond to singleton marginal distributions  $\mu_s(\cdot)$  defined at each vertex  $s \in V$ , and a joint pairwise marginal distribution  $\mu_{st}(\cdot, \cdot)$  defined for each edge  $(s, t) \in E(T)$ . As discussed earlier (e.g., in Section 4), the factorization (4.8) of any tree-structured probability distribution yields the entropy decomposition

$$H(\mu(T)) = \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E(T)} I_{st}(\mu_{st}). \quad (7.10)$$

Now consider the averaged form of the bound (7.6). Since the trees are all spanning, the entropy term  $H_s$  for node  $s \in V$  receives a weight of one in this average. On the other hand, the mutual information term  $I_{st}$  for edge  $(s, t)$  receives the weight  $\rho_{st} = \mathbb{E}_\rho[\mathbb{I}[(s, t) \in E(T)]]$ , where  $\mathbb{I}[(s, t) \in E(T)]$  is an indicator function for the event that edge  $(s, t)$  is included in the edge set  $E(T)$  of a given tree  $T$ . Overall, we obtain the following upper bound on the exact entropy:

$$H(\mu) \leq \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E} \rho_{st} I_{st}(\mu_{st}). \quad (7.11)$$

We refer to the edge weight  $\rho_{st}$  as the *edge appearance probability* of edge  $(s, t)$ , since it reflects the probability mass associated with edge  $(s, t)$ . We discuss these edge appearances at more length below (following Theorem 4).

Let us now consider the form of the outer bound  $\mathcal{L}(G; \mathfrak{T})$  on the set  $\mathcal{M}$ . For the pairwise MRF with the overcomplete parameterization under consideration, the set  $\mathcal{M}$  is simply the marginal polytope  $\mathbb{M}(G)$ . On the other hand, the set  $\mathbb{M}(T)$  is simply the marginal polytope for the tree  $T$ , which from our earlier development (see Proposition 4), is equivalent to  $\mathbb{L}(T)$ . Consequently, the constraint  $\mu(T) \in \mathbb{M}(T)$  is equivalent to enforcing non-negativity constraints, normalization (at each vertex) and marginalization (across each edge) of the tree. Enforcing the inclusion  $\mu(T) \in \mathbb{M}(T)$  for *all* trees  $T \in \mathfrak{T}$  is equivalent to enforcing the marginalization on *every* edge of the full graph  $G$ . We conclude that in this particular case, the set  $\mathcal{L}(G; \mathfrak{T})$  is equivalent to the set  $\mathbb{L}(G)$  of locally consistent pseudomarginals, as defined earlier (4.7).

Overall, then, we obtain a variational problem that can be viewed as a “convexified” form of the Bethe variational problem. We summarize our findings in the following result [238, 241]:

**Theorem 4** (Tree-reweighted Bethe and sum-product). *(a) For any choice of edge appearance vector  $\rho_e$  in the spanning tree polytope  $\mathcal{S}(G)$ , the cumulant function  $A(\theta)$  at evaluated at  $\theta$  is upper bounded by the solution of the tree-reweighted Bethe variational problem (BVP):*

$$B_{\mathfrak{T}}(\theta; \rho_e) := \max_{\tau \in \mathbb{L}(G)} \left\{ \langle \tau, \theta \rangle + \sum_{s \in V} H_s(\tau_s) - \sum_{(s,t) \in E} \rho_{st} I_{st}(\tau_{st}) \right\}. \quad (7.12)$$

This problem is strictly convex with a unique optimum for any valid vector of edge appearance probabilities with  $\rho_e > 0$  for all edges  $e$ .

(b) The tree-reweighted BVP can be solved using the tree-reweighted sum-product updates

$$M_{ts}(x_s) \leftarrow \kappa \sum_{x'_t \in \mathcal{X}_t} \exp\left(\frac{1}{\rho_{st}}\theta_{st}(x_s, x'_t) + \theta_t(x'_t)\right) \left\{ \frac{\prod_{v \in N(t) \setminus s} [M_{vt}(x'_t)]^{\rho_{vt}}}{[M_{st}(x'_t)]^{(1-\rho_{ts})}} \right\}, \quad (7.13)$$

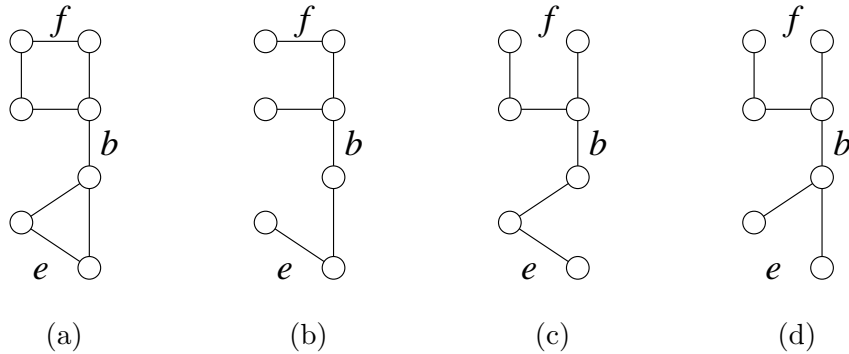
which has a unique fixed point under the conditions in (a).

We make a few comments on Theorem 4, before providing the proof.

**Valid edge weights:** Observe that the tree-reweighted BVP (7.12) is closely related to the ordinary Bethe problem (see equation (4.16)). In particular, if we set  $\rho_{st} = 1$  for all edges  $(s, t) \in E$ , then the two formulations are equivalent. However, the condition  $\rho_{st} = 1$  implies that every edge appears in every spanning tree of the graph with probability one, which can happen if and only if the graph is actually tree-structured. More generally, the set of valid edge appearance vectors  $\rho_e$  must belong to the so-called *spanning tree polytope* [70, 51] associated with  $G$ , which we denote by  $\mathcal{S}(G)$ . Note that these edge appearance probabilities must satisfy various constraints, depending on the structure of the graph. A simple example should help to provide intuition.

**Example 33 (Edge appearance probabilities).** Figure 7.1(a) shows a graph, and panels (b) through (d) show three of its spanning trees  $\{T^1, T^2, T^3\}$ . Suppose that we form a uniform distribution  $\rho$  over these trees by assigning probability  $\rho(T^i) = 1/3$  to each  $T^i$ ,  $i = 1, 2, 3$ . Consider the edge with label  $f$ ; notice that it appears in  $T^1$ , but in neither of  $T^2$  and  $T^3$ . Therefore, under the uniform distribution  $\rho$ , the associated edge appearance probability is  $\rho_f = 1/3$ . Since edge  $e$  appears in two of the three spanning trees, similar reasoning establishes that  $\rho_e = 2/3$ . Finally, observe that edge  $b$  appears in any spanning tree (i.e., it is a bridge), so that it must have edge appearance probability  $\rho_b = 1$ . ♣





**Fig. 7.1.** Illustration of valid edge appearance probabilities. Original graph is shown in panel (a). Probability  $1/3$  is assigned to each of the three spanning trees  $\{T_i \mid i = 1, 2, 3\}$  shown in panels (b)–(d). Edge  $b$  appears in all three trees so that  $\rho_b = 1$ . Edges  $e$  and  $f$  appear in two and one of the spanning trees respectively, which gives rise to edge appearance probabilities  $\rho_e = 2/3$  and  $\rho_f = 1/3$ .

In their work on fractional belief propagation, Wiegerinck and Heskes [256] examined the class of reweighted Bethe problems of the form (7.12), but without the requirement that the weights  $\rho_{st}$  belong to the spanning tree polytope  $\mathcal{S}(G)$ . Although loosening this requirement does yield a richer family of variational problems, in general one loses the guarantee of convexity, and (hence) that of a unique global optimum. On the other hand, Weiss et al. [246] have pointed out that other choices of weights  $\rho_{st}$ , not necessarily in the spanning tree polytope, can also lead to convex variational problems. In general, convexity and the upper bounding property are not equivalent. For instance, for any single cycle graph, setting  $\rho_{st} = 1$  for all edges (i.e., the ordinary BVP choice) yields a convex variational problem [246], but the value of the Bethe variational problem does not upper bound the cumulant function value. Various other researchers [107, 163, 184, 185] also discuss the choice edge/cliq weights in Bethe/Kikuchi approximations, and its consequences for convexity.

**Properties of tree-reweighted sum-product:** In analogy to the ordinary Bethe problem and sum-product algorithm, the fixed point of

tree-reweighted sum-product (TRW) message-passing (7.13) specifies the optimal solution of the variational problem (7.12) as follows:

$$\begin{aligned} \tau_s^*(x_s) &= \kappa \exp\{\theta_s(x_s)\} \prod_{v \in N(s)} [M_{vs}^*(x_s)]^{\rho_{vs}} & (7.14a) \\ \tau_{st}^*(x_s, x_t) &= \kappa \varphi_{st}(x_s, x_t; \theta) \frac{\prod_{v \in N(s) \setminus t} [M_{vs}^*(x_s)]^{\rho_{vs}}}{[M_{ts}^*(x_s)]^{(1-\rho_{st})}} \frac{\prod_{v \in N(t) \setminus s} [M_{vt}^*(x_t)]^{\rho_{vt}}}{[M_{st}^*(x_t)]^{(1-\rho_{ts})}} & (7.14b) \end{aligned}$$

where  $\varphi_{st}(x_s, x_t; \theta) := \exp\{\frac{1}{\rho_{st}}\theta_{st}(x_s, x_t) + \theta_s(x_s) + \theta_t(x_t)\}$ . In contrast to the ordinary sum-product algorithm, the fixed point (and associated optimum  $(\tau_s^*, \tau_{st}^*)$ ) is unique for any valid vector of edge appearances. Roosta et al. [199] provide sufficient conditions for convergence, based on contraction arguments such as those used for ordinary sum-product [225, 87, 175, 115]. In practical terms, the updates (7.13) appear to always converge if damped forms of the updates are used (i.e., setting  $\log M^{new} = (1 - \lambda) \log M^{old} + \lambda \log M$ , where  $M^{old}$  is the previous vector of messages, and  $\lambda \in (0, 1]$  is the damping parameter). As an alternative, Globerson and Jaakkola [93] proposed a related message-passing algorithm based on oriented trees that is guaranteed to converge, but appears to do so more slowly than damped TRW-message passing. Another possibility would be to adapt other double-loop algorithms [265, 249, 108, 107], originally developed for the ordinary Bethe/Kikuchi problems, to solve these convex minimization problems; see Hazan and Shashua [105] for some recent work along these lines.

**Optimal choice of edge appearance probabilities:** Note that equation (7.12) actually describes a family of variational problems, one for each choice of probability distribution over the set of spanning trees  $\mathfrak{T}$ . A natural question thus arises: what is the “optimal” choice of edge appearance? One notion of optimality is the distribution that leads to the tightest upper bound on the cumulant function (i.e., for which the gap  $B_{\mathfrak{T}}(\theta; \rho) - A(\theta)$  is smallest). In principle, it appears as if the computation of this optimum might be difficult, since it seems to entail searching over all probability distributions over the set  $\mathfrak{T}$  of all spanning trees of the graph. The number  $|\mathfrak{T}|$  of spanning trees can be computed via the matrix-tree theorem [217], and for sufficiently dense graphs, it

grows rapidly enough to preclude direct optimization over  $\mathfrak{T}$ . However, Wainwright et al. [241] show that it is possible to directly optimize the vector  $\rho_e$  of edge appearance probabilities, subject to the constraint of membership in the spanning tree polytope  $\mathcal{S}(G)$ . Although the polytope  $\mathcal{S}(G)$  has a prohibitively large number of inequalities, it is possible to maximize a linear function over it—equivalently, to solving a maximum weight spanning tree problem—by a greedy algorithm (see, for instance, Schrijver [207]). Using this fact, it is possible [241] to develop a version of the conditional gradient algorithm [19] for efficiently computing the optimal choice of edge appearance probabilities.

**Proof of Theorem 4:** (a) Our discussion preceding the statement of the result establishes that  $B_{\mathfrak{T}}$  is an upper bound on  $A$ . It remains to establish the strict convexity, and resulting uniqueness of the global optimum. Note that the cost function defining the function  $B_{\mathfrak{T}}$  consists of a linear term  $\langle \theta, \mu \rangle$  and a convex combination— $\mathbb{E}_{\rho}[A^*(\mu(T))]$  of tree entropies, and hence is concave. Moreover, the constraint set  $\mathbb{L}(G)$  is a polytope, and hence convex. Therefore, the variational problem (7.12) is always convex. We now establish uniqueness of the optimum when  $\rho_e > 0$ . To simplify details of the proof, we assume without loss of generality that we are working in a minimal representation. (If the variational problem is formulated in an overcomplete representation, it can be reduced to an equivalent minimal formulation.) To establish uniqueness, it suffices to establish that the function  $\mathbb{E}_{\rho}[A^*(\mu(T))]$  is strictly convex when  $\rho_e > 0$ . This function is a convex combination of functions of the form  $A^*(\mu(T))$ , each of which is strictly convex in the (non-zero) components of  $\mu(T)$ , but independent of the other components in the full vector  $\mu$ . For any vector  $\lambda \in \mathbb{R}^d$ , define  $\Pi^T(\lambda)_{\alpha} = \lambda_{\alpha}$  if  $\alpha \in \mathcal{I}(T)$ , and  $\Pi^T(\lambda)_{\alpha} = 0$  otherwise. We then have

$$\langle \lambda, \nabla^2 A^*(\mu(T)) \lambda \rangle = \langle \Pi^T(\lambda), \nabla^2 A^*(\mu(T)) \Pi^T(\lambda) \rangle \geq 0,$$

with strict inequality unless  $\Pi^T(\lambda) = 0$ . Now the condition  $\rho_e > 0$  for all  $e \in E$  ensures that  $\lambda \neq 0$  implies that  $\Pi^T(\lambda)$  must be different from zero for at least one tree  $T'$ . Therefore, for any  $\lambda \neq 0$ , we have

$$\langle \lambda, \mathbb{E}_{\rho}[\nabla^2 A^*(\mu(T))] \lambda \rangle \geq \langle \Pi^{T'}(\lambda), \nabla^2 A^*(\mu(T')) \Pi^{T'}(\lambda) \rangle > 0,$$

which establishes the assertion of strict convexity.

(b) Derivation of the TRW-S updates (7.13) is entirely analogous to the proof of Theorem 3: setting up the Lagrangian, taking derivatives, and performing some algebra yields the result (see Wainwright et al. [241] for full details). The uniqueness follows from the strict convexity established in part (a).  $\square$

### 7.2.2 Reweighted Kikuchi approximations

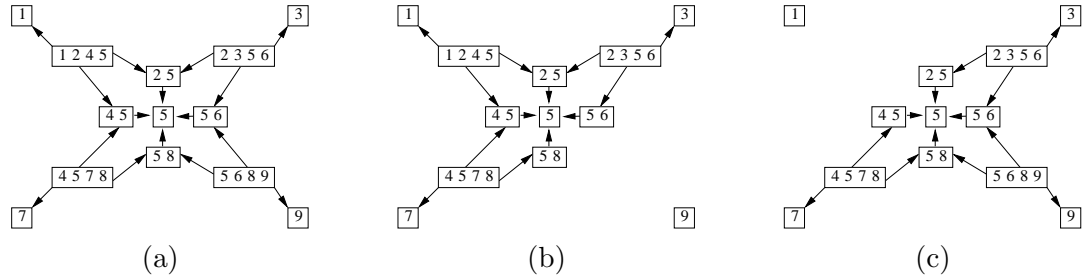
A natural extension of convex combinations of trees, entirely analogous to the transition from the Bethe to Kikuchi variational problems, is to take convex combinations of hypertrees. This extension was sketched out in Wainwright et al. [241], and has been further studied by various researchers [255, 246]. Here we describe the basic idea with one illustrative example. For a given treewidth  $t$ , consider the set of all hypertrees of width  $\mathfrak{T}(t)$  of width less than or equal to  $t$ . Of course, the underlying assumption is that  $t$  is sufficiently small that performing exact computations on hypertrees of this width is feasible. Applying Proposition 8 to a convex combination based on a probability distribution  $\rho = \{\rho(T)\}$  over the set of all hypertrees of width at most  $t$ , we obtain an upper bound on the entropy of the form

$$A^*(\mu) \leq -\mathbb{E}_\rho[A^*(\mu(T))] = -\sum_T \rho(T) A^*(\mu(T)). \quad (7.15)$$

For a fixed  $\rho$ , our strategy is to optimize the right-hand side of this upper bound over all pseudomarginals that are consistent on each hypertree. The resulting constraint set is precisely the polytope  $\mathbb{L}_t(G)$  defined in equation (4.46). Overall, the hypertree analog of Theorem 4 asserts that the log partition function  $A$  is upper bounded by the hypertree-based convex surrogate:

$$A(\theta) \leq B_{\mathfrak{T}(t)}(\theta; \rho) := \max_{\tau \in \mathbb{L}(G)} \{\langle \tau, \theta \rangle - \mathbb{E}_\rho[A^*(\mu(T))]\}. \quad (7.16)$$

Moreover, the cost function in this variational problem is concave for all choices of distributions  $\rho$  over the hypertrees. Equation (7.16) is the hypertree analog of equation (7.12); indeed, it reduces to the latter equation in the special case  $t = 1$ .



**Fig. 7.2.** Hyperforests embedded within augmented hypergraphs. (a) An augmented hypergraph for the  $3 \times 3$  grid with maximal hyperedges of size 4 that satisfies the single counting criterion. (b) One hyperforest of width three embedded within (a). (c) A second hyperforest of width three.

**Example 34 (Convex combinations of hypertrees).** Let us derive an explicit form of equation (7.16) for a particular hypergraph and choice of hypertrees. The original graph is the  $3 \times 3$  grid, as illustrated in the earlier Figure D.1(a). Based on this grid, we construct the augmented hypergraph shown in Figure 7.2(a), which has the hyperedge set  $E$  given by

$$\{(1245), (2356), (4578), (5689), (25), (45), (56), (58), (5), (1), (3), (7), (9)\}.$$

It is straightforward to verify that it satisfies the single counting criterion (see Appendix D for background).

Now consider a convex combination of four hypertrees, each obtained by removing one of the 4-hyperedges from the edge set. For instance, shown in Figure 7.2(b) is one particular acyclic substructure  $T^1$  with hyperedge set  $E(T^1)$

$$\{(1245), (2356), (4578), (25), (45), (56), (58), (5), (1), (3), (7), (9)\},$$

obtained by removing  $(5689)$  from the full hyperedge set  $E$ . To be precise, the structure  $T^1$  so defined is a spanning hyperforest, since it consists of two connected components (namely, the isolated hyperedge  $(9)$  along with the larger hypertree). This choice, as opposed to a spanning hypertree, turns out to be simplify the development to follow. Figure 7.2(c) shows the analogous spanning hyperforest  $T^2$  obtained by removing hyperedge  $(1245)$ ; the final two hyperforests  $T^3$  and  $T^4$  are defined analogously.

To specify the associated hypertree factorization, we first compute the form of  $\varphi_h$  for the maximal hyperedges (i.e., of size four). For instance, looking at the hyperedge  $h = (1245)$ , we see that hyperedges (25), (45), (5), and (1) are contained within it. Thus, using the definition in equation (4.40), we write (suppressing the functional dependence on  $x$ ):

$$\varphi_{1245} = \frac{\tau_{1245}}{\varphi_{25} \varphi_{45} \varphi_5 \varphi_1} = \frac{\tau_{1245}}{\frac{\tau_{25} \tau_{45}}{\tau_5} \tau_5 \tau_1} = \frac{\tau_{1245} \tau_5}{\tau_{25} \tau_{45} \tau_1}.$$

Having calculated all the functions  $\varphi_h$ , we can combine them, using the hypertree equation (4.41), in order to obtain the following factorization for a distribution on  $T^1$ :

$$p_{\tau(T^1)}(x) = \left[ \frac{\tau_{1245} \tau_5}{\tau_{25} \tau_{45} \tau_1} \right] \left[ \frac{\tau_{2356} \tau_5}{\tau_{25} \tau_{56} \tau_3} \right] \left[ \frac{\tau_{4578} \tau_5}{\tau_{45} \tau_{58} \tau_7} \right] \\ \times \left[ \frac{\tau_{25}}{\tau_5} \right] \left[ \frac{\tau_{45}}{\tau_5} \right] \left[ \frac{\tau_{56}}{\tau_5} \right] \left[ \frac{\tau_{58}}{\tau_5} \right] \left[ \tau_1 \right] \left[ \tau_3 \right] \left[ \tau_5 \right] \left[ \tau_7 \right] \left[ \tau_9 \right]. \quad (7.17)$$

Here each term within square brackets corresponds to  $\varphi_h$  for some hyperedge  $h \in E(T^1)$ ; for instance, the first three terms correspond to the three maximal 4-hyperedges in  $T^1$ . Although this factorization could be simplified, leaving it in its current form makes the connection to Kikuchi approximations more explicit; in particular, the factorization (7.17) leads immediately to a decomposition of the entropy. In an analogous manner, it is straightforward to derive factorizations and entropy decompositions for the remaining three hyperforests  $\{T^i, i = 2, 3, 4\}$ .

Now let  $E_4 = \{(1245), (2356), (5689), (4578)\}$  denote the set of all 4-hyperedges. We then form the convex combination of the four (negative) entropies with uniform weight  $1/4$  on each  $T^i$ ; this weighted sum  $\sum_{i=1}^4 \frac{1}{4} A^*(\tau(T^i))$  takes the form

$$\frac{3}{4} \sum_{h \in E_4} \sum_{x_h} \tau_h(x_h) \log \varphi_h(x_h) + \sum_{s \in \{2,4,6,8\}} \sum_{x_{s5}} \tau_{s5}(x_{s5}) \log \frac{\tau_{s5}(x_{s5})}{\tau_5(x_5)} \\ + \sum_{s \in \{1,3,5,7,9\}} \sum_{x_s} \tau_s(x_s) \log \tau_s(x_s). \quad (7.18)$$

The weight  $3/4$  arises because each of the maximal hyperedges  $h \in E_4$  appears in three of the four hypertrees. All of the (non-maximal) hy-

peredge terms receive a weight of one, because they appear in all four hypertrees. Overall, then, these weights represent hyperedge appearance probabilities for this particular example, in analogy to ordinary edge appearance probabilities in the tree case. We now simplify the expression in equation (7.18) by expanding and collecting terms; doing so yields that the sum  $-\sum_{i=1}^4 \frac{1}{4} A^*(\tau(T^i))$  is equal to the following weighted combination of entropies:

$$\begin{aligned} \frac{3}{4} [H_{1245} + H_{2356} + H_{5689} + H_{4578}] - \frac{1}{2} [H_{25} + H_{45} + H_{56} + H_{58}] \\ + \frac{1}{4} [H_1 + H_3 + H_7 + H_9]. \end{aligned} \quad (7.19)$$

If, on the other hand, starting from equation (7.18) again, suppose that we included each maximal hyperedge with a weight of 1, instead of  $3/4$ . Then, after some simplification, we would find that the (negative of) equation (7.18) is equal to the following combination of local entropies

$$[H_{1245} + H_{2356} + H_{5689} + H_{4578}] - [H_{25} + H_{45} + H_{56} + H_{58}] + H_5,$$

which is equivalent to the Kikuchi approximation derived in Example 19. However, the choice of all ones for the hyperedge appearance probabilities is invalid—that is, it could never arise from taking a convex combination of hypertree entropies. ♣

More generally, any entropy approximation formed by taking such convex combinations of hypertree entropies will necessarily be convex. In contrast, with the exception of certain special cases [107, 184, 185, 163], Kikuchi and other hypergraph-based entropy approximations are typically not convex. In analogy to the tree-reweighted sum-product algorithm, it is possible to develop hypertree-reweighted forms of generalized sum-product updates [255, 246]. With a suitable choice of convex combination, the underlying variational problem will be strictly convex, and so the associated hypertree-reweighted sum-product algorithms will have a unique fixed point. An important difference from the case of ordinary trees, however, is that optimization of the hyperedge weights  $\rho_e$  cannot be performed exactly using a greedy algorithm.

Indeed, the hypertree analog of the maximum weight spanning tree problem is known to be NP-hard [126].

### 7.2.3 Convex forms of expectation-propagation

As discussed in Section 4.3, the family of expectation-propagation algorithms [169, 172], as well as related moment-matching algorithms [181, 61, 112], can be understood in terms of term-by-term entropy approximations. From this perspective, it is clear how to derive convex analogs of the variational problem (4.63) that underlies expectation-propagation.

In particular, recall from equation (4.69) the form of the term-by-term entropy approximation

$$H_{\text{ep}}(\tau, \tilde{\tau}) := H(\tau) + \sum_{\ell=1}^{d_I} [H(\tau, \tilde{\tau}^\ell) - H(\tau)],$$

where  $H(\tau)$  represents the entropy of the base distribution, and  $H(\tau, \tilde{\tau}^\ell)$  represents the entropy of the base plus  $\ell^{\text{th}}$  term. Implicitly, this entropy approximation is weighting each term  $i = 1, \dots, d_I$  with weight one. Suppose instead that we consider a family of non-negative weights  $\{\rho(1), \dots, \rho(d_I)\}$  over the different terms, and use the reweighted term-by-term entropy approximation

$$H_{\text{ep}}(\tau, \tilde{\tau}; \rho) := H(\tau) + \sum_{\ell=1}^{d_I} \rho(\ell) [H(\tau, \tilde{\tau}^\ell) - H(\tau)]. \quad (7.20)$$

For suitable choice of weights—for instance, if we enforce the constraint  $\sum_{\ell=1}^{d_I} \rho(\ell) = 1$ —this reweighted entropy approximation is concave. Other choices of non-negative weights could also produce concave entropy approximations. The EP outer bound  $\mathcal{L}(\phi; \Phi)$  on the set  $\mathcal{M}(\phi)$ , as defined in equation (4.61), is a convex set by definition. Consequently, if we combine a convex entropy approximation (7.20) with this outer bound, we obtain “convexified” forms of the EP variational problem. Following the line of development in Section 4.3, it is also possible to derive Lagrangian algorithms for solving these optimization problems. Due to underlying structure, the updates can again be case in terms of



moment-matching steps. In contrast to standard forms of EP, a benefit of these algorithms would be the existence of unique fixed points, assured by the convexity of the underlying variational principle. With different motivations, not directly related to convexity and uniqueness, Minka [170] has discussed reweighted forms of EP, referred to as “power EP.” Seeger et al. [209] reported that reweighted forms of EP appear to have better empirical convergence properties than standard EP. It would be interesting to study connections between the convexity of the EP entropy approximation, and the convergence and stability properties of the EP updates.

#### 7.2.4 Planar graph decomposition

Trees (and the associated hierarchy of hypertrees) represent the best known classes of tractable graphical models. However, exact computation is also tractable for certain other classes of graphical models. Unlike hypertrees, certain models may be intractable in general, but tractable for specific settings of the parameters. For instance, consider the subclass of planar graphs (which can be drawn in the 2-D plane without any crossing of the edges, for instance 2-D grid models). Although exact computation for an arbitrary Markov random field on a planar graph is known to be intractable, if one considers only binary random variables without any observation potentials (i.e., Ising models of the form  $p_{\theta}(x) \propto \exp\{\sum_{(s,t) \in E} \theta_{st} x_s x_t\}$ ), then it is possible to compute the cumulant function  $A(\theta)$  exactly, using clever combinatorial reductions due independently to Fisher [78] and Kastelyn [129]. Globerson and Jaakkola [91] exploit these facts in order to derive a variational relaxation based on convex combination of planar subgraphs, as well as iterative algorithms for computing the optimal bound. They provide experimental results for various types graphs (both grids and fully connected graphs), with results superior to the TRW relaxation (7.12) for most coupling strengths, albeit with higher computational cost. As they point out, it should also be possible to combine hypertrees with planar graphs to form composite approximations.

### 7.3 Other convex variational methods

The previous section focused exclusively on variational methods based on the idea of convex combinations, as stated in equation (7.8). However, not all convex variational methods need arise in this way; more generally, a convex variational method requires only (a) some type of convex approximation to the dual function  $A^*$  (negative entropy), and (b) convex outer bound on the set  $\mathcal{M}$ . This basic two-step procedure can be used to derive many convex relaxations of the marginalization problem.

#### 7.3.1 Semidefinite constraints and log-determinant bound

We illustrate this two-step procedure by describing a log-determinant relaxation for discrete Markov random fields [243]. This method differs qualitatively from the previously described Bethe/Kikuchi methods, in that it uses a semidefinite outer bound on the marginal polytope  $\mathbb{M}(G)$ . Although the basic ideas are more generally applicable, consider a binary random vector  $X \in \{0, 1\}^m$ , with a distribution in the Ising form

$$p_\theta(x) = \exp \left\{ \sum_{s \in V} \theta_s x_s + \sum_{(s,t)} \theta_{st} x_s x_t - A(\theta) \right\}. \quad (7.21)$$

Without loss of generality, we assume that the underlying graph is the complete graph  $K_m$ , so that the marginal polytope of interest is  $\mathbb{M}(K_m)$ . Of course, a problem defined on an arbitrary  $G = (V, E)$  can be embedded into the complete graph by setting  $\theta_{st} = 0$  for all  $(s, t) \notin E$ . With this set-up, the mean parameterization consists of the vector  $\mu \in \mathbb{M}(K_m) \subset \mathbb{R}^{m + \binom{m}{2}}$ , consisting of the elements  $\mu_s = \mathbb{E}[X_s]$  for each node, and quantities  $\mu_{st} = \mathbb{E}[X_s X_t]$  for each edge.

First, we describe how to approximate the negative dual function  $-A^*$ , or entropy function. An important characterization of the multivariate Gaussian is as the maximum entropy distribution subject to covariance constraints [54]. In particular, the differential entropy  $h(\tilde{X})$  of any continuous random vector  $\tilde{X}$  is upper bounded by the entropy

of a Gaussian with matched covariance. In analytical terms, we have

$$h(\tilde{X}) \leq \frac{1}{2} \log \det \text{cov}(\tilde{X}) + \frac{m}{2} \log(2\pi e), \quad (7.22)$$

where  $\text{cov}(\tilde{X})$  is the covariance matrix of  $\tilde{X}$ . The upper bound (7.22) is not directly applicable to a random vector taking values in a discrete space (since differential entropy in this case diverges to minus infinity). Therefore, in order to exploit this bound for the discrete random vector  $X \in \{0, 1\}^m$  of interest, it is necessary to construct a suitably matched continuous version of  $X$ . One method to do so is by the addition of an independent random vector  $U$ , such that the delta functions in the density of  $X$  are smoothed out. In particular, let us define the continuous random vector  $\tilde{X} := X + U$ , where  $U$  is independent of  $X$ , with independent components distributed uniformly as  $U_s \sim \mathcal{U}[-\frac{1}{2}, \frac{1}{2}]$ . A key property of this construction is that it matches the discrete entropy of  $X$  with the differential entropy of  $\tilde{X}$ —that is,  $h(\tilde{X}) = H(X)$ ; see Wainwright and Jordan [243] for the proof. Since  $X$  and  $U$  are continuous by construction, a straightforward computation yields  $\text{cov}(\tilde{X}) = \text{cov}(X) + \frac{1}{12}I_m$ , so equation (7.22) yields the upper bound

$$H(X) \leq \frac{1}{2} \log \det[\text{cov}(X) + \frac{1}{12}I_m] + \frac{m}{2} \log(2\pi e).$$

Second, we need to provide an outer bound on the marginal polytope for which the entropy approximation above is well defined. The simplest such outer bound is obtained by observing that the covariance matrix  $\text{cov}(X)$  must always be positive semidefinite. (Indeed, for any vector  $a \in \mathbb{R}^m$ , we have  $\langle a, \text{cov}(X)a \rangle = \text{var}(\langle a, X \rangle) \geq 0$ .) Note that elements of the covariance matrix  $\text{cov}(X)$  can be expressed in terms of the mean parameters  $\mu_s$  and  $\mu_{st}$ . In particular, consider the

$(m + 1) \times (m + 1)$  matrix

$$\begin{aligned}
 M_1[\mu] &= \mathbb{E} \left[ \begin{bmatrix} 1 \\ X \end{bmatrix} \begin{bmatrix} 1 & X \end{bmatrix} \right] \\
 &= \begin{bmatrix} 1 & \mu_1 & \mu_2 & \cdots & \cdots & \mu_{m-1} & \mu_m \\ \mu_1 & \mu_1 & \mu_{12} & \cdots & \cdots & \mu_{1(m-1)} & \mu_{1m} \\ \mu_2 & \mu_{21} & \mu_2 & \mu_{23} & \cdots & \cdots & \mu_{2m} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mu_{m-1} & \mu_{(m-1)1} & \cdots & \cdots & \cdots & \mu_{(m-1)} & \mu_{(m-1)m} \\ \mu_m & \mu_{m1} & \cdots & \cdots & \cdots & \mu_{m(m-1)} & \mu_m \end{bmatrix}.
 \end{aligned}$$

By the Schur complement formula [111], the constraint  $M_1[\mu] \succeq 0$  is equivalent to the constraint  $\text{cov}(X) \succeq 0$ . Consequently, the set

$$\mathbb{S}_1(K_m) := \left\{ \mu \in \mathbb{R}^{m+\binom{m}{2}} \mid M_1[\mu] \succeq 0 \right\}$$

is an outer bound on the marginal polytope  $\mathbb{M}(K_m)$ . It is not difficult to see that the set  $\mathbb{S}_1(K_m)$  is closed, convex, and also bounded, since it is contained within the hypercube  $[0, 1]^{m+\binom{m}{2}}$ .

We now have the two necessary ingredients to define a convex variational principle based on a log-determinant relaxation [243]. In addition to the constraint  $\mu \in \mathbb{S}_1(K_m)$ , one might imagine enforcing other constraints (e.g., the linear constraints defining the set  $\mathbb{L}(G)$ ). Accordingly, let  $\mathbb{B}(K_m)$  denote any convex and compact outer bound on  $\mathbb{M}(K_m)$  that is contained within  $\mathbb{S}_1(K_m)$ . With this notation, the cumulant function  $A(\theta)$  is upper bounded by the log-determinant convex surrogate

$$\begin{aligned}
 B_{\text{LD}}(\theta) &:= \\
 \max_{\mu \in \mathbb{B}(K_m)} &\left\{ \langle \theta, \mu \rangle + \frac{1}{2} \log \det \left[ M_1(\mu) + \frac{1}{12} \text{blkdiag}[0, I_m] \right] \right\} + \frac{m}{2} \log(2\pi e),
 \end{aligned} \tag{7.23}$$

where  $\text{blkdiag}[0, I_m]$  is a  $(m + 1) \times (m + 1)$  block-diagonal matrix. By enforcing the inclusion  $\mathbb{B}(K_m) \subseteq \mathbb{S}_1(K_m)$  we guarantee that the matrix  $M_1(\mu)$ , and hence the matrix sum

$$M_1(\mu) + \frac{1}{12} \text{blkdiag}[0, I_m]$$

will always be positive semidefinite. Importantly, the optimization problem in equation (7.23) is a standard determinant maximization problem, for which efficient interior point methods have been developed [e.g., 230]. Wainwright and Jordan [243] derive an efficient algorithm for solving a slightly weakened form of the log-determinant problem (7.23), and provide empirical results for various types of graphs. Banerjee et al. [7] exploit the log-determinant relaxation (7.23) for a sparse model selection procedure, and develop coordinate-wise algorithms for solving a broader class of convex programs.

### 7.3.2 Other entropy approximations and polytope bounds

There are many other approaches, using different approximations to the entropy and outer bounds on the marginal polytope, that lead to interesting convex variational principles. For instance, Globerson and Jaakkola [92] exploit the notion of conditional entropy decomposition in order to construct whole families of upper bounds on the entropy, including as a special case the hypertree-based entropy bounds that underlie the reweighted Bethe/Kikuchi approaches. Sontag and Jaakkola [214] examine the effect of tightened outer bounds on the marginal polytope, including the so-called cycle inequalities [66, 183]. Among other issues, they examine the differing effects of changing the entropy approximation versus refining the outer bound on the marginal polytope.

## 7.4 Algorithmic stability

As discussed earlier, an obvious benefit of approximate marginalization methods based on convex variational principles is the resulting uniqueness of the minima, and hence lack of dependence on algorithmic details, or initial conditions. In this section, we turn to a less obvious but equally important benefit of convexity, of particular relevance in a statistical setting. More specifically, in typical applications of graphical models, the model parameters (or mode structure) themselves are uncertain. (For instance, they may have been estimated from noisy or incomplete data.) In such a setting, then, it is natural to ask whether the output of a given variational method remains stable under small

perturbations to the model parameters. Not all variational algorithms are stable in this sense; for instance, both mean field methods and the ordinary sum-product algorithm can be highly unstable, with small perturbations in the model parameters leading to very large changes in the fixed point(s) and outputs of the algorithm.

In this section, we show that any variational method based on a suitably convex optimization problem—in contrast to mean field and ordinary sum-product—has an associated guarantee of Lipschitz stability. In particular, let  $\tau(\theta)$  denote the output when some variational method is applied to the model  $p_\theta$ . Given two norms  $\|\cdot\|_a$  and  $\|\cdot\|_b$ , we say that the method is *globally Lipschitz stable* if there exists some finite constant  $L$  such that

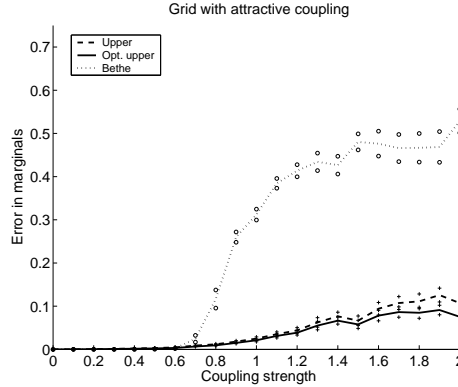
$$\|\tau(\theta) - \tau(\theta')\|_a \leq L\|\theta - \theta'\|_b, \quad (7.24)$$

for any  $\theta, \theta' \in \Omega$ . This definition is one way of formalizing the notion of *algorithmic stability*: namely, that if the difference  $\|\theta - \theta'\|_b$  between models indexed by  $\theta$  and  $\theta'$  is small, then the difference  $\|\tau(\theta) - \tau(\theta')\|_a$  between the algorithm's outputs is also correspondingly small. As an illustration, Figure 7.3 compares the behavior of the ordinary sum-product algorithm—known to be unstable in certain regimes—to the tree-reweighted sum-product algorithm. Note how the sum-product algorithm degrades rapidly beyond a certain critical coupling strength, whereas the performance of the tree-reweighted algorithm varies smoothly as a function of the coupling strength.

Let us now try to gain some theoretical insight into which variational methods satisfy this Lipschitz property. Consider a generic variational method of the form

$$B(\theta) = \sup_{\tau \in \mathcal{L}} \{\langle \theta, \tau \rangle - B^*(\tau)\}, \quad (7.25)$$

where  $\mathcal{L}$  is a convex outer bound on the set  $\mathcal{M}$ , and  $B^*$  is a convex approximation to the dual function  $A^*$ . Let us consider the properties of the surrogate function  $B$  that we have defined. As mentioned previously, it is always a convex function of  $\theta$ . In terms of properties beyond convexity, they are determined by the nature of the function  $B^*$ . For instance, suppose that  $B^*$  is strictly convex; then the optimum (7.25)



**Fig. 7.3.** Contrast between the instability of ordinary sum-product and stability of tree-reweighted sum-product [241]. Plots show the error between the true marginals and correct marginals versus the coupling strength in a binary pairwise Markov random field. Note that the ordinary sum-product algorithm is very accurate up to a critical coupling ( $\approx 0.70$ ), after which it degrades rapidly. On the other hand, the performance of TRW message-passing varies smoothly as a function of the coupling strength.

is uniquely attained, so that Danskin's theorem [19] ensures that  $B$  is differentiable.

In terms of Lipschitz stability, it turns out that a somewhat stronger notion of convexity is required; in particular, a differentiable function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is *strongly convex* with parameter  $c$  if for all  $y, z \in \mathbb{R}^d$ ,

$$f(z) \geq f(y) + \langle \nabla f(y), z - y \rangle + \frac{c}{2} \|z - y\|^2. \quad (7.26)$$

Note that any differentiable convex function satisfies this inequality (7.26) with  $c = 0$ ; the essence of strong convexity is that the same inequality is satisfied with slack  $\frac{c}{2} \|z - y\|^2$ . (In fact, an equivalent characterization of strong convexity is to require that the function  $f(z) - \frac{c}{2} \|z\|^2$  be convex [109].) In the setting of the convex surrogate (7.25), suppose that the function  $B$  is strictly convex, and the (negative) entropy approximation  $B^*$  is strongly convex, say with parameter  $1/L > 0$ . Under this assumption, it can be shown [236, 109] that the output  $\tau(\theta) = \nabla B(\theta)$  of the variational method is Lipschitz stable, in the sense of definition (7.24), with parameter  $L$  and norms  $\|\cdot\|_a = \|\cdot\|_b = \|\cdot\|_2$ .

Thus, when using variational methods of the form (7.25), the issue of algorithmic stability reduces to strong convexity of the negative entropy approximation  $B^*$ . What existing variational methods are based on strongly convex entropy approximations? For instance, the tree-reweighted Bethe entropy (7.11) is strongly convex [236] for any pairwise Markov random field, so that as a corollary, the tree-reweighted sum-product algorithm is guaranteed to be globally Lipschitz stable. This Lipschitz stability provides a theoretical explanation of the behavior illustrated in Figure 7.3. Similarly, empirical results due to Wiegerinck [255] show that suitably reweighted forms of generalized belief propagation (GBP) also behave in a stable manner (unlike standard GBP), and this stability can be confirmed theoretically via strong convexity. It can also be shown that the log-determinant relaxation (7.23) is Lipschitz stable.

## 7.5 Convex surrogates in parameter estimation

Until now, we have discussed convex variational methods in the context of computing approximations to the cumulant function  $A(\theta)$ , as well as approximate mean parameters  $\mu$ . Here we discuss an alternative use of convex surrogates of the form (7.25), namely for approximate parameter estimation. Recall from Section 6 our discussion of maximum likelihood (ML) estimation in exponential families: it entails maximizing the log likelihood of the data, given by

$$\ell(\theta; X_1^n) = \ell(\theta) = \langle \theta, \hat{\mu} \rangle - A(\theta), \quad (7.27)$$

where  $\hat{\mu} := \frac{1}{n} \sum_{i=1}^n \phi(X^i)$  is the empirical expectation of the sufficient statistics given the data  $X_1^n = (X^1, \dots, X^n)$ . Recall also that the derivative of the log likelihood takes the form  $\nabla \ell(\theta) = \hat{\mu} - \mu(\theta)$ , where  $\mu(\theta) = \mathbb{E}_\theta[\phi(X)]$  are the mean parameters under the current model parameter. As discussed in Section 6, this fact motivates a natural approach to approximate parameter estimation, in which the true model parameters  $\mu(\theta)$  are approximated by the output  $\tau(\theta)$  of some variational method. For instance, this approach has frequently been suggested, using sum-product to compute the approximate mean parameters  $\tau(\theta)$ . Such heuristics may yield good results in some settings;



however, if the underlying variational method is non-convex, such a heuristic cannot be interpreted as minimizing an approximation to the likelihood. Indeed, given multiple fixed points, its behavior can be unstable and erratic; the papers [238, 236] provide some cautionary instances of poor behavior with ordinary sum-product.

### 7.5.1 Surrogate likelihoods

On the other hand, the perspective of convex surrogates to  $A$ , of the form (7.25), suggests a related but more principled method for approximate ML estimation, based on maximizing a concave approximation to the likelihood. Let us assume that the outer bound  $\mathcal{L}$  is compact and the negative entropy approximation  $B^*$  is strictly convex, so that the maximum (7.25) is uniquely attained, say at some  $\tau(\theta) \in \mathcal{L}$ . Recall that Danskin's theorem [19] shows that the convex surrogate  $B$  is differentiable, with  $\nabla B(\theta) = \tau(\theta)$ . (Note the parallel with the gradient  $\nabla A$ , and its link to the mean parameterization, as summarized in Proposition 2). Given a convex surrogate  $B$ , let us define the associated *surrogate likelihood* as

$$\ell_B(\theta; X_1^n) = \ell_B(\theta) := \langle \theta, \hat{\mu} \rangle - B(\theta). \quad (7.28)$$

Assuming that  $B$  is strictly convex, then this surrogate likelihood is strictly concave, and so (at least for compact parameter spaces  $\Omega$ ), we have the well-defined *surrogate ML estimate*

$$\tilde{\theta}_B := \arg \max_{\theta \in \Omega} \ell_B(\theta; X_1^n).$$

Note that when  $B$  upper bounds the true cumulant function (as in the convex variational methods described in Section 7.1), then the surrogate likelihood  $\ell_B$  lower bounds the true likelihood, and  $\tilde{\theta}_B$  is obtained from maximizing this lower bound.

### 7.5.2 Optimizing surrogate likelihoods

In general, it is always relatively straightforward to compute the surrogate estimate  $\tilde{\theta} = \tilde{\theta}_B$ . Indeed, from our discussion above, the derivative of the surrogate likelihood is  $\nabla \ell_B(\theta) = \hat{\mu} - \tau(\theta)$ , where  $\tau(\theta) = \nabla B(\theta)$

are the approximate mean parameters computed by the variational method associated with  $B$ . Since the optimization problem is concave, a standard gradient ascent algorithm (among other possibilities) could be used to compute  $\tilde{\theta}_B$ .

Interestingly, in certain cases, the unregularized surrogate MLE can be specified in closed-form, without the need for any optimization. As an example, consider the reweighted Bethe surrogate (7.12), and the associated surrogate likelihood. The surrogate MLE is defined by the fixed point relation  $\nabla \ell_B(\tilde{\theta}) = 0$ , which has a very intuitive interpretation. Namely, the model parameters  $\tilde{\theta}$  are chosen such that, when the TRW sum-product updates (7.13) are applied to the model  $p_{\tilde{\theta}}$ , the resulting pseudo mean parameters  $\tau(\tilde{\theta})$  computed by the algorithm are equal to the empirical mean parameters  $\hat{\mu}$ . Assuming that the empirical mean vector has strictly positive elements, the unique vector  $\tilde{\theta}$  with this property has the closed-form expression:

$$\tilde{\theta}_{s;j} = \log \hat{\mu}_{s;j}, \quad \text{and} \quad (7.29a)$$

$$\tilde{\theta}_{st;jk} = \rho_{st} \log \frac{\hat{\mu}_{st;jk}}{\hat{\mu}_{s;j} \hat{\mu}_{t;k}}. \quad (7.29b)$$

Here  $\hat{\mu}_{s;j} = \widehat{\mathbb{E}}[\mathbb{I}_j(X_s)]$  is the empirical probability of the event  $\{X_s = j\}$ , and  $\hat{\mu}_{st;jk}$  is the empirical probability of the event  $\{X_s = j, X_t = k\}$ . By construction, the vector  $\tilde{\theta}$  defined by equation (7.29) has the property that  $\hat{\mu}$  are the pseudomarginals associated with mode  $p_{\tilde{\theta}}$ . The proof of this fact relies on the tree-based reparameterization interpretation of the ordinary sum-product algorithm, as well as its reweighted extensions [238, 241]; see Section 4.1.4 for details on reparameterization. Similarly, closed form solutions can be obtained for reweighted forms of Kikuchi and region graph approximations (Section 7.2.2), and convexified forms of expectation-propagation and other moment-matching algorithms. Such a closed-form solution eliminates the need for any iterative algorithms in the complete data setting (as described here), and can substantially reduce computational overhead if a variational method is used within an inner loop in the incomplete data setting.

### 7.5.3 Penalized surrogate likelihoods

Rather than maximizing the log likelihood (7.27) itself, it is often better to maximize a penalized likelihood function  $\tilde{\ell}(\theta; \lambda) := \ell(\theta) - \lambda R(\theta)$ , where  $\lambda > 0$  is a regularization constant and  $R$  is a penalty function, assumed to be convex but not necessarily differentiable (e.g.,  $\ell_1$  penalization). Such procedures can be motivated either from a frequentist perspective (and justified in terms of improved loss and risk), or from a Bayesian perspective, where the quantity  $\lambda R(\theta)$  might arise from a prior over the parameter space. If maximization of the likelihood is intractable then regularized likelihood is still intractable, so let us consider approximate procedures based on maximizing a regularized surrogate likelihood (RSL)  $\tilde{\ell}_B(\theta; \lambda) := \ell_B(\theta) - \lambda R(\theta)$ . Again, given that  $\ell_B$  is convex and differentiable, this optimization problem could be solved by a variety of standard methods.

Here we describe an alternative formulation of optimizing the penalized surrogate likelihood, which exploits the variational manner in which the surrogate  $B$  was defined (7.25). If we reformulate the problem as one of minimizing the negative RSL, substituting the definition (7.25) yields the equivalent saddlepoint problem:

$$\begin{aligned} \inf_{\theta \in \Omega} \{-\langle \theta, \hat{\mu} \rangle + B(\theta) + \lambda R(\theta)\} &= \inf_{\theta \in \Omega} \left\{ -\langle \theta, \hat{\mu} \rangle + \sup_{\tau \in \mathcal{L}} \{\langle \theta, \tau \rangle - B^*(\tau)\} + \lambda R(\theta) \right\} \\ &= \inf_{\theta \in \Omega} \sup_{\tau \in \mathcal{L}} \left\{ \langle \theta, \tau - \hat{\mu} \rangle - B^*(\tau) + \lambda R(\theta) \right\}. \end{aligned}$$

Note that this saddlepoint problem involves convex minimization and concave maximization; therefore, under standard regularity assumptions (e.g., either compactness of the constraint sets, or coercivity of the function; see Section VII of Hiriart-Urruty and Lemaréchal [109]), we can exchange the sup and inf to obtain

$$\begin{aligned} \inf_{\theta \in \Omega} \{-\ell_B(\theta) + \lambda R(\theta)\} &= \sup_{\tau \in \mathcal{L}} \inf_{\theta \in \Omega} \left\{ \langle \theta, \tau - \hat{\mu} \rangle - B^*(\tau) + \lambda R(\theta) \right\} \\ &= \sup_{\tau \in \mathcal{L}} \left\{ -B^*(\tau) - \lambda \sup_{\theta \in \Omega} \left\{ \langle \theta, \frac{\tau - \hat{\mu}}{\lambda} \rangle - R(\theta) \right\} \right\} \\ &= \sup_{\tau \in \mathcal{L}} \left\{ -B^*(\tau) - \lambda R_{\Omega}^* \left( \frac{\tau - \hat{\mu}}{\lambda} \right) \right\}, \quad (7.30) \end{aligned}$$

where  $R_\Omega^*$  is the conjugate dual function of  $R(\theta) + \mathbb{I}_\Omega(\theta)$ . Consequently, assuming that it is feasible to compute the dual function  $R_\Omega^*$ , maximizing the RSL reduces to an unconstrained optimization problem involving the function  $B^*$  used to define the surrogate function  $B$ . We illustrate with a particular family of examples:

**Example 35 (Regularized surrogate Bethe likelihood).** Suppose that we use the surrogate function  $B_{\mathfrak{T}}(\theta; \rho_e)$  defined in Theorem 4 to form a surrogate likelihood  $\ell_B$ , and that for our regularizing function, we use an  $\ell_q$  norm  $\|\theta\|_q$  with  $q \geq 1$ . For the discrete Markov random fields to which the reweighted Bethe surrogate  $B_{\mathfrak{T}}(\theta; \rho_e)$  applies, the parameter space is simply  $\Omega = \mathbb{R}^d$ , so that we can calculate

$$\begin{aligned} R^*(\tau) &= \sup_{\theta \in \mathbb{R}^d} \{ \langle \theta, \tau \rangle - \|\theta\|_q \} \\ &= \begin{cases} 0 & \text{if } \|\tau\|_r \leq 1, \text{ where } 1/r = 1 - 1/q, \\ +\infty & \text{otherwise,} \end{cases} \end{aligned}$$

using the fact that  $\|\theta\|_q = \sup_{\|\tau\|_r \leq 1} \langle \tau, \theta \rangle$ . Thus, the dual RSL problem (7.30) takes the form

$$\sup_{\|\tau - \hat{\mu}\|_r \leq \lambda} [ -B^*(\tau) ] = \sup_{\|\tau - \hat{\mu}\|_r \leq \lambda} \left\{ \sum_{s \in V} H_s(\tau_s) - \sum_{(s,t) \in E} \rho_{st} I_{st}(\tau_{st}) \right\}.$$

Note that this problem has a very intuitive interpretation: it involves maximizing the Bethe entropy, subject to an  $\ell_r$ -norm constraint on the distance to the empirical mean parameters  $\hat{\mu}$ . As pointed out by various authors [60, 68], a similar dual interpretation also exists for regularized exact maximum likelihood. The choice of regularization translates into a certain uncertainty or robustness constraint in the dual domain. For instance, given  $\ell_2$  regularization ( $q = 2$ ) on  $\theta$ , we have  $r = 2$  so that the dual regularization is also  $\ell_2$ -based. On the other hand, for  $\ell_1$ -regularization, the corresponding dual norm is  $\ell_\infty$ , so that we have box constraints in the dual domain. ♣

# 8

---

## Integer programming, max-product, linear programming and conic relaxations

---

Thus far, the bulk of our discussion has focused on variational methods that, when applied to a given model  $p_\theta$ , yield approximations to the mean parameters  $\mu = \mathbb{E}_\theta[\phi(X)]$ , as well as to the log partition or cumulant function  $A(\theta)$ . In this section, we turn our attention to a related but distinct problem—namely, that of computing *a mode or most probable configuration* associated with the distribution  $p_\theta(x)$ . It turns out that the mode problem has a variational formulation in which the set  $\mathcal{M}$  of mean parameters once again plays a central role. More specifically, since mode computation is of significant interest for discrete random vectors, the marginal polytopes  $\mathbb{M}(G)$  play a central role through much of this chapter.

### 8.1 Variational formulation of computing modes

Given a member  $p_\theta$  of some exponential family, the mode computation problem refers to computing an element  $x^*$  that belongs to the set

$$\arg \max_{x \in \mathcal{X}^m} p_\theta(x) := \{y \in \mathcal{X}^m \mid p_\theta(y) \geq p_\theta(x) \quad \forall x \in \mathcal{X}^m\}. \quad (8.1)$$

We assume that at least one mode exists, so that this  $\arg \max$  set is non-empty. Moreover, it is convenient for development in the sequel to

note the equivalence

$$\arg \max_{x \in \mathcal{X}^m} p_\theta(x) = \arg \max_{x \in \mathcal{X}^m} \langle \theta, \phi(x) \rangle, \quad (8.2)$$

which follows from the exponential form of the density  $p_\theta(x)$ , and the fact that the cumulant function  $A(\theta)$  does not depend on  $x$ .

It is not immediately obvious how the MAP problem (8.2) is related to the variational principle from Theorem 2. As it turns out, the link lies in the notion of the “zero-temperature” limit. We begin by providing intuition for the more formal result to follow. From its definition and Theorem 2, for any parameter vector  $\theta \in \Omega$ , the cumulant function has the following two representations:

$$\begin{aligned} A(\theta) &:= \log \int \exp \{ \langle \theta, \phi(x) \rangle \} \nu(dx) \\ &= \sup_{\mu \in \mathcal{M}} \{ \langle \theta, \mu \rangle - A^*(\mu) \}. \end{aligned} \quad (8.3)$$

Now suppose that we rescale the canonical parameter  $\theta$  by some scalar  $\beta > 0$ . For the sake of this argument, let us assume that  $\beta\theta \in \Omega$  for all  $\beta > 0$ . Such a rescaling will put more weight, in a relative sense, on regions of the sample space  $\mathcal{X}^m$  for which  $\langle \theta, \phi(x) \rangle$  is large. Ultimately, as  $\beta \rightarrow +\infty$ , probability mass should remain only on configurations  $x^*$  in the set  $\arg \max_x \langle \theta, \phi(x) \rangle$ .

This intuition suggests that the behavior of the function  $A(\beta\theta)$  should have a close connection to the problem of computing  $\max_x \langle \theta, \phi(x) \rangle$ . Since  $A(\beta\theta)$  may diverge as  $\beta \rightarrow +\infty$ , it is most natural to consider the limiting behavior of the scaled quantity  $A(\beta\theta)/\beta$ . More formally, we state the following:

**Theorem 5.** *For all  $\theta \in \Omega$ , the problem of mode computation has the following alternative representations:*

$$\max_{x \in \mathcal{X}^m} \langle \theta, \phi(x) \rangle = \max_{\mu \in \mathcal{M}} \langle \theta, \mu \rangle, \quad \text{and} \quad (8.4a)$$

$$\lim_{\beta \rightarrow +\infty} \frac{A(\beta\theta)}{\beta} = \max_{x \in \mathcal{X}^m} \langle \theta, \phi(x) \rangle \quad (8.4b)$$

We provide the proof of this result in Appendix B.5; here we elaborate on its connections to the general variational principle (3.45), stated

in Theorem 2. As a particular consequence of the general variational principle, we have

$$\begin{aligned} \lim_{\beta \rightarrow +\infty} \frac{A(\beta\theta)}{\beta} &= \lim_{\beta \rightarrow +\infty} \frac{1}{\beta} \sup_{\mu \in \mathcal{M}} \{ \langle \beta\theta, \mu \rangle - A^*(\mu) \} \\ &= \lim_{\beta \rightarrow +\infty} \sup_{\mu \in \mathcal{M}} \left\{ \langle \theta, \mu \rangle - \frac{1}{\beta} A^*(\mu) \right\}. \end{aligned}$$

In this context, one implication of Theorem 5 is that the order of the limit over  $\beta$  and the supremum over  $\mu$  can be exchanged. It is important that such an exchange is *not* correct in general; in this case, justification is provided by the convexity of  $A^*$  (see Appendix B.5 for details).

As with the variational principle of Theorem 2, only certain special cases of the variational problem  $\max_{\mu \in \bar{\mathcal{M}}} \langle \theta, \mu \rangle$  can be solved exactly in a computationally efficient manner. The case of tree-structured Markov random fields on discrete random variables, discussed in Section 8.2, is one such case; another important exactly solvable case is the Gauss Markov random field, discussed in Section 8.3 and Appendix C.3. With reference to other (typically intractable) models, since the objective function itself is linear, the sole source of difficulty is the set  $\mathcal{M}$ , and in particular, the complexity of characterizing it with a small number of constraints. In the particular case of discrete Markov random fields (say associated with a graph  $G$ ), the set  $\mathcal{M}$  corresponds to a marginal polytope, as discussed previously in Chapters 3 and 4. For a discrete random vector, the mode-finding problem  $\max_{x \in \mathcal{X}^m} \langle \theta, \phi(x) \rangle$  is an integer program (IP), since it involves searching over a finite discrete space. An implication of Theorem 5, when specialized to this setting, is that this IP is equivalent to a linear program over the marginal polytope. Since integer programming problems are NP-hard in general, this equivalence underscores the inherent complexity of marginal polytopes.

This type of transformation—namely, from an integer program to a linear program over the convex hull of its solutions—is a standard technique in integer programming and combinatorial optimization [e.g., 22, 101]. The field of polyhedral combinatorics [66, 101, 177, 207] is devoted to understanding the structure of such polytopes arising from various classes of discrete problems. As we describe in this section, the perspective of graphical models provides additional insight into the

structure of such polytopes.

## 8.2 Max-product and linear programming on trees

As discussed in Section 4, an important class of exactly solvable models are Markov random fields defined on tree-structured graphs, say  $T = (V, E)$ . Our earlier discussion focused on the computation of mean parameters and the cumulant function, showing that the sum-product algorithm can be understood as an iterative algorithm for exactly solving the Bethe variational problem on trees [263, 264]. In this section, we discuss the parallel link between the ordinary max-product algorithm and linear programming [240].

For a discrete MRF on the tree, the set  $\mathcal{M}$  is given by the marginal polytope  $\mathbb{M}(T)$ , whose elements consist of a marginal probability vector  $\mu_s(\cdot)$  for each node, and joint probability matrix  $\mu_{st}(\cdot, \cdot)$  for each edge  $(s, t) \in E$ . Recall from Proposition 4 that for a tree, the polytope  $\mathbb{M}(T)$  is equivalent to the constraint set  $\mathbb{L}(T)$ , and so has a compact description in terms of non-negativity, local normalization, and edge marginalization conditions. Using this fact, we now show how the mode-finding problem for a tree-structured problem can be reformulated as a simple linear program (LP) over the set  $\mathbb{L}(G)$ .

Using  $\mathbb{E}_{\mu_s}[\theta_s(x_s)] := \sum_{x_s} \mu_s(x_s)\theta_s(x_s)$  to expectation under the marginal distribution  $\mu_s$  and with a similar definition for  $\mathbb{E}_{\mu_{st}}$ , let us define the cost function

$$\langle \tau, \theta \rangle := \sum_{s \in V} \mathbb{E}_{\mu_s}[\theta_s(x_s)] + \sum_{(s,t) \in E} \mathbb{E}_{\mu_{st}}[\theta_{st}(x_s, x_t)].$$

With this notation, Theorem 5 implies that the MAP problem for a tree is equivalent to

$$\max_{x \in \mathcal{X}^m} \left[ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \right] = \max_{\mu \in \mathbb{L}(T)} \langle \tau, \theta \rangle. \quad (8.5)$$

The left-hand side is the standard representation of the MAP problem as an integer program, whereas the right-hand side—an optimization problem over the polytope  $\mathbb{L}(G)$  with a linear objective function—is a linear program.



We now have two links: first, the equivalence between the original mode-finding integer program and the LP (8.5), and secondly, via Theorem 5, a connection between this linear program and the zero-temperature limit of the Bethe variational principle. These links raise the intriguing possibility: *is there a precise connection between the max-product algorithm and the linear program (8.5)?* The max-product algorithm—like its relative the sum-product algorithm—is a distributed graph-based algorithm, which operates by passing a message  $M_{st}$  along each direction of each edge of the tree. For discrete random variables  $X_t \in \{0, 1, \dots, r-1\}$ , each message is a  $r$ -vector, and the messages are updated according to recursion

$$M_{ts}(x_s) \leftarrow \kappa \max_{x_t \in \mathcal{X}_t} \left[ \exp \{ \theta_{st}(x_s, x_t) + \theta_t(x_t) \} \prod_{u \in N(t) \setminus s} M_{ut}(x_t) \right]. \quad (8.6)$$

The following result [240] gives an affirmative answer to the question above: for tree-structured graphs, the max-product updates (8.6) are a Lagrangian method for solving the dual of the linear program (8.5). This result can be viewed as the max-product analog to the connection between sum-product and the Bethe variational problem on trees, and its exactness for computing marginals, as summarized in Theorem 3.

To set up the Lagrangian, for each  $x_s \in \mathcal{X}_s$ , let  $\lambda_{st}(x_s)$  be a Lagrange multiplier associated with the marginalization constraint  $C_{ts}(x_s) = 0$ , where

$$C_{ts}(x_s) := \mu_s(x_s) - \sum_{x_t} \mu_{st}(x_s, x_t). \quad (8.7)$$

Let  $\mathbb{N} \subset \mathbb{R}^d$  be the set of  $\mu$  that are non-negative and appropriately normalized:

$$\mathbb{N} := \left\{ \mu \in \mathbb{R}^d \mid \mu \geq 0, \sum_{x_s} \mu_s(x_s) = 1, \sum_{x_s, x_t} \mu_{st}(x_s, x_t) = 1 \right\} \quad (8.8)$$

With this notation, we have the following result:

**Proposition 9 (Max-product and LP duality).** *Consider the dual function  $Q$  defined by the following partial Lagrangian formulation of*

the tree-structured LP (8.5):

$$\begin{aligned} \mathcal{Q}(\lambda) &:= \max_{\mu \in \mathbb{N}} \mathcal{L}(\mu; \lambda), \\ \mathcal{L}(\mu; \lambda) &:= \langle \theta, \mu \rangle + \sum_{(s,t) \in E(T)} \left[ \sum_{x_s} \lambda_{ts}(x_s) C_{ts}(x_s) + \sum_{x_t} \lambda_{st}(x_t) C_{st}(x_t) \right]. \end{aligned}$$

For any fixed point  $M^*$  of the max-product updates (8.6), the vector  $\lambda^* := \log M^*$ , where the logarithm is taken elementwise, is an optimal solution of the dual problem  $\min_{\lambda} \mathcal{Q}(\lambda)$ .

*Proof.* We begin by deriving a convenient representation of the dual function  $\mathcal{Q}$ . First, let us convert the tree to a directed version by first designating some node  $r \in V$  as the root, and then directing all the edges from parent to child  $t \rightarrow s$ . With regard to this rooted tree, the objective function  $\langle \theta, \mu \rangle$  has the alternative decomposition:

$$\sum_{x_r} \mu_r(x_r) \theta_r(x_r) + \sum_{t \rightarrow s} \sum_{x_t, x_s} \mu_{st}(x_s, x_t) [\theta_{st}(x_s, x_t) + \theta_s(x_s)].$$

With this form of the cost function, the dual function can be put into the form:

$$\begin{aligned} \mathcal{Q}(\lambda) &:= \max_{\mu \in \mathbb{N}} \left\{ \sum_{x_r} \mu_r(x_r) \nu_s(x_s) \right. \\ &\quad \left. + \sum_{t \rightarrow s} \sum_{x_t, x_s} \mu_{st}(x_s, x_t) [\nu_{st}(x_s, x_t) - \nu_t(x_t)] \right\}, \quad (8.9) \end{aligned}$$

where the quantities  $\nu_s$  and  $\nu_{st}$  are defined in terms of  $\lambda$  and  $\theta$  as:

$$\nu_t(x_t) = \theta_t(x_t) + \sum_{u \in N(t)} \lambda_{ut}(x_t), \quad \text{and} \quad (8.10a)$$

$$\begin{aligned} \nu_{st}(x_s, x_t) &= \theta_{st}(x_s, x_t) + \theta_s(x_s) + \theta_t(x_t) \\ &\quad + \sum_{u \in N(s) \setminus t} \lambda_{us}(x_s) + \sum_{u \in N(t) \setminus s} \lambda_{ut}(x_t). \quad (8.10b) \end{aligned}$$

Taking the maximum over  $\mu \in \mathbb{N}$  in equation (8.9) yields that the dual function has the explicit form

$$\mathcal{Q}(\lambda) = \max_{x_r} \nu_r(x_r) + \sum_{t \rightarrow s} \max_{x_s, x_t} [\nu_{st}(x_s, x_t) - \nu_t(x_t)]. \quad (8.11)$$

Using the representation (8.11), we now have the necessary ingredients to make the connection to the max-product algorithm. Given a vector of messages  $M$  in the max-product algorithm, we use it to define a vector of Lagrange multipliers via  $\lambda = \log M$ , where the logarithm is taken elementwise. With a bit of algebra, it can be seen that a message vector  $M^*$  is a fixed point of the max-product updates (8.6) if and only if the associated  $\nu_s^*$  and  $\nu_{st}^*$ , as defined by  $\lambda^* := \log M^*$ , satisfy the *edgewise consistency* condition

$$\max_{x_s} \nu_{st}^*(x_s, x_t) = \nu_t^*(x_t) + C_{st} \quad (8.12)$$

for all  $x_t \in \mathcal{X}_t$ , where  $C_{st}$  is a constant independent of  $x$ . We now show that any such  $\lambda^*$  is a dual optimal solution.

We first claim that under the edgewise consistency condition on a tree-structured graph, we can always find at least one configuration  $x^*$  that satisfies

$$x_s^* \in \arg \max_{x_s} \nu_s^*(x_s) \quad \forall s \in V, \text{ and} \quad (8.13a)$$

$$(x_s^*, x_t^*) \in \arg \max_{x_s, x_t} \nu_{st}^*(x_s, x_t) \quad \forall (s, t) \in E. \quad (8.13b)$$

Indeed, such a configuration can be constructed recursively as follows. First, for the root node  $r$ , choose  $x_r^*$  to achieve the maximum  $\max_{x_r} \nu_r^*(x_r)$ . Second, for each child  $t \in N(r)$ , choose  $x_t$  to maximize  $\nu_{tr}(x_t, x_r^*)$ , and then iterate the process. The edgewise consistency (8.12) guarantees that any configuration  $x^*$  constructed in this way satisfies the conditions (8.13a) and (8.13b).

The edgewise consistency condition (8.12) also guarantees the following equalities:

$$\begin{aligned} \max_{x_s, x_t} [\nu_{st}^*(x_s, x_t) - \nu_t^*(x_t)] &= \max_{x_t} [\nu_t^*(x_t) + C_{st} - \nu_t^*(x_t)] \\ &= C_{st} \\ &= \nu_{st}^*(x_s^*, x_t^*) - \nu_t^*(x_t^*). \end{aligned}$$

Combining these relations yields the following expression for the dual value at  $\lambda^*$ :

$$\mathcal{Q}(\lambda^*) = \nu_r^*(x_r^*) + \sum_{t \rightarrow s} [\nu_{st}^*(x_s^*, x_t^*) - \nu_t^*(x_t^*)]$$

Next, applying the definition (8.10) of  $\nu^*$  and simplifying, we find that

$$\mathcal{Q}(\lambda^*) = \theta_r(x_r^*) + \sum_{t \rightarrow s} [\theta_{st}(x_s^*, x_t^*) + \theta_s(x_s^*)].$$

Now consider the primal solution defined by  $\mu_s^*(x_s) := \mathbb{I}_{x_s^*}[x_s]$  and  $\mu_{st}^*(x_s, x_t) = \mathbb{I}_{x_s^*}[x_s] \mathbb{I}_{x_t^*}[x_t]$ , where  $\mathbb{I}_{x_s^*}[x_s]$  is an indicator function for the event  $\{x_s = x_s^*\}$ . It is clear that  $\mu^*$  is primal feasible; moreover, the primal cost is equal to

$$\begin{aligned} \sum_{x_r} \mu_r^*(x_r) \theta_r(x_r) + \sum_{t \rightarrow s} \sum_{x_t, x_s} \mu_{st}^*(x_s, x_t) [\theta_{st}(x_s, x_t) + \theta_s(x_s)] \\ = \theta_r(x_r^*) + \sum_{t \rightarrow s} [\theta_{st}(x_s^*, x_t^*) + \theta_s(x_s^*)], \end{aligned}$$

which is precisely equal to  $\mathcal{Q}(\lambda^*)$ . Therefore, by strong duality for linear programs [22], the pair  $(\mu^*, \lambda^*)$  is primal-dual optimal.  $\square$

**Remark:** A careful examination of the proof of Proposition 9 shows that several steps rely heavily on the fact that the underlying graph is a tree. In fact, the corresponding result for a graph with cycles *fails* to hold, as we discuss at length in Section 8.4.

### 8.3 Max-product for Gaussians and other convex problems

Multivariate Gaussians are another special class of Markov random fields for which there are various correctness guarantees associated with the max-product (or min-sum) algorithm. Given an undirected graph  $G = (V, E)$ , consider a Gaussian MRF in exponential form

$$\begin{aligned} p_{(\theta, \Theta)}(x) &= \exp \{ \langle \theta, x \rangle + \langle \langle \Theta, xx^T \rangle \rangle - A(\theta, \Theta) \} \\ &= \exp \left\{ \sum_{s \in V} \theta_s x_s + \sum_{s, t} \Theta_{st} x_s x_t - A(\theta, \Theta) \right\}. \end{aligned} \quad (8.14)$$

Recall that the graph structure is reflected in the sparsity of the matrix  $\Theta$ , with  $\Theta_{uv} = 0$  whenever  $(u, v) \notin E$ .

In Appendix C.3, we describe how taking the zero-temperature limit for a multivariate Gaussian, according to Theorem 5, leads to either a quadratic program, or a semidefinite program. There are various ways

in which these convex programs could be solved; here we discuss known results for the application of the max-product algorithm.

To describe the max-product in application to a multivariate Gaussian, we first need to convert the model into pairwise form. Let us define potential functions of the form

$$\gamma_s(x_s) = \theta_s x_s + \theta_s x_s^2 \quad (8.15a)$$

$$\gamma_{st}(x_s, x_t) = \gamma_{ss;t} x_s^2 + 2\Theta_{st} x_s x_t + \gamma_{tt;s} x_t^2, \quad (8.15b)$$

where the free parameters  $\{\gamma_s, \gamma_{ss;t}\}$  must satisfy the constraint

$$\gamma_s + \sum_{t \in N(s)} \gamma_{ss;t} = \Theta_{ss}.$$

Under this condition, the multivariate Gaussian density (8.14) can equivalently be expressed as the pairwise Markov random field

$$p_\gamma(x) \propto \exp \left\{ \sum_{s \in V} \gamma_s(x_s) + \sum_{(s,t) \in E} \gamma_{st}(x_s, x_t) \right\}. \quad (8.16)$$

With this set-up, we can now describe the max-product message-passing updates. For graphical models involving continuous random variables, each message  $M_{ts}(x_s)$  is a real-valued function of the indeterminate  $x_s$ . These functional messages are updated according to the usual max-product recursion (8.6), with the potentials  $\theta_s$  and  $\theta_{st}$  in equation (8.6) replaced by their  $\gamma_s$  and  $\gamma_{st}$  analogues from the pairwise factorization (8.16). In general continuous models, it is computationally challenging to represent and compute with these functions, as they are infinite-dimensional quantities. An attractive feature of the Gaussian problems these messages can be compactly parameterized; in particular, assuming a scalar Gaussian random variable  $X_s$  at each node, any message must be an exponentiated-quadratic function of its argument, of the form  $M_{ts}(x_s) \propto \exp(ax_s + bx_s^2)$ . Consequently, max-product message-passing updates (8.6) can be efficiently implemented with one recursion for the mean term (number  $a$ ), and a second recursion for the variance component (see the papers [245, 257] for further details).

For Gaussian max-product applied to a tree-structured problem, the updates are guaranteed to converge, and compute both the correct means  $\mu_s = \mathbb{E}[X_s]$  and variances  $\sigma_s^2 = \mathbb{E}[X_s^2] - \mu_s^2$  at each node [257].

For general graphs with cycles, the max-product algorithm is no longer guaranteed to converge. However, Weiss and Freeman [245] and independently Rusmevichientong and van Roy [200] showed that if Gaussian max-product (or equivalently, Gaussian sum-product) converges, then the fixed point specifies the correct Gaussian means  $\mu_s$ , but the estimates of the node variances  $\sigma_s^2$  need not be incorrect. The correctness of Gaussian max-product for mean computation also follows as a consequence of the reparameterization properties of the sum- and max-product algorithms [237, 239]. Weiss and Freeman [245] also showed that Gaussian max-product converges for arbitrary graphs if the precision matrix ( $-\Theta$  in our notation) satisfies a certain diagonal dominance condition. This sufficient condition for convergence was substantially tightened in later work by Malioutov et al. [158], using the notions of walk-summability and pairwise normalizability; see also the paper [173] for further refinements. Moallemi and van Roy [174] consider the more general problem of maximizing an arbitrary function of the form (8.16), where the potentials  $\{\gamma_s, \gamma_{st}\}$  are required only to define a suitably convex function. The quadratic functions (8.15) for the multivariate Gaussian are a special case of this general set-up. For this family of pairwise separable convex programs, Moallemi and van Roy [174] established convergence of the max-product updates under a certain scaled diagonal dominance condition.

#### 8.4 First-order LP relaxation and reweighted max-product

In this section, we return to Markov random fields with discrete random variables, for which the mode-finding problem is an integer program, as opposed to the quadratic program that arises in the Gaussian case. For a general graph with cycles, this discrete mode-finding problem is known to be computationally intractable, since it includes as special cases many problems known to be NP-complete, among them MAX-CUT and related satisfiability problems [128]. Given the computational difficulties associated with exact solutions, it is appropriate to consider approximate algorithms.

An important class of approximate methods for solving integer and combinatorial optimization problems are based on *linear programming*

*relaxations.* The basic idea is to approximate the convex hull of the solution set by a set of linear constraints, and solve the resulting linear program. If the obtained solution is integral, then the LP relaxation is tight, whereas in other cases, the obtained solution may be fractional, corresponding to looseness of the relaxation. Frequently, LP relaxations are developed in a manner tailored to specific subclasses of combinatorial problems [177, 231].

In this section, we discuss the linear programming (LP) relaxation that emerges from the Bethe variational principle, or its zero-temperature limit. This LP relaxation, when specialized to particular combinatorial problems—among them node cover, independent set, matching, and satisfiability problems—recovers various classical LP relaxations as special cases. Finally, we discuss connections between this LP relaxation and max-product message-passing. For general MRFs, max-product itself is *not* a method for solving this LP, as we demonstrate with a concrete counterexample [240]. However, it turns out that suitably reweighted versions of max-product are intimately linked to this LP relaxation, which provides an entry point to an ongoing line of research on LP relaxations and message-passing algorithms on graphs.

#### 8.4.1 Basic properties of first-order LP relaxation

For a general graph, the set  $\mathbb{L}(G)$  no longer provides an exact characterization of the marginal polytope  $\mathbb{M}(G)$ , but is always guaranteed to outer bound it. Consequently, the following inequality is valid for any graph:

$$\max_{x \in \mathcal{X}^m} \langle \theta, \phi(x) \rangle = \max_{\mu \in \mathbb{M}(G)} \langle \theta, \mu \rangle \leq \max_{\tau \in \mathbb{L}(G)} \langle \theta, \tau \rangle. \quad (8.17)$$

Since the relaxed constraint set  $\mathbb{L}(G)$  is a polytope, the right-hand side of equation (8.17) is a linear program, which we refer to as the *first-order LP relaxation*.<sup>1</sup> Most importantly, the right-hand side is a linear program, where the number of constraints defining  $\mathbb{L}(G)$  grows only linearly in the graph size, so that it can be solved in polynomial-time for any graph [206].

<sup>1</sup>The term “first-order” refers to its status as the first in a natural hierarchy of relaxations, based on the treewidth of the underlying graph, as discussed at more length in Section 8.5.

Special cases of this LP relaxation (8.17) have studied in past and ongoing work by various authors, including the special cases of  $\{0, 1\}$ -quadratic programs [102], metric labeling with Potts models [132, 45], error-control coding problems [73, 75, 234, 223, 47, 59], independent set problems [177, 202], and various types of matching problems [13, 113, 201]. We begin by discussing some generic properties of the LP relaxation, before discussing some of these particular examples in more detail.

By standard properties of linear programs [206], the optimal value of the relaxed LP must be attained by at least one vertex of the polytope  $\mathbb{L}(G)$ . We say that a vertex of  $\mathbb{L}(G)$  is *integral* if all of its components are zero or one, and *fractional* otherwise. By definition, any vertex of  $\mathbb{M}(G)$  is of the form  $\mu_y := \phi(y)$ , where  $y \in \mathcal{X}^m$  is a fixed configuration. Recall that in the canonical overcomplete representation (3.34), the sufficient statistic vector  $\phi$  consists of  $\{0, 1\}$ -valued indicator functions, so that  $\mu_y = \phi(y)$  is vector with  $\{0, 1\}$  components. The following result specifies the relation between vertices of  $\mathbb{M}(G)$  and those of  $\mathbb{L}(G)$ :

**Proposition 10.** *The vertices of  $\mathbb{L}(G)$  and  $\mathbb{M}(G)$  are related as follows:*

- (a) *All the vertices of  $\mathbb{M}(G)$  are integral, and each one is also a vertex of  $\mathbb{L}(G)$ .*
- (b) *For any graph with cycles,  $\mathbb{L}(G)$  also includes additional vertices with fractional elements that lie strictly outside  $\mathbb{M}(G)$ .*

*Proof.* (a) Any vertex of  $\mathbb{M}(G)$  is of the form  $\phi(y)$ , for some configuration  $y \in \mathcal{X}^m$ . Each of these vertices has 0 – 1 components, and so is integral. In order to show that  $\phi(y)$  is also a vertex of  $\mathbb{L}(G)$ , it suffices [22] to show that there are  $d$  constraints of  $\mathbb{L}(G)$  that are active at  $\phi(y)$ , and are also linearly independent. For any  $y \in \mathcal{X}^m$ , we have  $\mathbb{I}_k(x_s) = 0$  for all  $k \in \mathcal{X} \setminus \{y_s\}$ , and  $\mathbb{I}_k(x_s)\mathbb{I}_l(x_t) = 0$  for all  $(k, l) \in (\mathcal{X} \times \mathcal{X}) \setminus \{y_s, y_t\}$ . All of these active inequality constraints are linearly independent, and there are a total of  $d' = (r-1)m + (r^2-1)|E|$ . All of the normalization and marginalization constraints are also satisfied by the vector  $\mu_y = \phi(y)$ , but not all of them are linearly independent (when added to the active inequality constraints). However,



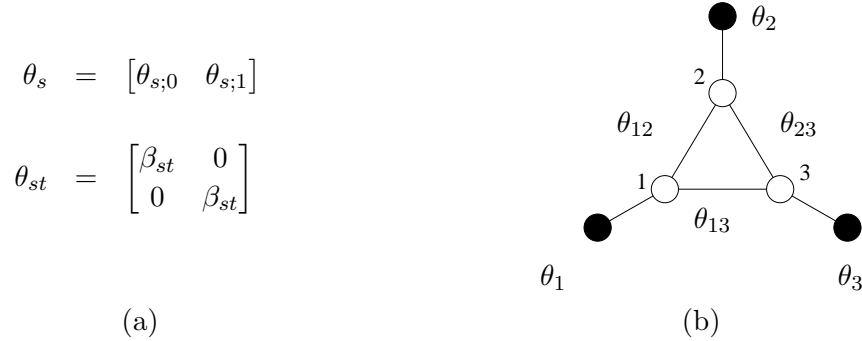
we can add the normalization constraints for each  $s = 1, \dots, m$  and for each  $(s, t) \in E$ , while still preserving linear independence. Adding these  $m + |E|$  equality constraints to the  $d'$  inequality constraints yields a total of  $d$  linearly independent constraints of  $\mathbb{L}(G)$  that are satisfied by  $\mu_J$ , so that it is a vertex. (b) Example 36 provides a constructive procedure for demonstrating fractional vertices for  $\mathbb{L}(G)$  for graphs with cycles.  $\square$

The distinction between fractional and integral vertices is crucial, because it determines whether or not the LP relaxation (8.17) specified by  $\mathbb{L}(G)$  is tight. In particular, there are only two possible outcomes to solving the relaxation:

- (a) the optimum is attained at a vertex of  $\mathbb{M}(G)$ , in which case the upper bound in equation (8.17) is tight, and a mode can be obtained.
- (b) the optimum is attained only at one or more fractional vertices of  $\mathbb{L}(G)$ , which lie strictly outside  $\mathbb{M}(G)$ . In this case, the upper bound of equation (8.17) is loose, and the optimal solution to the LP relaxation does not specify the optimal configuration. In this case, one can imagine various types of rounding procedures for producing near-optimal solutions [231].

When the graph has cycles, it is possible to explicitly construct a fractional vertex of the polytope  $\mathbb{L}(G)$ .

**Example 36 (Fractional vertices of  $\mathbb{L}(G)$ ).** Here we explicitly construct a fractional vertex for the simplest problem on which the relaxation fails to be tight in general: a binary random vector  $X \in \{0, 1\}^3$  following an Ising model on the complete graph  $K_3$ . Consider the canonical parameter  $\theta$  shown in matrix form in Figure 8.1(a). When  $\beta_{st} < 0$ , then configurations with  $x_s \neq x_t$  are favored, so that the interaction is repulsive. In contrast, when  $\beta_{st} > 0$ , the interaction is attractive, because it favors configurations with  $x_s = x_t$ . When  $\beta_{st} > 0$  for all  $(s, t) \in E$ , it can be shown [135] that the first-order LP relaxation (8.17) is tight, for any choice of the single node parameters  $\{\theta_s, s \in V\}$ . In contrast, when  $\beta_{st} < 0$  for all edges, then there are choices of  $\theta_s, s \in V$



**Fig. 8.1.** The smallest graph  $G = (V, E)$  on which the relaxation (8.17) can fail to be tight. For  $\beta_{st} \geq 0$  for all  $(s, t) \in E$ , the relaxation is tight for any choice of  $\theta_s, s \in V$ . On the other hand, if  $\beta_{st} < 0$  for all edges  $(s, t)$ , the relaxation will fail for certain choices of  $\theta_s, s \in V$ .

for which the relaxation breaks down.

The following canonical parameter corresponds to a direction for which the relaxation (8.17) is *not* tight, and hence exposes a fractional vertex. First choose  $\theta_s = [0 \ 0]^T$  for  $s = 1, 2, 3$ , and then set  $\beta_{st} = \beta < 0$  for all edges  $(s, t)$ , and use these values to define the pairwise potentials  $\theta_{st}$  via the construction in Figure 8.1(a). Observe that for any configuration  $x \in \{0, 1\}^3$ , we must have  $x_s \neq x_t$  for at least one edge  $(s, t) \in E$ . Therefore, any  $\mu \in \mathbb{M}(G)$  must place non-zero mass on at least one term of  $\theta$  involving  $\beta$ , whence  $\max_{\mu \in \mathbb{M}(G)} \langle \theta, \mu \rangle < 0$ . In fact, this optimal value is exactly equal to  $\beta < 0$ . On the other, consider the pseudomarginal  $\tau^* \in \mathbb{L}(G)$  formed by the singleton and pairwise pseudomarginals defined as follows:

$$\begin{aligned} \tau_s^* &:= [0.5 \quad 0.5]^T \quad \text{for } s \in V, \text{ and} \\ \tau_{st}^* &= \begin{bmatrix} 0 & 0.5 \\ 0.5 & 0 \end{bmatrix} \quad \text{for } (s, t) \in E. \end{aligned}$$

Observe that  $\langle \theta, \tau^* \rangle = 0$ . Since  $\theta_\alpha \leq 0$  for all elements  $\alpha$ , this value is the optimum of  $\langle \theta, \tau \rangle$  over  $\mathbb{L}(G)$ , thereby showing that the relaxation (8.17) is *not* tight. Indeed, we have  $\max_{\mu \in \mathbb{M}(G)} \langle \theta, \mu \rangle = \beta < 0$ .

Finally, to establish that  $\tau^*$  is a vertex of  $\mathbb{L}(G)$ , we will show that  $\langle \theta, \tau \rangle < 0$  for all  $\tau \neq \tau^*$ . If  $\langle \theta, \tau \rangle = 0$ , then for all  $(s, t) \in E$  the

pairwise pseudomarginals must be of the form

$$\tau_{st} := \begin{bmatrix} 0 & \alpha_{st} \\ 1 - \alpha_{st} & 0 \end{bmatrix}$$

for some  $\alpha_{st} \in [0, 1]$ . Enforcing the marginalization constraints on these pairwise pseudomarginals yields the constraints  $\alpha_{12} = \alpha_{13} = \alpha_{23}$  and  $1 - \alpha_{12} = \alpha_{23}$ , whence  $\alpha_{st} = 0.5$  is the only possibility. Therefore, we conclude that the pseudomarginal vector  $\tau^*$  is a fractional vertex of the polytope  $\mathbb{L}(G)$ .  $\clubsuit$

Given the possibility of fractional vertices in the first-order LP relaxation (8.5), it is natural to ask the question: do fractional solutions yield partial information about the set of optimal solutions to the original integer program? One way in which to formalize this question is through the notion of persistence [102]. In particular, letting  $\mathcal{O}^* := \arg \max_{x \in \mathcal{X}^m} \langle \theta, \phi(x) \rangle$  denote the set of optima to an integer program, we have:

**Definition 2.** Given a fractional solution  $\tau$  to the LP relaxation (8.5), let  $I \subset V$  represent the subset of vertices for which  $\tau_s$  has only integral elements, say fixing  $x_s = x_s^*$  for all  $s \in I$ . The fractional solution is *strongly persistent* if any optimal integral solution  $y^* \in \mathcal{O}^*$  satisfies  $y_s^* = x_s^*$  for all  $s \in I$ . The fractional solution is *weakly persistent* if there exists some  $y^* \in \mathcal{O}^*$  such that  $y_s^* = x_s^*$  for all  $s \in I$ .

Thus, any persistent fractional solution (whether weak or strong) can be used to fix a subset of elements in the integer program, while still being assured that there exists an optimal integral solution that is consistent. Strong persistency ensures that no candidate solutions are eliminated by the fixing procedure. Hammer et al. [102] studied the roof-dual relaxation for binary quadratic programs, an LP relaxation which is equivalent to specializing the first-order LP to binary variables, and proved the following result:

**Proposition 11.** *Suppose that the first-order LP relaxation (8.5) is applied to the binary quadratic program*

$$\max_{x \in \{0,1\}^m} \left\{ \sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t \right\}. \quad (8.18)$$

*Then any fractional solution is strongly persistent.*

As will be discussed at more length in Example 39, the class of binary QPs includes as special cases various classical problems from the combinatorics literature, among them MAX-2SAT, independent set, MAX-CUT, and vertex cover. For all of these problems, then, the first-order LP relaxation is strongly persistent. Unfortunately, this strong persistency fails to extend to the first-order relaxation (8.5) with non-binary variables; see the discussion following Example 38 for a counterexample.

#### 8.4.2 Connection to max-product message-passing

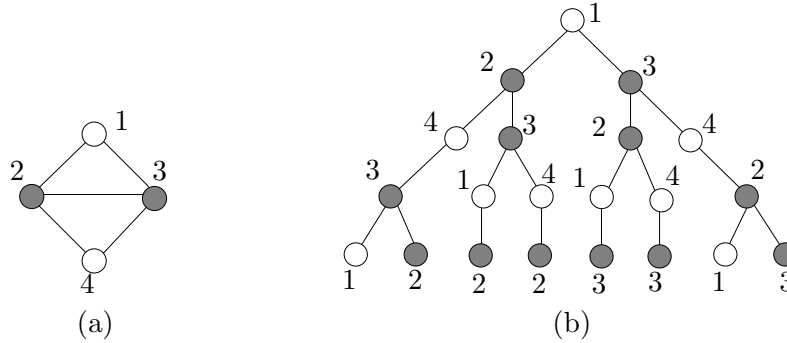
In analogy to the general connection between the Bethe variational problem and the sum-product algorithm (see Section 4), one might postulate that Proposition 9 could be extended to graphs with cycles—specifically, that the max-product algorithm solves the dual of the tree-based relaxation (8.17). In general, this conjecture is false, as the following counterexample [240] shows:

**Example 37 (Max-product does not solve the LP).** Consider the diamond graph  $G_{\text{dia}}$  shown in Figure 8.2, and suppose that we wish to maximize a cost function of the form

$$\alpha(x_1 + x_4) + \beta(x_2 + x_3) + \gamma \sum_{(s,t) \in E} \mathbb{I}[x_s \neq x_t]. \quad (8.19)$$

Here the maximization is over all binary vectors  $x \in \{0, 1\}^4$ , and  $\alpha, \beta$  and  $\gamma$  are parameters to be specified. By design, the cost function (8.19) is such that if we make  $\gamma$  sufficiently negative, then any optimal solution will either be  $0^4 := [0 \ 0 \ 0 \ 0]$  or  $1^4 := [1 \ 1 \ 1 \ 1]$ . As an extreme example, if we set  $\alpha = 0.31$ ,  $\beta = -0.30$ , and  $\gamma = -\infty$ , we see immediately that the optimal solution is  $1^4$ . (Note that setting  $\gamma = -\infty$  is equivalent to imposing the “hard-core” constraint that  $x_s = x_t$  for all  $(s, t) \in E$ .)

A classical way of studying the ordinary sum- and max-product algorithms, dating back to the work of Gallager [84] and Wiberg et al. [253], is via the computation tree associated with the message-passing updates. As illustrated in Figure 8.2(b), the computation tree



**Fig. 8.2.** (a) Simple diamond graph  $G_{\text{dia}}$ . (b) Associated computation tree after four rounds of message-passing. Max-product solves exactly the modified integer program defined by the computation tree.

is rooted at a particular vertex (1 in this case), and it tracks the paths of messages that reach this root node. In general, the  $(n + 1)^{\text{th}}$  level of tree includes all vertices  $t$  such that a path of length  $n$  joins  $t$  to the root. For the particular example shown in Figure 8.2(b), after one iteration, the root node 1 receives messages from nodes 2 and 3 (level 2 of the tree in panel (b)), and after two iterations, it receives messages from nodes 3 (via 2), and two messages from node 4, one via node 2 and the other via node 3, corresponding to the third level in panel (b), and so on.

The significance of this computation tree is based on the following observation: by definition, the decision of the max-product algorithm at the root node 1 after  $n$  iterations is optimal with respect to the *modified integer program* defined by the computation tree with  $n + 1$  levels. That is, the max-product decision will be  $\hat{x}_1 = 1$  if and only if the optimal configuration in the tree with  $x_1 = 1$  has higher probability than the optimal configuration with  $x_1 = 0$ . For the diamond graph under consideration and given the “hard-core” constraints imposed by setting  $\gamma = -\infty$ , the only two possible configurations in any computation tree are all-zeroes, or all-ones. Thus, the max-product reduces to comparing the total weight on the computation tree associated with all-ones to that associated with all-zeroes.

However, the computation tree in Figure 8.2(b) has a curious prop-

erty: due to the inhomogeneous node degrees—more specifically, with nodes 2 and 3 having three neighbors, and 1 and 4 having only two neighbors—nodes 2 and 3 receive a disproportionate representation in the computation tree. This fact is clear by inspection, and can be verified rigorously by setting up a simple recursion to compute the asymptotic appearance fractions of nodes 2 and 3 versus nodes 1 and 4. Doing so shows that nodes  $\{2, 3\}$  appear roughly  $\rho \approx 1.0893$  more frequently than nodes  $\{1, 4\}$ . As a consequence, the ordinary max-product algorithm makes its decision according to the threshold rule

$$\hat{x}_{\text{MP}} = \begin{cases} 1^4 & \text{if } \alpha + \rho\beta > 0 \\ 0^4 & \text{otherwise,} \end{cases} \quad (8.20)$$

whereas the correct decision rule is based on thresholding  $\alpha + \beta$ . Consequently, for any pair  $(\alpha, \beta)$  such that  $\alpha + \beta > 0$  but  $\alpha + \rho\beta < 0$ , the max-product algorithm outputs an incorrect configuration. For instance, setting  $\alpha = 0.31$  and  $\beta = -0.30$  yields one such counterexample, as discussed in Wainwright et al. [240]. Kulesza and Pereira [141] provide an detailed analysis of the max-product message-passing updates for this example, analytically deriving the updates and explicitly demonstrating convergence to incorrect configurations.

Thus far, we have demonstrated a simple problem for which the max-product algorithm fails. What is the connection to the LP relaxation? It turns out that the problem instance that we have constructed—in particular, with the hard core constraint  $\gamma = -\infty$ —generates an instance of super-modular maximization for binary variables [155]. In particular, a pairwise interaction  $\theta_{st}$  involving binary variables is said to be supermodular if  $\theta_{st}(1, 1) + \theta_{st}(0, 0) \geq \theta_{st}(1, 0) + \theta_{st}(0, 1)$ ; see Example 39 for further discussion of this property. It is known that the first-order LP relaxation is tight for maximizing any binary quadratic cost function—that is, of the form (8.19)—where the interactions between variables are supermodular [102, 135]. The cost function (8.19) is supermodular for all weights  $\gamma \leq 0$ , so that in particular, the first-order LP relaxation is tight for the cost function specified by  $(\alpha, \beta, \gamma) = (0.31, -0.30, -\infty)$ . However, as we have just shown, ordinary max-product fails for this problem, so it is not solving the LP relaxation. ♣

As discussed at more length in Section 8.4.4, for some problems with special combinatorial structure, max-product algorithm does solve the first-order LP relaxation. These instances include the problem of bipartite matching and weighted  $b$ -matching. However, establishing a general connection to the LP relaxation requires related but different message-passing algorithms, the topic to which we now turn.

### 8.4.3 Reweighted max-product and other modified message-passing schemes

In this section, we begin by presenting the tree-reweighted max-product updates [240], and describing their connection to the first-order LP relaxation (8.5). Recall from Theorem 4 the connection between the tree-reweighted Bethe variational problem (7.12), and the tree-reweighted sum-product updates (7.13). Note that the constraint set in the tree-reweighted Bethe variational problem is simply the polytope  $\mathbb{L}(G)$  defining the first-order LP relaxation. This fact suggests that there should be a connection between the “zero-temperature” limit of the tree-reweighted Bethe variational problem, and the first-order LP relaxation. In particular, recalling the convex surrogate  $B(\theta) = B_{\mathfrak{T}}(\theta; \rho_e)$  defined by the tree-reweighted Bethe variational problem from Theorem 4, let us consider the limit  $B(\beta\theta)/\beta$  as  $\beta \rightarrow +\infty$ . From the variational definition (7.12), we have

$$\lim_{\beta \rightarrow +\infty} \frac{B(\beta\theta)}{\beta} = \lim_{\beta \rightarrow +\infty} \left[ \sup_{\tau \in \mathbb{L}(G)} \left\{ \langle \theta, \tau \rangle - \frac{1}{\beta} B^*(\tau) \right\} \right].$$

As discussed previously, convexity allows us to exchange the order of the limit and supremum, so that we conclude that the zero-temperature limit of the convex surrogate  $B$  is simply the first-order LP relaxation (8.5).

Based on this intuition, it is natural to suspect that the *tree-reweighted max-product* algorithm should have a general connection to the first-order LP relaxation. In analogy to the TRW-sum-product

updates (7.13), the reweighted max-product updates take the form

$$M_{ts}(x_s) \leftarrow \kappa \max_{x'_t \in \mathcal{X}_t} \exp\left(\frac{1}{\rho_{st}} \theta_{st}(x_s, x'_t) + \theta_t(x'_t)\right) \left\{ \frac{\prod_{v \in N(t) \setminus s} [M_{vt}(x'_t)]^{\rho_{vt}}}{[M_{st}(x'_t)]^{(1-\rho_{ts})}} \right\}, \quad (8.21)$$

where  $\{\rho_{st}, (s, t) \in E\}$  are a collection of positive edge weights in the spanning tree polytope. See Theorem 4 on the reweighted sum-product algorithm and the accompanying discussion for more details on the choice of these edge weights.

As with the reweighted sum-product updates, these messages define a collection of *pseudo-max-marginals* of the form

$$\begin{aligned} \nu_s(x_s) &\propto \exp(\theta_s(x_s)) \prod_{t \in N(s)} M_{ts}(x_s), \\ \nu_s(x_s, x_t) &\propto \exp(\gamma_{st}(x_s, x_t)) \frac{\prod_{u \in N(s) \setminus t} [M_{us}(x_s)]^{\rho_{us}} \prod_{u \in N(t) \setminus s} [M_{ut}(x_t)]^{\rho_{ut}}}{[M_{ts}(x_s)]^{1-\rho_{st}} [M_{st}(x_t)]^{1-\rho_{st}}}, \end{aligned}$$

where  $\gamma_{st}(x_s, x_t) := \theta_s(x_s) + \theta_t(x_t) + \frac{\theta_{st}(x_s, x_t)}{\rho_{st}}$ .

Under appropriate conditions, these pseudo-max-marginals can be used to specify an optimal configuration  $\hat{x}_{\text{TRW}}$ . As detailed in the paper [240], the most general sufficient condition is that there exists a configuration  $\hat{x} = \hat{x}_{\text{TRW}}$  that is nodewise and edgewise optimal across the entire graph, meaning that

$$\hat{x}_s \in \arg \max_{x_s} \nu_s(x_s) \quad \forall s \in V, \text{ and} \quad (8.23a)$$

$$(\hat{x}_s, \hat{x}_t) \in \arg \max_{x_s, x_t} \nu_{st}(x_s, x_t) \quad \forall (s, t) \in E. \quad (8.23b)$$

In this case, we say that the pseudo-max-marginals  $\nu$  satisfy the *strong tree agreement condition*.

Under this condition, Wainwright et al. [240] showed the following:

**Proposition 12.** *For a weight vector  $\rho_e$  in the spanning tree polytope. Then any fixed point  $(M^*, \nu^*)$ , any STA fixed point of TRW-max-product specifies an optimal dual solution for the first-order tree LP relaxation (8.5).*

Fixed points  $\nu^*$  satisfying conditions (8.23) are the most practically relevant, since in this case, the fixed point can be used to determine



the configuration  $\hat{x}_{\text{TRW}}$ , which is guaranteed to be globally optimal for the original problem—that is, an element of  $\arg \max_{x \in \mathcal{X}^m} \langle \theta, \phi(x) \rangle$ —so that the LP relaxation solves the original MAP problem. An interesting theoretical question is whether *any* fixed point  $(M^*, \nu^*)$ , regardless of whether it satisfies the criteria (8.23), specifies an optimal solution to the dual of the first-order LP relaxation (8.5). This question was left open by Wainwright et al. [240], and later resolved by Kolmogorov [134], who provided a counterexample, involving non-binary variables, for which a TRW max-product fixed point does not correspond to a dual-optimal solution. Kolmogorov and Wainwright [135] showed that for pairwise MRFs with binary variables, the equivalence between TRW message-passing and the LP relaxation is exact in all cases: any fixed point of TRW max-product specifies a dual-optimal solution to the first-order LP relaxation (8.5). Examples of pairwise MRFs with binary variables for which TRW message-passing always solves the LP relaxation include the Ising ground state problem, as well as various combinatorial problems such as the independent set problem and the vertex cover problem; see Example 39 for further discussion of these examples.

Kolmogorov [134] also developed a clever sequential scheduling of TRW updates, and established certain convergence guarantees for these so-called TRW-S updates, and showed empirically that the sequential updates tend to outperform a standard parallel scheduling of the TRW max-product updates (8.21). Various forms of these reweighted max-product algorithms have been applied in problems such as segmentation and disparity problems in computer vision [165, 134, 136, 246, 222, 260], error-control coding [73], side-chain prediction [261, 246], sensor fusion [43, 46]. There also turn out to be a number of interesting connections between TRW max-product and a line of research, due to Schlesinger and collaborators, previously published in the Russian literature [205, 137]; Werner [251] provides a detailed overview of this line of work, and some connections to reweighted max-product and LP relaxation.

In addition to the basic TRW algorithm [240] and the TRW-S scheduling studied by Kolmogorov [134], other researchers have proposed distributed algorithms for solving the tree-based LP relax-

ation (8.5), including subgradient methods [73, 136], dual coordinate ascent methods [94, 234], annealing-type methods [121, 246], proximal optimization schemes [193], and adaptive LP solvers [223]. Weiss et al. [246] discuss connections between the zero-temperature limits of convex free energy problems, including the tree-reweighted Bethe problem from Theorem 4 as a special case, and optima of the first-order LP (8.5).

#### 8.4.4 Examples of the first-order LP relaxation

In this section, we discuss various examples of the first-order LP relaxation (8.5). One line of work ongoing in communication and information theory studies the behavior of the LP relaxation for decoding in error-control coding. When applied to combinatorial problems, the first-order LP relaxation recovers various known methods from the integer programming and approximation literature [177, 231]. In the special case of binary variables, there are a number of links to work in the literature on pseudo-Boolean optimization [102, 34].

**Example 38 (LP relaxations for error-control coding).** We begin by discussing an instance of the first-order LP relaxation (8.5), introduced by Feldman et al. [75] for decoding low-density parity check (LDPC) codes. We recall here from Example 8 the notion of an error-correcting code, and its definition as a graphical model. Any binary linear code can be viewed as a Markov random field defined by a hypergraph  $G = (V, F)$  with vertex set  $V = \{1, \dots, m\}$  and hyperedges  $F$ , representing a set of parity checks. Codewords  $x \in \{0, 1\}^m$  are those binary sequences that satisfy a set of parity checks, say of the form  $\bigoplus_{i \in N(a)} x_i = 0$ , where  $N(a) \subset V$  is a subset of bits, and  $\oplus$  denotes addition in modulo two arithmetic. These parity checks are represented by the 0 – 1-valued constraint functions

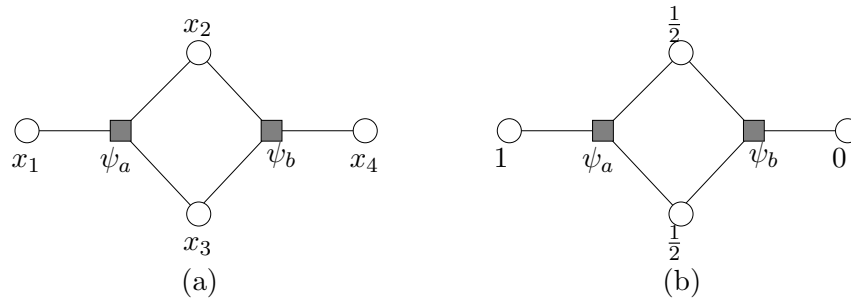
$$\psi_a(x_{N(a)}) = \begin{cases} 1 & \bigoplus_{i \in A} x_i = 0 \\ 0 & \text{otherwise.} \end{cases}$$

Overall, the code can be viewed as an exponential family of the form

$$p_{\theta}(x) = \exp \left\{ \sum_{i \in V} \theta_i x_i \right\} \quad (8.24)$$

with respect to the base measure  $\nu(dx) := [\prod_{a \in F} \psi_a(x_{N(a)})] dx$ . For this model, the set  $\mathcal{M}$  is given by the set of all vectors  $\mu \in \mathbb{R}^m$  that can be expressed as expectations of the form  $\mu_i = \mathbb{E}_f[X_i]$ , where  $f$  is some density with respect to the base measure  $\nu(dx)$ . More concretely, this marginal polytope is simply the convex hull of all possible codewords, which is known as the *codeword polytope*, as discussed in Example 11.

The construction is called a low-density parity check code when the sizes  $|N(a)|$  of the factor neighborhoods stay bounded even as the code length  $m$  increases. Figure 8.3 provides a toy example of an LDPC code over bits  $(x_1, x_2, x_3, x_4) \in \{0, 1\}^4$ , and with two parity checks  $\psi_a$  and  $\psi_b$ , corresponding to the constraints  $x_1 \oplus x_2 \oplus x_3 = 0$  and  $x_2 \oplus x_3 \oplus x_4 = 0$  respectively.



**Fig. 8.3.** (a) The factor graph or hypergraph representation of a toy binary linear code on four bits  $(x_1, x_2, x_3, x_4) \in \{0, 1\}^4$ . Each codeword must satisfy the two parity check constraints  $x_1 \oplus x_2 \oplus x_3 = 0$  and  $x_2 \oplus x_3 \oplus x_4 = 0$ , as defined by the constraint functions  $\psi_a$  and  $\psi_b$ . (b) Construction of a fractional vertex, also known as a pseudocodeword, for the code shown in panel (a). The pseudocodeword has elements  $\tilde{\tau} = [1 \ \frac{1}{2} \ \frac{1}{2} \ 0]$ , which satisfy all the constraints (8.26) defining the LP relaxation. However, the vector  $\tilde{\tau}$  cannot be expressed as a convex combination of codewords, and so lies strictly outside the codeword polytope, which corresponds to the marginal polytope for the graphical model (8.24).

In the form (8.24), the code is not immediately recognizable as a pairwise MRF to which the first-order LP relaxation (8.5) can be ap-

plied. However, any Markov random field over discrete variables can be converted into pairwise form by the selective introduction of auxiliary variables. In particular, suppose that for each check  $a \in F$ , we introduce an auxiliary variable  $z_a$  taking values in the space  $\{0, 1\}^{|N(a)|}$ . Defining pairwise interactions between  $z_a$  and each bit  $x_i$  for nodes  $i \in N(a)$ , we can use  $z_a$  as a device to enforce constraints among the variables  $\{x_i, i \in N(a)\}$ . See Appendix E.2 for further details on converting a general discrete graphical model to an equivalent pairwise form.

If we apply the first-order LP relaxation (8.5) to this pairwise MRF, the relevant variables consist of a pseudomarginal vector  $[1 - \tau_i \ \tau_i]$  for each  $i \in V$ , and for each check  $a \in F$ , a set of pseudomarginals  $\{\tau_{a;J}, J \text{ an even-sized subset of } N(a)\}$ . In this case, the pairwise marginalization conditions that define the set  $\mathbb{L}(G)$  reduce to

$$\sum_{J \ni i} \tau_{a;J} = \tau_i, \quad \text{for each } a, \text{ and } i \in N(a), \quad (8.25)$$

along with the box constraints  $\tau_i, \tau_{a;J} \in [0, 1]$ . It is also possible to remove the variables  $\tau_{a;J}$  by Fourier-Motzkin elimination, so as to obtain an equivalent LP relaxation that is described only in terms of the vector  $[\tau_1 \ \tau_2 \ \dots \ \tau_m]$ . Doing so shows that the first-order relaxation (8.5), when applied to the coding problem, is characterized by the box constraints  $\tau_i \in [0, 1]$  for each  $i \in V$ , and for each  $a \in F$ , the *forbidden set* constraints

$$\sum_{k \in K} (1 - \tau_k) + \sum_{k \in N(a) \setminus K} \tau_k \geq 1 \quad \forall K \subset N(a) \text{ with } |K| \text{ odd}. \quad (8.26)$$

The interpretation of this inequality is very intuitive: it enforces that for each parity check  $a \in F$ , the subvector  $\tau_{N(a)} = \{\tau_i \mid i \in N(a)\}$  must be at Hamming distance at least one from any odd-parity configuration over the bits  $N(a)$ .

The problem of maximum likelihood (ML) decoding is an integer program that is well known to be computationally intractable [16]. Suppose that a codeword  $(x_1, x_2, \dots, x_m)$  is transmitted over a stochastic channel that is memoryless, in that it acts independently on each bit  $x_i$  in the codeword. The channel action is characterized by the conditional distributions  $p(y_i \mid x_i)$ . As a simple example, the binary symmetric

channel simply flips bit  $X_i$  with probability  $\epsilon$ , so that

$$p(y_i | x_i) = \begin{cases} 1 - \epsilon & \text{for } x_i = y_i \\ \epsilon & \text{for } x_i \neq y_i. \end{cases}$$

Given a particular received sequence  $(y_1, \dots, y_m)$  from the channel, define the vector  $\theta \in \mathbb{R}^m$  of log likelihoods, with components  $\theta_i = \log \frac{p(y_i|1)}{p(y_i|0)}$ . The ML decoding problem corresponds to the integer program of maximizing the likelihood  $\sum_{i \in V} \theta_i x_i$  over the discrete set of codewords  $x \in \mathbb{C}$ .

On the other hand, the LP decoding algorithm of Feldman et al. [75] is based on maximizing the objective  $\sum_{i \in V} \theta_i \tau_i$  subject to the box constraints  $(\tau_1, \tau_2, \dots, \tau_m) \in [0, 1]^m$  as well as the forbidden set constraints (8.26). Since its introduction [73, 75], the performance of this LP relaxation has been extensively studied [e.g., 133, 74, 233, 234, 223, 67, 50, 47, 59]. Not surprisingly, given the role of the constraint set  $\mathbb{L}(G)$  in the Bethe variational problem, there are close connections between LP decoding and standard iterative algorithms like sum-product decoding [73, 133, 75]. Among other connections, the fractional vertices of the first-order LP relaxation have a very specific interpretation as *pseudocodewords* of the underlying code, studied in earlier work on iterative decoding [80, 252, 83]. Figure 8.3(b) provides a concrete illustration of a fractional vertex or pseudocodeword that arises when the relaxation is applied to the toy code shown in Figure 8.3(b). Consider the vector  $\tilde{\tau} = [1 \ \frac{1}{2} \ \frac{1}{2} \ 0]$ ; it is easy to verify that it satisfies the box inequalities and the forbidden set constraints (8.26) that define the LP relaxation. (In fact, with a little more work, it can be shown that  $\tilde{\tau}$  is a vertex of the relaxed LP decoding polytope.) However, we claim that  $\tilde{\tau}$  does *not* belong to the marginal polytope for this graphical model—that is, it cannot be written as a convex combination of codewords. To see this fact, note that by taking a mod two sum of the parity check constraints  $x_1 \oplus x_2 \oplus x_3 = 0$  and  $x_2 \oplus x_3 \oplus x_4 = 0$ , we obtain that  $x_1 \oplus x_4 = 0$  for any codeword. Therefore, for any vector  $\mu$  in the codeword polytope, we must have  $\mu_1 = \mu_4$ , which implies that  $\tilde{\tau}$  lies outside the codeword polytope.



Apart from its interest in the context of error-control coding, the fractional vertex in Figure 8.3(b) also provides a counterexample regarding the *persistence* of fractional vertices of the polytope  $\mathbb{L}(G)$ . Recall from our discussion at the end of Section 8.4.1 that the first-order LP relaxation, when applied to integer programs with binary variables, has the strong persistence property [102], as summarized in Proposition 11. It is natural to ask whether strong persistence also holds for higher-order variables as well. Note that after conversion to a pairwise Markov random field (to which the first-order LP relaxation applies), the coding problem illustrated in Figure 8.3(a) includes non-binary variables  $z_a$  and  $z_b$  at each of the factor nodes  $a, b \in F$ . We now claim that the fractional vertex  $\tilde{\tau}$  illustrated in Figure 8.3 constitutes a failure of strong persistence for the first-order LP relaxation with higher-order variables. Consider applying the LP decoder to the cost function  $\theta = (2, 0, 0, -1)$ ; using the previously constructed  $\tilde{\tau} = [1 \ \frac{1}{2} \ \frac{1}{2} \ 0]$ , a little calculation shows that  $\langle \theta, \tilde{\tau} \rangle = 2$ . In contrast, the optimal codewords are  $x^* = (1, 1, 0, 1)$  and  $y^* = (1, 0, 1, 1)$ , both with value  $\langle \theta, x \rangle = \langle \theta, y^* \rangle = 1$ . Thus, neither of the optimal integral solutions have  $x_4 = 0$ , even though  $\tilde{\tau}_4 = 0$  in the fractional vertex. Consequently, strong persistence is violated for this instance.

We conclude that relaxation (8.5) is strongly persistent only for pairwise Markov random fields over binary random variables, otherwise known as binary quadratic programs. We now turn to an in-depth consideration of this binary pairwise case:

**Example 39 (Binary quadratic programs and combinatorial problems).** Recall the Ising model, as first introduced in Example 3: it is an exponential family over a vector  $X$  of binary random variables, which may take either “spin” values  $\{-1, +1\}^m$ , or zero-one values  $\{0, 1\}^m$ . Note that the mapping  $x_s \mapsto 2x_s - 1$  and its inverse  $z_s \mapsto \frac{1}{2}(z_s + 1)$  may be used to convert freely back and forth from the  $\{0, 1\}$  form to the  $\{-1, +1\}$  form.

Let us consider the  $\{0, 1\}$ -case, and the mode-finding problem associated with the canonical overcomplete set of potential functions—

namely:

$$\max_{x \in \{0,1\}^m} \langle \theta, \phi(x) \rangle = \max_{x \in \{0,1\}^m} \left\{ \sum_{s \in V} \theta_s(x_s) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \right\} \quad (8.27)$$

where

$$\theta_s(x_s) := \sum_{j=0}^1 \theta_{s;j} \mathbb{I}_j[x_s] \quad \text{and} \quad \theta_{st}(x_s, x_t) := \sum_{j,k=0}^1 \theta_{st;jk} \mathbb{I}_j[x_s] \mathbb{I}_k[x_t]$$

are weighted sums of indicator functions.

This problem is a binary quadratic program, and includes as special cases various types of classical problems in combinatorial optimization.

**Independent set and vertex cover:** Given an undirected graph  $G = (V, E)$ , an independent set  $I$  is a subset of vertices such that  $(s, t) \notin E$  for all  $s, t \in I$ . Given a set of vertex weights  $w_s \geq 0$ , the maximum weight independent set (MWIS) problem is to find the independent set  $I$  with maximal weight,  $w(I) := \sum_{s \in I} w_s$ . To model this combinatorial problem as an instance of the binary quadratic program, let  $X \in \{0, 1\}^m$  be an indicator vector for membership in  $S$ , meaning that  $X_s = 1$  if and only if  $s \in S$ . Then define  $\theta_s(x_s) = [0 \ w_s]$ , and the pairwise interaction

$$\theta_{st}(x_s, x_t) = \begin{bmatrix} 0 & 0 \\ 0 & -\infty \end{bmatrix}.$$

With these definitions, the binary quadratic program corresponds to the maximum weight independent set (MWIS) problem.

A related problem is that of finding a vertex cover—a set  $C$  of vertices that for any edge  $(s, t) \in E$ , at least one of  $s$  or  $t$  belongs to  $C$ —with minimum weight  $w(C) = \sum_{s \in C} w_s$ . This minimum weight vertex cover (MWVC) problem can be cast as another instance of the binary QP with the same choice of  $\theta_s$ , and the pairwise interactions  $\theta_{st}(0, 0) = +\infty$  and  $\theta_{st}(x_s, x_t) = 0$  for all other binary pairs  $(x_s, x_t) \neq (0, 0)$ .

Let us now consider how the first-order LP relaxation (8.5) specializes to these problems. Recall that the general LP relaxation is in terms of the two-vector of singleton pseudomarginals  $\tau_s = [\tau_{s;0} \ \tau_{s;1}]$ , and an

analogous  $2 \times 2$  matrix  $\tau_{st}$  of pairwise pseudomarginals. In the independent set problem, the parameter setting  $\theta_{st;11} = -\infty$  is tantamount to enforcing the constraint  $\tau_{st;11} = 0$ . After simplification, the constraints defining  $\mathbb{L}(G)$  (see Proposition 4) can be reduced to  $\mu_s \equiv \tau_{s;1} \geq 0$  for all nodes  $s \in V$ , and  $\mu_s + \mu_t \leq 1$  for all edges  $(s, t) \in E$ . Thus, for the MWIS problem, the first-order relaxation (8.5) reduces to the linear program

$$\max_{\mu \geq 0} \sum_{s \in V} w_s \mu_s \quad \text{such that } \mu_s + \mu_t \leq 1 \text{ for all } (s, t) \in E.$$

This LP relaxation is the classical one for the independent set problem [177, 231]. Sanghavi et al. [202] discuss some connections between the ordinary max-product algorithm and this LP relaxation, as well as to auction algorithms [20].

In a similar way, specializing the first-order LP relaxation (8.5) to the MWVC problem yields the linear program

$$\min_{\mu \geq 0} \sum_{s \in V} w_s \mu_s \quad \text{such that } \mu_s + \mu_t \geq 1 \text{ for all } (s, t) \in E,$$

which is another classical LP relaxation from the combinatorics literature [177].

**MAX-CUT:** Consider the following graph-theoretic problem: given a non-negative weight  $w_{st} \geq 0$  for each edge of an undirected graph  $G = (V, E)$ , find a partition  $(U, U^c)$  of the vertex set such that the associated weight

$$w(U, U^c) := \sum_{\{(s,t) \mid s \in U, t \in U^c\}} w_{st}$$

of edges across the partition is maximized. To model this MAX-CUT problem as an instance of the binary quadratic program, let  $X \in \{0, 1\}^m$  be an indicator vector for membership in  $U$ , meaning that  $X_s = 1$  if and only if  $s \in U$ . Then define  $\theta_s(x_s) = 0$  for all vertices, and define the pairwise interaction

$$\theta_{st}(x_s, x_t) = \begin{bmatrix} 0 & w_{st} \\ w_{st} & 0 \end{bmatrix}. \quad (8.28)$$



With these definitions, a little algebra shows that problem (8.27) is equivalent to the MAX-CUT problem, a canonical example of an NP-complete problem. As before, the first-order LP relaxation (8.5) can be specialized to this problem; later, we also describe the celebrated semidefinite program (SDP) relaxation for MAX-CUT due to Goemans and Williamson [95].

**Supermodular and submodular interactions:** An important subclass of binary quadratic programs are those based on supermodular potential functions. The interaction  $\theta_{st}$  is *supermodular* if  $\theta_{st}(1, 1) + \theta_{st}(0, 0) \geq \theta_{st}(1, 0) + \theta_{st}(0, 1)$ ; on the other hand, it is *submodular* if  $\theta_{st}(1, 1) + \theta_{st}(0, 0) \leq \theta_{st}(1, 0) + \theta_{st}(0, 1)$ . Note that the MAX-CUT problem and the independent set problem both involve submodular potential functions, whereas the vertex cover problem involves supermodular potentials. It is well known that the class of *regular binary QPs*—meaning supermodular maximization problems or submodular minimization problems—can be solved in polynomial time. As a particular instance, consider the problem of finding the minimum s-t cut in a graph. This can be formulated as a minimization problem in terms of the potential functions (8.28), with additional non-zero singleton potentials  $\theta_s$ , yielding an instance of submodular minimization. It is easily solved by conversion to a maximum flow problem, using the classical Ford-Fulkerson duality theorem [20, 98]. However, the MAX-CUT, independent set, and vertex cover problems all fall outside this class, and indeed are canonical instances of intractable problems.

Among other results, Kolmogorov and Wainwright [135] establish that the tree-reweighted max-product is exact for any regular binary QP. The same statement fails to hold for the ordinary max-product updates, since it fails on the regular binary QP discussed in Example 37. The exactness of tree-reweighted max-product stems from the tightness of the first-order LP relaxation (8.5) for regular binary QPs. This tightness can be established via results due to Hammer et al. [102] on the so-called roof dual relaxation in pseudo-Boolean optimization, which turns out to be equivalent to the tree-based relaxation (8.5) for the special case of binary variables.



Finally, we discuss a related class of combinatorial problems for which some recent work has studied message-passing and linear programming:

**Example 40 (Maximum weight matching problems).** Given an undirected graph  $G = (V, E)$ , the matching problem is to find a subset  $F$  of edges, such that each vertex is adjacent to at most one edge  $e \in F$ . In the weighted variant, each edge is assigned a weight  $w_e$ , and the goal is to find the matching  $F$  that maximizes the weight function  $\sum_{e \in F} w_e$ . This maximum weight matching (MWM) problem is well known to be solvable in polynomial time for any graph [207]. For bipartite graphs, the MWM problem can be reduced to an especially simple linear program, which we derive here as a special case of the first-order relaxation (8.5).

In order to apply the first-order relaxation, it is convenient to first reformulate the matching problem as a mode-finding problem in an MRF described by a factor graph, and then convert the factor graph to a pairwise form, as in Example 38 above. We begin by associating with the original graph  $G = (V, E)$  a hypergraph  $\tilde{G}$ , in which each edge  $e \in E$  corresponds to a vertex of  $\tilde{G}$ , and each vertex  $s \in V$  corresponds to a hyperedge. The hyperedge indexed by  $s$  connects to all those vertices of  $\tilde{G}$  (edges of the original graph) in the set  $E(s) = \{e \in E \mid s \in e\}$ . Finally, we define a Markov random field over the hypergraph as follows. First, let each  $e \in E$  be associated with a binary variable  $x_e \in \{0, 1\}$ , which acts as an indicator variable for whether edge  $e$  participates in the matching. We define the weight function  $\theta_e(x_e) = w_e x_e$ , where  $w_e$  is the weight specified for edge  $e$  in the matching problem. Second, we define an interaction potential over the variables  $x_{E(s)} = \{x_e \mid e \in E(s)\}$  according to

$$\theta_s(x_{E(s)}) := \begin{cases} 0 & \text{if } \sum_{e \ni s} x_e \leq 1 \\ -\infty & \text{otherwise.} \end{cases}$$

With this definition, the mode-finding problem

$$\max_{x \in \{0,1\}^{|E|}} \left\{ \sum_{e \in E} \theta_e(x_e) + \sum_{s \in V} \theta(x_{E(s)}) \right\}$$

in the hypergraph is equivalent to the original matching problem. Moreover, this hypergraph problem can be converted to an equivalent mode-finding problem in a pairwise MRF by following the generic recipe described in Appendix E.2. Doing so and applying the first-order LP relaxation (8.5) to the resulting pairwise MRF yields the following LP relaxation of the maximum weight matching  $M^*$ :

$$\begin{aligned} M^* &\leq \max_{\tau \in \mathbb{R}^{|E|}} \sum_{e \in E} w_e \tau_e \\ \text{s.t. } &x_e \geq 0 \quad \forall e \in E, \quad \text{and} \quad \sum_{e \ni s} x_e \leq 1 \quad \forall s \in V. \end{aligned} \quad (8.29)$$

This LP relaxation is a classical one for the matching problem, known to be tight for any bipartite graph but loose for non-bipartite graphs [207]. A line of recent research has established close links between the LP relaxation and the ordinary max-product algorithm, including the case of bipartite weighted matching [13], bipartite weighted  $b$ -matching [113], weighted matching on general graphs [201], and weighted  $b$ -matching on general graphs [12]. ♣

## 8.5 Higher-order LP relaxations

The tree-based relaxation (8.17) can be extended to hypertrees of higher treewidth  $t$ , by using the hypertree-based outer bounds  $\mathbb{L}_t(G)$  on marginal polytopes described in Section 4.2.2. This extension produces a sequence of progressively tighter LP relaxations, which we describe here. Given a hypergraph  $G = (V, E)$ , we use the shorthand notation  $x_h := \{x_i \mid i \in h\}$  to denote the set of variables associated with hyperedge  $h \in E$ . We define interaction potentials  $\theta_h(x_h) := \sum_{J \subseteq h} \theta_{h,J} \mathbb{I}[x_h = J]$  and consider the mode-finding problem of computing

$$x^* \in \arg \max_{x \in \mathcal{X}^m} \left\{ \sum_{h \in E} \theta_h(x_h) \right\}. \quad (8.30)$$

Note that problem (8.30) generalizes the analogous problem for pairwise MRFs, a special case in which the hyperedge set consists of only vertices and pairwise edges.

Letting  $t + 1$  denote the maximal cardinality of any hyperedge, the relaxation based on  $\mathbb{L}_t(G)$  involves the collection of pseudomarginals  $\{\tau_h \mid h \in E\}$  subject to local consistency constraints

$$\mathbb{L}_t(G) := \left\{ \tau \geq 0 \mid \sum_{x'_h} \tau_h(x'_h) = 1 \quad \forall h \in E, \text{ and} \right. \\ \left. \sum_{\{x'_h \mid x'_g = x_g\}} \tau_h(x'_h) = \tau_g(x_g) \quad \forall g \subset h \right\}. \quad (8.31)$$

Following the same reasoning as in Section 8.4.1, we have the upper bound

$$\max_{x \in \mathcal{X}^m} \left\{ \sum_{h \in E} \theta_h(x_h) \right\} \leq \underbrace{\max_{\tau \in \mathbb{L}_t(G)} \sum_{h \in E} \left[ \sum_{x_h} \tau_h(x_h) \theta_h(x_h) \right]}_{\max_{\tau \in \mathbb{L}_t(G)} \langle \tau, \theta \rangle}. \quad (8.32)$$

Notice that the relaxation (8.32) can be applied directly to a graphical model that involves higher-order interactions, obviating the need to convert higher-order interactions into pairwise form, as done in illustrating the first-order LP relaxation (see Example 38). In fact, in certain cases, this direct application can yield a tighter relaxation than that based on conversion to the pairwise case, as illustrated in Example 41 below. Of course, the relaxation (8.32) can also be applied directly to a pairwise Markov random field, which can always be embedded into a hypergraph with higher-order interactions ( $t + 1 > 2$ ). Thus, equation (8.32) actually describes a sequence of relaxations, with increasing accuracy as the interaction size  $t + 1$  is increased. The trade-off, of course, is that computational complexity of the relaxation also increases in the parameter  $t$ . The following result is an immediate consequence of our development thus far:

**Proposition 13** (Hypertree tightness). *The LP relaxation (8.32) is tight for any hypergraph  $G$  of treewidth  $t$ .*

*Proof.* This assertion is equivalent to establishing that  $\mathbb{L}_t(G) = \mathbb{M}(G)$  for any hypergraph  $G$  of treewidth  $t$ . Here  $\mathbb{M}(G)$  denotes the marginal polytope, corresponding to all marginal probability distributions  $\{\mu_h \mid h \in E\}$  that are globally consistent. First, the inclusion  $\mathbb{L}_t(G) \supseteq \mathbb{M}(G)$  is immediate, since any set of globally consistent marginals must satisfy the local consistency conditions defining  $\mathbb{L}_t(G)$ .

In the reverse direction, consider a locally consistent set of pseudo-marginals  $\tau \in \mathbb{L}_t(G)$ . Recall the factorization (4.41) of any hypertree factorization in terms of marginals on its hyperedges. Given  $\tau$ , let us define the distribution

$$\begin{aligned} p_\tau(x_1, x_2, \dots, x_m) &= \prod_{h \in E} \varphi_h(x_h; \tau) \\ &= \prod_{h \in E} \left( \prod_{g \in \mathcal{D}^+(h)} [\mu_g(x_g)]^{\omega(g,h)} \right), \end{aligned}$$

where  $\omega$  is the Möbius function associated with the hypertree; see the discussion prior to equation (4.41) and Appendix E.1 for more background.

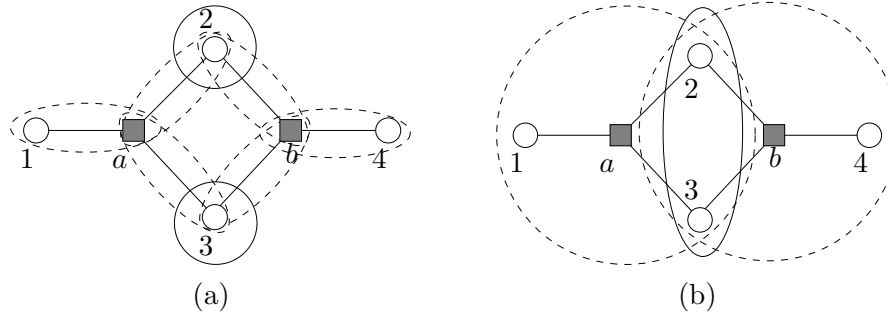
By the nature of this construction and the junction tree theorem, this distribution has marginal distribution  $\tau_h$  for every hyperedge  $h \in E$ . Consequently, it provides a certificate of the membership of  $\tau$  in the marginal polytope  $\mathbb{M}(G)$ .  $\square$

In the binary  $\{0, 1\}$  case, this sequence of relaxations has been proposed and studied previously by Hammer et al. [102], Boros et al. [33], and Sherali and Adams [212], although without the connections to the underlying graphical structure provided by Proposition 13.

**Example 41 (Tighter relaxations for higher-order interactions).** We begin by illustrating how higher-order LP relaxations yield tighter constraints with a continuation of Example 38. Consider the factor graph shown in Figure 8.4(a), corresponding to a hypergraph-structured MRF of the form

$$p_\theta(x) \propto \exp \left\{ \sum_{i=1}^4 \theta_i(x_i) + \theta_a(x_1, x_2, x_3) + \theta_b(x_2, x_3, x_4) \right\}, \quad (8.33)$$

for suitable potential functions  $\{\theta_i\}$ ,  $\theta_a$  and  $\theta_b$ . In this example, we consider two different procedures for obtaining an LP relaxation: (a) first convert this hypergraph MRF (8.33) into a pairwise MRF, and then apply the first-order relaxation (8.5), and (b) apply the second-order relaxation based on  $\mathbb{L}_2(G)$  directly to the original hypergraph MRF. After some algebraic manipulation, both relaxations can be ex-



**Fig. 8.4.** (a) A hypergraph-structured Markov random field with two sets of triplet interactions over  $a = \{1, 2, 3\}$  and  $b = \{2, 3, 4\}$ . The first-order LP relaxation (8.5) applied to this graph first converts it into an equivalent pairwise MRF, and thus enforces only consistency only via the singleton marginals  $\tau_2$  and  $\tau_3$ . (b) The second-order LP relaxation enforces additional consistency over the shared pairwise pseudomarginal  $\tau_{23}$ , and is exact for this graph.

pressed purely in terms of the triplet pseudomarginals  $\tau_{123}(x_1, x_2, x_3)$  and  $\tau_{234}(x_2, x_3, x_4)$ , and in particular their consistency on the overlap  $(x_2, x_3)$ . The basic first-order LP relaxation (procedure (a)) imposes only the two marginalization conditions

$$\sum_{x'_1, x'_2} \tau_{123}(x'_1, x'_2, x_3) = \sum_{x'_2, x'_4} \tau_{234}(x'_2, x_3, x'_4), \quad \text{and} \quad (8.34a)$$

$$\sum_{x'_1, x'_3} \tau_{123}(x'_1, x_2, x'_3) = \sum_{x'_3, x'_4} \tau_{234}(x_2, x'_3, x'_4), \quad (8.34b)$$

which amount to ensuring that the *singleton* pseudomarginals  $\tau_2$  and  $\tau_3$  induced by  $\tau_{123}$  and  $\tau_{234}$  agree. Note, however, that the first-order relaxation imposes no constraint on the pairwise marginal(s)  $\tau_{23}$  induced by the triplet.

In contrast, the second-order relaxation treats the triplets  $(x_1, x_2, x_3)$  and  $(x_2, x_3, x_4)$  exactly, and so addition to the singleton conditions (8.34a), also requires agreement on the overlap—viz.

$$\sum_{x'_1} \tau_{123}(x'_1, x_2, x_3) = \sum_{x'_4} \tau_{234}(x_2, x_3, x'_4). \quad (8.35)$$

As a special case of Proposition 13, this second-order relaxation is tight, since the hypergraph in Figure 8.4(a) has treewidth two.

In Example 38, we considered the first-order LP relaxation applied to the MRF in Figure 8.4(a) for the special case of linear code defined over binary random variables  $x \in \{0, 1\}^4$ . There we constructed the fractional vector

$$\tilde{\tau} = [\tau_1(1) \quad \tau_2(1) \quad \tau_3(1) \quad \tau_4(1)] = [1 \quad \frac{1}{2} \quad \frac{1}{2} \quad 0],$$

and showed that it was a fractional vertex for the relaxed polytope. Although the vector  $\tilde{\tau}$  is permitted under the pairwise relaxation in Figure 8.4(a), we claim that it is forbidden by the second-order relaxation in Figure 8.4(b). Indeed, the singleton pseudomarginals  $\tilde{\tau}$  are consistent with triplet pseudomarginals  $\tau_{123}$  and  $\tau_{234}$  defined as follows: let  $\tau_{123}$  assign mass  $\frac{1}{2}$  to the configurations  $(x_1, x_2, x_3) = (101)$  and  $(110)$ , with zero mass elsewhere, and let  $\tau_{234}$  assign mass  $\frac{1}{2}$  to the configurations  $(x_2, x_3, x_4) = (000)$  and  $(110)$ , and zero elsewhere. By computing the induced marginals, it can be verified that  $\tau_{123}$  and  $\tau_{234}$  marginalize down to  $\tilde{\tau}$ , and satisfy the first order consistency condition (8.34a) required by the first-order relaxation. However, the second-order relaxation also requires agreement over the overlap  $(x_2, x_3)$ ; note that the condition (8.35) is not satisfied, since  $\sum_{x'_1} \tau_{123}(x'_1, x_2, x_3) = \frac{1}{2}\mathbb{I}[(x_2, x_3) = (0, 1)] + \frac{1}{2}\mathbb{I}[(x_2, x_3) = (1, 0)]$ , whereas

$$\sum_{x'_4} \tau_{234}(x_2, x_3, x'_4) = \frac{1}{2}\mathbb{I}[(x_2, x_3) = (1, 1)] + \frac{1}{2}\mathbb{I}[(x_2, x_3) = (1, 1)].$$

More generally, the second-order LP relaxation is exact for Figure 8.4(b), meaning that it is impossible to find any triplet pseudomarginals  $\tau_{123}$  and  $\tau_{234}$  that marginalize down to  $\tilde{\tau}$ , and agree over the overlap  $(x_2, x_3)$ .



Of course, the treewidth-two relaxation can be applied directly to MRFs with cost functions already in pairwise form. In explicit terms, the second-order LP relaxation applied to a pairwise MRF has the form

$$\max_{\tau} \left\{ \sum_{s \in V} \left[ \sum_{x_s} \tau_s(x_s) \theta_s(x_s) \right] + \sum_{(s,t) \in E} \left[ \sum_{x_s, x_t} \tau_{st}(x_s, x_t) \theta_{st}(x_s, x_t) \right] \right\}, \quad (8.36)$$

subject to the constraints

$$\tau_s(x_s) = \sum_{x'_t, x'_u} \tau_{stu}(x_s, x'_t, x'_u) \quad \forall (s, t, u) \ni s, \quad \forall s \in V \quad (8.37a)$$

$$\tau_{st}(x_s, x_t) = \sum_{x'_u} \tau_{stu}(x_s, x_t, x'_u) \quad \forall (s, t, u) \ni (s, t), \quad \forall (s, t) \in E. \quad (8.37b)$$

Assuming that all edges and vertices are involved in the cost function, this relaxation involves  $m$  singleton marginals,  $\binom{m}{2}$  pairwise marginals, and  $\binom{m}{3}$  triplet marginals. Note that even though the triplet pseudomarginals  $\tau_{stu}$  play no role in the cost function (8.36) itself, they nonetheless play a central role in the relaxation via the consistency conditions (8.37) that they impose. Among other implications, equations (8.37a) and (8.37b) imply that the pairwise marginals must be consistent with the singleton marginals (i.e.,  $\sum_{x'_t} \tau_{st}(x_s, x'_t) = \tau_s(x_s)$ ). Since these pairwise conditions define the first-order LP relaxation (8.5), this fact confirms that the second-order LP relaxation is at least as good as the first-order one. In general, the triplet consistency also imposes additional constraints not ensured by pairwise consistency alone. For instance, Example 14 shows that the pairwise constraints are insufficient to characterize the marginal polytope of a single cycle on three nodes, whereas Proposition 13 implies that the triplet constraints (8.37) provide a complete characterization, since a single cycle on three nodes has treewidth two.

This technique—namely, introducing additional parameters in order to tighten a relaxation, such as the triplet pseudomarginals  $\tau_{stu}$ —is known as a *lifting operation*. It is always possible, at least in principle, to project the lifted polytope down to the original space, thereby yielding an explicit representation of the tightened relaxation in the original space. This combination is known as a *lift-and-project*



method [212, 147, 156]. In the case of the lifted relaxation based on the triplet consistency (8.37) defining  $\mathbb{L}_2(G)$ , the projection is from the full space down to the set of pairwise and singleton marginals.

**Example 42 (Lift-and-project for binary MRFs).** Here we illustrate how to project the triplet consistency constraints (8.37) for any binary Markov random field, thereby deriving the so-called cycle inequalities for the binary marginal polytope. (We presented these cycle inequalities earlier in Example 14; see §27.1 of Deza and Laurent [66] for a complete discussion.) In order to do so, it is convenient to work with a minimal set of pseudomarginal parameters: in particular, for a binary Markov random field, the seven numbers  $\{\tau_{stu}, \tau_{st}, \tau_{su}, \tau_{tu}, \tau_s, \tau_t, \tau_u\}$  suffice to characterize the triplet pseudomarginal over variables  $(X_s, X_t, X_u)$ . Following a little algebra, it can be shown that the singleton (8.37a) and pairwise consistency (8.37b) conditions are equivalent to the inequalities:

$$\tau_{stu} \geq 0 \tag{8.38a}$$

$$\tau_{stu} \geq -\tau_s + \tau_{st} + \tau_{su} \tag{8.38b}$$

$$\tau_{stu} \leq 1 - \tau_s - \tau_t - \tau_u + \tau_{st} + \tau_{su} + \tau_{tu} \tag{8.38c}$$

$$\tau_{stu} \leq \tau_{st}, \tau_{su}, \tau_{tu}. \tag{8.38d}$$

Note that following permutations of the triplet  $(s, t, u)$ , there are eight inequalities in total, since inequality (8.38b) has three distinct versions, as does inequality (8.38c).

The goal of projection is to eliminate the variable  $\tau_{stu}$  from the description, and obtain a set of constraints purely in terms of the singleton and pairwise pseudomarginal parameters. If we consider singleton, pairwise, and triplet pseudomarginals for all possible combinations, there are a total of  $T = m + \binom{m}{2} + \binom{m}{3}$  possible pseudomarginal parameters. We would like to project this subset of  $\mathbb{R}^T$  down to a lower-dimensional subset of  $\mathbb{R}^L$ , where  $L = m + \binom{m}{2}$  is the total number of singleton and pairwise parameters. More specifically, we would like to determine the

image of the projection mapping  $\Pi : (\mathbb{L}_2(G)) \rightarrow \mathbb{R}^L$  given by

$$\Pi(\mathbb{L}_2(G)) := \left\{ (\tau_s, \tau_t, \tau_u, \tau_{st}, \tau_{su}, \tau_{tu}) \mid \right. \\ \left. \exists \tau_{stu} \text{ such that inequalities (8.38) hold} \right\}.$$

A classical technique for computing such projections is Fourier-Motzkin elimination [22, 266]. It is based on the following two steps: (a) first express all the inequalities so that the variable  $\tau_{stu}$  to be eliminated appears on the left-hand side; and (b) then combine the ( $\leq$ ) constraints with the ( $\geq$ ) constraints in pairs, thereby yielding a new inequality in which the variable  $\tau_{stu}$  no longer plays a role. Note that  $\tau_{stu}$  appears on the left-hand side of all the inequalities (8.38); hence, performing step (b) yields the following inequalities

$$\begin{aligned} \tau_{st}, \tau_{su}, \tau_{tu} &\geq 0, && \text{combine (8.38a) and (8.38d)} \\ 1 + \tau_{tu} - \tau_t - \tau_u &\geq 0, && \text{combine (8.38b) and (8.38c)} \\ \tau_s - \tau_{su} &\geq 0, && \text{combine (8.38b) and (8.38d)} \\ \tau_s + \tau_{tu} - \tau_{su} - \tau_{st} &\geq 0, && \text{combine (8.38b) and (8.38d)} \\ 1 - \tau_s - \tau_t - \tau_u + \tau_{st} + \tau_{su} + \tau_{tu} &\geq 0, && \text{combine (8.38a) and (8.38c)}. \end{aligned}$$

The first three sets of inequalities should be familiar; in particular, they correspond to the constraints defining the first-order relaxation, in the special case of binary variables (see Example 14). Finally, the last two inequalities, and all permutations thereof, correspond to a known set of inequalities on the binary marginal polytope, usually referred to as the cycle inequalities.<sup>2</sup> In recent work, Sontag and Jaakkola [214] examined the use of these cycle inequalities, as well as their extensions to non-binary variables, in variational methods for approximate marginalization, and showed significantly more accurate results in many cases.



<sup>2</sup>To be clear, in the specified form, these inequalities are usually referred to as the triangle inequalities, for obvious reasons. Note that for a graph  $G = (V, E)$  with  $|V|$  variables, there are a total of  $\binom{|V|}{3}$  such groups of triangle inequalities. However, if the graph  $G$  is not fully connected, then some of these triangle inequalities involve mean parameters  $\mu_{uv}$  for pairs  $(u, v)$  not in the graph edge set. In order to obtain inequalities that involve only mean parameters  $\mu_{st}$  for edges  $(s, t) \in E$ , one can again perform Fourier-Motzkin elimination, which leads to the so-called cycle inequalities. See Deza and Laurent [66] for more details.

# 9

---

## Moment matrices, semidefinite constraints and conic programming

---

Although the linear constraints that we have considered thus far yield a broad class of relaxations, many problems require a more expressive framework for imposing constraints on parameters. In particular, as we have seen at several points in the preceding sections, semidefinite constraints on moment matrices arise naturally within the variational approach—for instance, in our discussion of Gaussian mean parameters from Example 9. This chapter is devoted to a more systematic study of moment matrices, in particular their use in constructing hierarchies of relaxations based on conic programming. Moment matrices and conic programming provide a very expressive language for the design of variational relaxations. In particular, we will see that the LP relaxations considered in earlier sections are special cases of conic relaxations where the underlying cone is the positive orthant. We will also see that moment matrices allow us to define a broad class of additional conic relaxations based on semidefinite programming (SDP) and second-order cone programming (SOCP).

The study of moment matrices and their properties has an extremely rich history [3, 127], particularly in the context of scalar random variables. The basis of our presentation is more recent work [e.g.,

145, 144, 147, 187] that applies to multivariate moment problems. While much of this work aims at general classes of problems in algebraic geometry, in our treatment we limit ourselves to considering marginal polytopes and we adopt the statistical perspective of imposing positive semidefiniteness on covariance and other moment matrices. The moment matrix perspective allows for a unified treatment of various relaxations of marginal polytopes. See Wainwright and Jordan [242] for additional material on the ideas presented here.

### 9.1 Moment matrices and their properties

Given a random vector  $Y \in \mathbb{R}^d$ , consider the collection  $\lambda_{st} = \mathbb{E}[Y_s Y_t]$ ,  $s, t = 1, \dots, d$  of its second-order moments. Using these moments, we can form the following symmetric  $d \times d$  matrix:

$$M[\lambda] = \mathbb{E}[YY^T] = \begin{bmatrix} \lambda_{11} & \lambda_{12} & \cdots & \lambda_{1n} \\ \lambda_{21} & \lambda_{22} & \cdots & \lambda_{2n} \\ \vdots & \vdots & \cdots & \vdots \\ \lambda_{n1} & \lambda_{n2} & \cdots & \lambda_{nn} \end{bmatrix}. \quad (9.1)$$

At first sight, this definition might seem limiting, because the matrix involves only second-order moments. However, given some random vector  $X$  of interest, we can expose any of its moments by defining  $Y = f(X)$  for a suitable choice of function  $f$ , and then considering the associated second-order moment matrix (9.1) for  $Y$ . For instance, by setting  $Y := [1 \ X] \in \mathbb{R} \times \mathbb{R}^m$ , the moment matrix (9.1) will include both first and second-order moments of  $X$ . Similarly, by including terms of the form  $X_s X_t$  in the definition of  $Y$ , we can expose third moments of  $X$ . The significance of the moment matrix (9.1) lies in the following simple result:

**Lemma 1** (Moment matrices). *Any valid moment matrix  $M[\lambda]$  is positive semidefinite.*

*Proof.* We must show that  $a^T M[\lambda] a \geq 0$  for an arbitrary vector  $a \in \mathbb{R}^d$ . If  $\lambda$  is a valid moment vector, then it arises by taking expectations under some distribution  $p$ . Accordingly, we can write

$$a^T M[\lambda] a = \mathbb{E}_p[a^T Y Y^T a] = \mathbb{E}_p[\|a^T Y\|^2],$$

which is clearly non-negative.  $\square$

Lemma 1 provides a *necessary* condition for a vector  $\lambda = \{\lambda_{st} \mid s, t = 1, \dots, d\}$  to be a valid set of second-order moments. Such a condition is both necessary and sufficient for certain classical moment problems involving scalar random variables [e.g., 127, 111]. This condition is also necessary and sufficient for a multivariate Gaussian random vector, as discussed in Example 9.

### 9.1.1 Multinomial and indicator bases

Now consider a discrete random vector  $X \in \{0, 1, \dots, r-1\}^m$ . In order to study moments associated with  $X$ , it is convenient to define some function bases. Our first basis involves multinomial functions over  $(x_1, \dots, x_m)$ . Each multinomial is associated with a *multi-index*, meaning a vector  $\alpha := (\alpha_1, \alpha_2, \dots, \alpha_m)$  of non-negative integers  $\alpha_s$ . For each such multi-index, we define the *multinomial function*

$$x^\alpha := \prod_{s=1}^m x_s^{\alpha_s}. \quad (9.2)$$

Our convention for the all-zeros multi-index  $(0, 0, \dots, 0)$  is that  $x^0 = 1$ . We claim that for discrete random variables in  $\mathcal{X}^m := \{0, 1, \dots, r-1\}^m$ , it suffices to consider multi-indices such that  $\alpha_s \leq r$  for all components  $s$ . Indeed, for any variable  $x \in \mathcal{X} = \{0, 1, \dots, r-1\}$ , note that there holds

$$\prod_{j=0}^{r-1} (x - j) = 0. \quad (9.3)$$

A minor re-arrangement of this relation yields an expression for  $x^r$  as a polynomial of degree  $r-1$ , which implies that any monomial  $x^i$  with  $i \geq r$  can be expressed as a linear combination of lower-order monomials. Therefore, without loss of generality, we can restrict our attention to multi-indices for which the maximum degree  $\|\alpha\|_\infty := \max_s \alpha_s$  is less than or equal to  $r-1$ . Consequently, our multinomial basis involves a total of  $r^m$  multinomial functions. We define the Hamming norm  $\|\alpha\|_0 := \text{card}\{i = 1, \dots, m \mid \alpha_i \neq 0\}$ , which counts the number of

non-zero elements in the multi-index  $\alpha$ . For each integer  $k = 1, \dots, m$ , we then define the multi-index set

$$\mathcal{I}_k := \{ \alpha \mid \|\alpha\|_0 \leq k \}. \quad (9.4)$$

This nested set of multi-index sets describes a hierarchy of models, which can be associated with hypergraphs with increasing sizes of hyperedges. To calculate the cardinality of  $\mathcal{I}_k$ , observe that for each  $i = 0, \dots, k$ , there are  $\binom{m}{i}$  possible subsets of size  $i$ . Moreover, for each member of each such subset, there are  $(r - 1)$  possible choices of the index value, so that  $\mathcal{I}_k$  has  $\sum_{i=0}^k \binom{m}{i} (r - 1)^i$  elements in total. The total number of all possible multi-indices (with  $\|\alpha\|_\infty \leq r - 1$ ) is given by  $|\mathcal{I}_m| = \sum_{i=0}^m \binom{m}{i} (r - 1)^i = r^m$ .

Our second basis is a generalization of the standard overcomplete potentials (3.34), based on indicator functions for events of the form  $\{X_s = j\}$  as well as their higher-order analogs. Recall the  $\{0, 1\}$ -valued indicator function  $\mathbb{I}_{j_s}(x_s)$  for the event  $\{X_s = j\}$ . Using these node-based functions, we define, for each global configuration  $J = (j_1, \dots, j_m) \in \{0, 1, \dots, r - 1\}^m$ , the indicator function

$$\mathbb{I}_J(x) := \prod_{s=1}^m \mathbb{I}_{j_s}(x_s) = \begin{cases} 1 & \text{if } (x_1, \dots, x_m) = (j_1, \dots, j_m) \\ 0 & \text{otherwise.} \end{cases} \quad (9.5)$$

In our discussion thus far, we have defined these function bases over the full vector  $(x_1, \dots, x_m)$ ; more generally, we can also consider the same bases for sub-vectors  $(x_1, \dots, x_k)$ . Overall, for any integer  $k \leq m$ , we have two function bases: the  $r^k$  vector of multinomial functions

$$\mathfrak{M}_k(x_1, \dots, x_k) := \left\{ \prod_{s=1}^k x_s^{\alpha_s} \mid \alpha \in \{0, 1, \dots, r - 1\}^k \right\},$$

and the  $r^k$  vector of indicator functions

$$\mathfrak{I}_k(x_1, \dots, x_k) := \{ \mathbb{I}_J(x) \mid J \in \{0, 1, \dots, r - 1\}^k \},$$

The following elementary lemma shows that these two bases are in one-to-one correspondence:

**Lemma 2.** For each  $k = 2, 3, \dots, m$ , there is an invertible  $r^k \times r^k$  matrix  $B$  such that

$$\mathfrak{M}_k(x) = B \mathfrak{J}_k(x) \quad \text{for all } x \in \mathcal{X}^k. \quad (9.6)$$

*Proof.* Our proof is based on explicitly constructing the matrix  $B$ . For each  $s = 1, \dots, k$  and  $j \in \{0, 1, 2, \dots, r-1\}$ , consider the following identities between the scalar indicator functions  $\mathbb{I}_j(x)$  and monomials  $x^j$ :

$$\mathbb{I}_j(x) = \prod_{\ell \neq j} \frac{x_s - \ell}{j - \ell}, \quad (x)^j = \sum_{\ell=0}^{r-1} (\ell)^j \mathbb{I}_\ell(x). \quad (9.7)$$

Using this identity multiple times (once for each  $s \in \{1, \dots, k\}$ ), we obtain that for each  $\alpha \in \mathcal{I}_k$ ,

$$x^\alpha = \prod_{s=1}^k (x_s)^{\alpha_s} = \prod_{s=1}^k \left[ \sum_{\ell=0}^{r-1} (\ell)^{\alpha_s} \mathbb{I}_\ell(x_s) \right].$$

This expression shows that the multinomial  $x^\alpha$  can be written as a linear combination of the indicator functions  $\{\mathbb{I}_J(x), J \in \mathcal{X}^k\}$ . Conversely, for each  $J \in \{0, 1, \dots, r-1\}^k$ , we have

$$\mathbb{I}_J(x) := \prod_{s=1}^k \mathbb{I}_{j_s}[x_s] = \prod_{s=1}^k \left[ \prod_{\ell \neq j_s} \frac{x_s - \ell}{j_s - \ell} \right],$$

which shows that the indicators can be written as linear combination of the monomials  $\{x^\alpha, \alpha \in \mathcal{I}_k\}$ . Thus, there is an invertible linear transformation between the indicator functions  $\{\mathbb{I}_J(x), J \in \mathcal{X}^k\}$  and the monomials  $\{x^\alpha, \alpha \in \mathcal{I}_k\}$ , and we let  $B \in \mathbb{R}^{r^k \times r^k}$  denote the invertible matrix that carries out this transformation.

Both bases are convenient for different purposes, and Lemma 2 allows us to move freely between them. The mean parameters associated with the indicator basis  $\mathfrak{J}$  are easily interpretable as probabilities—viz.  $\mathbb{E}[\mathbb{I}_J(X)] = \mathbb{P}[X = J]$ . However, the multinomial basis is convenient for defining hierarchies of semidefinite relaxations.

### 9.1.2 Marginal polytopes for hypergraphs

We now turn to the use of these function bases in order to define the notion of a marginal polytope for a general hypergraph. Given a hypergraph  $G = (V, E)$  where the maximal hyperedge has cardinality  $k$ , we may consider the multinomial Markov random field  $p_\theta(x) \propto \exp \left\{ \sum_{\alpha \in \mathcal{I}(G)} \theta_\alpha x^\alpha \right\}$ . Here  $\mathcal{I}(G) \subseteq \mathcal{I}_k$  corresponds to those multi-indices associated with the hyperedges of  $G$ . (For instance, if  $G$  includes the triplet hyperedge  $\{s, t, u\}$ , then  $\mathcal{I}(G)$  must include the set of  $r^3$  multi-indices  $\alpha$ , with  $(\alpha_s, \alpha_t, \alpha_u)$  ranging over  $\{0, 1, \dots, r-1\}^3$ , and  $\alpha_v = 0$  for all  $v \notin \{s, t, u\}$ .) As an important special case, for each integer  $t = 1, 2, \dots, m$ , we also define the multi-index sets  $\mathcal{I}_t := \mathcal{I}(K_{m,t})$ , where  $K_{m,t}$  denotes the hypergraph on  $m$  nodes that includes all hypergraphs up to size  $t$ . For instance, the hypergraph  $K_{m,2}$  is simply the ordinary complete graph on  $m$  nodes, including all  $\binom{m}{2}$  edges.

For any multi-index  $\alpha \in \mathbb{Z}_+^m$ , let

$$\mu_\alpha := \mathbb{E}[X^\alpha] = \mathbb{E} \left[ \prod_{i=1}^m X_i^{\alpha_i} \right] \tag{9.8}$$

denote the associated mean parameter or moment. We use  $\mathbb{M}(G)$  to denote the marginal polytope associated with hypergraph  $G$ —that is,

$$\mathbb{M}(G) = \{ \mu_\alpha, \alpha \in \mathcal{I}(G) \mid \mu_\alpha = \mathbb{E}_m[X^\alpha] \text{ for some } p \}. \tag{9.9}$$

To be precise, it should be noted our choice of notation is slightly inconsistent with previous sections, where we defined marginal polytopes in terms of the indicator functions  $\mathbb{I}_j(x_s)$  and  $\mathbb{I}_j[x_s]\mathbb{I}_k[x_t]$ . However, since there is a one-to-one correspondence between between these indicator functions and the multinomials  $\{x^\alpha\}$ , the corresponding marginal polytopes are isomorphic objects, regardless of the underlying potential functions chosen.

## 9.2 Semidefinite bounds on marginal polytopes

We now describe the Lasserre sequence [145, 147] of semidefinite outer bounds on the marginal polytope  $\mathbb{M}(G)$ . This sequence is defined in



terms of a hierarchy of moment matrices  $M_t[\mu]$ , indexed by an integer parameter  $t = 1, 2, \dots, m$ . Each moment matrix  $M_t$  is defined in terms of the subset  $\mathcal{I}_t$  of multi-indices associated with the hypergraph  $K_{m,t}$ , and so contains as a minor each matrix  $M_s$  for all integers  $s < t$ . Overall, these moment matrices generate a nested sequence of semidefinite outer bounds on any marginal polytope.

### 9.2.1 Lasserre sequence

To define the relevant moment matrices, for each  $t = 1, 2, \dots, m$ , consider the  $|\mathcal{I}_t| \times |\mathcal{I}_t|$  matrix of moments  $M_t[\mu]$  defined by the  $|\mathcal{I}_t| := r^t$ -dimensional random vector  $Y := \{X^\alpha \mid \alpha \in \mathcal{I}_t\}$ . Each row and column of  $M_m[\mu]$  is associated with some multi-index  $\alpha \in \mathcal{I}_t$ , and its entries are specified as follows:

$$(M_t[\mu])_{\alpha\beta} := \mu_{\alpha+\beta} = \mathbb{E}[X^\alpha X^\beta]. \quad (9.10)$$

As a particular example, Figure 9.1(a) provides an illustration of the matrix  $M_3[\mu]$  for the special case of binary variables ( $r = 2$ ) on three nodes ( $m = 3$ ), so that the overall matrix is eight-dimensional. The shaded region in Figure 9.1(b) shows the matrix  $M_2[\mu]$  for this same example; note that it has  $|\mathcal{I}_2| = 7$  rows and columns.

Some clarifying comments regarding these moment matrices: in both panels, so as to simplify notation in this binary case, we have used  $\mu_1$  as a short-hand for the first-order moment  $\mu_{1,0,0}$ , with similar short-hand for the other first-order moments  $\mu_2$  and  $\mu_3$ . Similarly, the quantity  $\mu_{12}$  is short-hand for the second-order moment  $\mu_{1,1,0} = \mathbb{E}[X_1^1 X_2^1 X_3^0]$ , and the quantity  $\mu_{123}$  denotes the triplet moment  $\mu_{1,1,1} = \mathbb{E}[X_1 X_2 X_3]$ . In both matrices, the element in the upper left-hand corner is  $\mu_\emptyset = \mu_{0,0,0} = \mathbb{E}[X^0]$ , which is always equal to one. Finally, in calculating the form of these moment matrices, we have repeatedly used the fact that  $X_i^2 = X_i$  for any binary variable to simplify moment calculations. For instance, in computing the (8, 7) element of  $M_3[\mu]$ , we write  $\mathbb{E}[(X_1 X_2 X_3)(X_1 X_3)] = \mathbb{E}[X_1 X_2 X_3] = \mu_{123}$ .

We now describe how the moment matrices  $M_t[\mu]$  induce outer bounds on the marginal polytope  $\mathbb{M}(G)$ . Given a hypergraph  $G$ , define the mapping  $\Pi_G : \mathbb{R}^{|\mathcal{I}_t|} \rightarrow \mathbb{R}^{|\mathcal{I}(G)|}$  that maps any vector  $\mu \in \mathbb{R}^{\mathcal{I}_m}$  to the

$$\begin{bmatrix} 1 & \mu_1 & \mu_2 & \mu_3 & \mu_{12} & \mu_{23} & \mu_{13} & \mu_{123} \\ \mu_1 & \mu_1 & \mu_{12} & \mu_{13} & \mu_{12} & \mu_{123} & \mu_{13} & \mu_{123} \\ \mu_2 & \mu_{12} & \mu_2 & \mu_{23} & \mu_{12} & \mu_{23} & \mu_{123} & \mu_{123} \\ \mu_3 & \mu_{13} & \mu_{23} & \mu_3 & \mu_{123} & \mu_{23} & \mu_{13} & \mu_{123} \\ \mu_{12} & \mu_{12} & \mu_{12} & \mu_{123} & \mu_{12} & \mu_{123} & \mu_{123} & \mu_{123} \\ \mu_{23} & \mu_{123} & \mu_{23} & \mu_{23} & \mu_{123} & \mu_{23} & \mu_{123} & \mu_{123} \\ \mu_{13} & \mu_{13} & \mu_{123} & \mu_{13} & \mu_{123} & \mu_{123} & \mu_{13} & \mu_{123} \\ \mu_{123} & \mu_{123} & \mu_{123} & \mu_{123} & \mu_{123} & \mu_{123} & \mu_{123} & \mu_{123} \end{bmatrix}$$

(a)

$$\begin{bmatrix} 1 & \mu_1 & \mu_2 & \mu_3 & \mu_{12} & \mu_{23} & \mu_{13} & \mu_{123} \\ \mu_1 & \mu_1 & \mu_{12} & \mu_{13} & \mu_{12} & \mu_{123} & \mu_{13} & \mu_{123} \\ \mu_2 & \mu_{12} & \mu_2 & \mu_{23} & \mu_{12} & \mu_{23} & \mu_{123} & \mu_{123} \\ \mu_3 & \mu_{13} & \mu_{23} & \mu_3 & \mu_{123} & \mu_{23} & \mu_{13} & \mu_{123} \\ \mu_{12} & \mu_{12} & \mu_{12} & \mu_{123} & \mu_{12} & \mu_{123} & \mu_{123} & \mu_{123} \\ \mu_{23} & \mu_{123} & \mu_{23} & \mu_{23} & \mu_{123} & \mu_{23} & \mu_{123} & \mu_{123} \\ \mu_{13} & \mu_{13} & \mu_{123} & \mu_{13} & \mu_{123} & \mu_{123} & \mu_{13} & \mu_{123} \\ \mu_{123} & \mu_{123} & \mu_{123} & \mu_{123} & \mu_{123} & \mu_{123} & \mu_{123} & \mu_{123} \end{bmatrix}$$

(b)

**Fig. 9.1.** Moment matrices and minors defining the Lasserre sequence of semidefinite relaxations. (a) Full matrix  $M_3[y]$ . (b) Shaded region:  $7 \times 7$  principal minor  $M_2[y]$  constrained by the Lasserre relaxation at order 1.

indices  $\{\alpha \in \mathcal{I}(G)\}$ . Then for each  $t = 1, 2, \dots$ , define the semidefinite constraint set

$$\mathbb{S}_t(G) := \Pi_G \left[ \{ \mu \in \mathbb{R}^{|\mathcal{I}_t|} \mid M_t[\mu] \succeq 0 \} \right] \tag{9.11}$$

As a consequence of Lemma 1, each set  $\mathbb{S}_t(G)$  is an outer bound on the marginal polytope  $\mathbb{M}(G)$ , and by definition of the matrices  $M_t[\mu]$ , these outer bounds form a nested sequence (i.e.,  $\mathbb{S}_1(G) \supseteq \mathbb{S}_2(G) \supseteq \mathbb{S}_3(G) \supseteq \dots \mathbb{M}(G)$ ). This sequence is known as the Lasserre sequence of relaxations [144, 147]; Lovász and Schrijver [156] describe a related class of lifting procedure.

**Example 43 (First-order semidefinite bound on  $\mathbb{M}(K_3)$ ).** To illustrate the power of semidefinite bounds, recall Example 16, in which we considered the fully connected graph  $K_3$  on three

nodes. In terms of the six-dimensional vector of sufficient statistics  $(X_1, X_2, X_3, X_1X_2, X_1X_3, X_2X_3)$ , we considered the pseudomarginal vector

$$(\tau_1, \tau_2, \tau_3, \tau_{12}, \tau_{23}, \tau_{13}) = (0.5, 0.5, 0.5, 0.4, 0.4, 0.1).$$

In Example 16, we used the cycle inequalities, as derived in Example 42, to show that  $\tau$  does *not* belong to  $\mathbb{M}(K_3)$ . Here we provide an alternative, direct proof of this fact based on semidefinite constraints. In particular, the moment matrix  $M_1[\tau]$  associated with these putative mean parameters has the form

$$M_1[\tau] = \begin{bmatrix} 1 & 0.5 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0.4 & 0.1 \\ 0.5 & 0.4 & 0.5 & 0.4 \\ 0.5 & 0.1 & 0.4 & 0.5 \end{bmatrix}.$$

A simple calculation shows that this matrix  $M_1[\tau]$  is not positive definite, whence  $\tau \notin \mathbb{S}_1(K_3)$ , so that Lemma 1 implies that  $\tau \notin \mathbb{M}(K_3)$ .

♣

### 9.2.2 Tightness of semidefinite outer bounds

Given the nested sequence  $\mathbb{S}_t(G)$  of outer bounds on the marginal polytope  $\mathbb{M}(G)$ , we now turn to a natural question: what is the minimal order (if any) for which  $\mathbb{S}_t(G)$  provides an exact characterization of the marginal polytope? It turns out that for an arbitrary hypergraph,  $t = m$  is the minimal required (although see Section 9.3 and Proposition 16 for sharper guarantees based on treewidth). This property of finite termination for semidefinite relaxations in a general setting was proved by Lasserre [144], and also by Laurent [147, 146] using different methods. Here we provide an alternative and arguably more direct proof of exact semidefinite characterizations of marginal polytopes:

**Proposition 14** (Tightness of semidefinite constraints). *For a random vector  $X \in \{0, 1, \dots, r-1\}^m$  Markov with respect to any hypergraph  $G$ , the semidefinite constraint set  $\mathbb{S}_m(G)$  provides an exact description of the associated marginal polytope  $\mathbb{M}(G)$ .*

*Proof.* Lemma 2 shows that the multinomial basis  $\mathfrak{M}_m$  and the indicator basis  $\mathfrak{I}_m$  are related by a bijective linear mapping. Accordingly, let us consider the  $r^m \times r^m$  moment matrix associated with the indicator basis  $\mathfrak{I}_m(x) = \{\mathbb{I}_J(X), J \in \mathcal{X}^m\}$ . Its form is very simple: since the product  $\mathbb{I}_J(x)\mathbb{I}_{J'}(x)$  vanishes for all  $J \neq J'$ , it is a diagonal matrix  $D = \text{diag}(\mu_J)$ , where  $\mu_J = \mathbb{P}[X = J]$  is the probability of the configuration  $J \in \mathcal{X}^m$ . Given the constraint  $\sum_J \mathbb{I}_J(x) = 1$ , the positive semidefinite constraint  $D \succeq 0$  is necessary and sufficient to ensure that  $\{\mu_J, J \in \mathcal{X}^m\}$  specifies a valid probability distribution. Moreover, by the linear bijection established above and the linearity of expectation, we have  $M_m[\mu] = BDB^T$  with  $B$  invertible, so that  $D \succeq 0$  if and only if  $M_m[\mu] \succeq 0$ .  $\square$

This result shows that imposing a semidefinite constraint on the largest possible moment matrix  $M_m[\mu]$  is sufficient to fully characterize all marginal polytopes. From a practical point of view, however, the consequences of this result are limited, because  $M_m[\mu]$  is a  $|\mathcal{I}_m| \times |\mathcal{I}_m|$  matrix, where  $|\mathcal{I}_m| = |\mathcal{X}^m| = r^m$  is exponentially large. Moreover, the  $t = m$  condition in Proposition 14 is not only sufficient for exactness, but also necessary in a worst case setting (meaning that for any  $t < m$ , there exists some hypergraph  $G$  such that  $\mathbb{M}(G)$  is strictly contained within  $\mathbb{S}_t(G)$ ). The following example illustrates both the sufficiency and necessity of Proposition 14:

**Example 44 ((In)exactness of semidefinite constraints).** Consider a pair of binary random variables  $(X_1, X_2) \in \{0, 1\}^2$ . With respect to the monomials  $(X_1, X_2, X_1X_2)$ , the marginal polytope  $\mathbb{M}(K_2)$  consists of three moments  $\{\mu_1, \mu_2, \mu_{12}\}$ . Figure 9.2(a) provides an illustration of this three-dimensional polytope; as derived previously in Example 10, this set is characterized by the four constraints

$$\mu_{12} \geq 0, \quad 1 + \mu_1 + \mu_2 - \mu_{12} \geq 0, \quad \text{and} \quad \mu_s - \mu_{12} \geq 0 \quad \text{for } s = 1, 2. \quad (9.12)$$

The first-order semidefinite constraint set  $\mathbb{S}_1(K_2)$  is defined by the

semidefinite moment matrix constraint

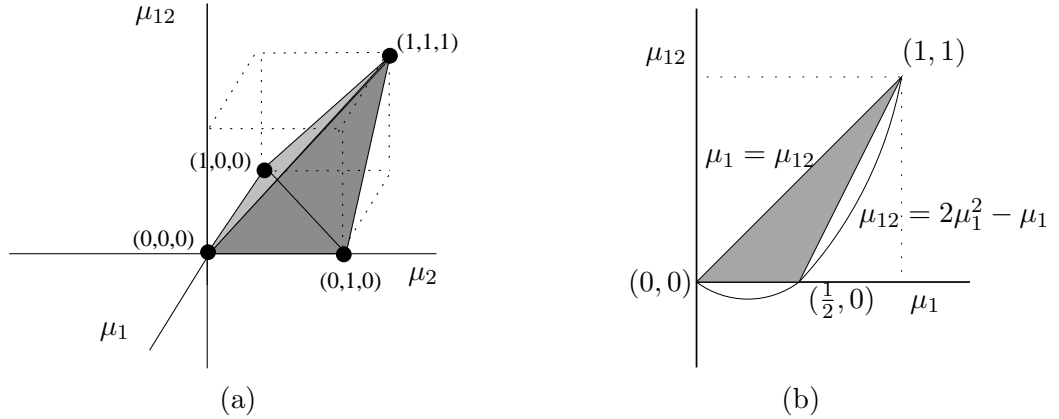
$$M_1[\mu] = \begin{bmatrix} 1 & \mu_1 & \mu_2 \\ \mu_1 & \mu_1 & \mu_{12} \\ \mu_2 & \mu_{12} & \mu_2 \end{bmatrix} \succeq 0. \quad (9.13)$$

In order to deduce the various constraints implied by this semidefinite inequality, recall a slight extension of Sylvester's criterion for assessing positive semidefiniteness: a square matrix is positive semidefinite if and only if the determinant of any principal submatrix is non-negative (see Horn and Johnson [111], p. 405). So as to facilitate both our calculation and subsequent visualization, it is convenient to focus on the intersection of both the marginal polytope and the constraint set  $\mathbb{S}_1(K_2)$  with the hyperplane  $\mu_1 = \mu_2$ . Stepping through the requirements of Sylvester's criterion, positivity of the (1, 1) element is trivial ( $1 > 0$ ), and the remaining singleton principal submatrices imply that  $\mu_1, \mu_2 \geq 0$ . The (1, 2) and (1, 3) principal submatrices are equivalent to the interval constraints  $\mu_1, \mu_2 \in [0, 1]$ . The (2, 3) principal submatrix yields  $\mu_1\mu_2 \geq \mu_{12}^2$ , which after setting  $\mu_1 = \mu_2$  reduces to  $\mu_1 \geq |\mu_{12}|$ . Finally, after some simplification (and setting  $\mu_1 = \mu_2$ ), non-negativity of the full determinant leads to the constraint  $\mu_{12}^2 - (2\mu_1^2)\mu_{12} + (2\mu_1^3 - \mu_2^2) \leq 0$ . Viewing the left-hand side as a quadratic in  $\mu_{12}$ , we can factor it into the product  $[\mu_{12} - \mu_1][\mu_{12} - \mu_1(2\mu_1 - 1)]$ . For  $\mu_1 \in [0, 1]$ , this quadratic inequality is equivalent to the pair of constraints

$$\mu_{12} \leq \mu_1, \quad \mu_{12} \geq \mu_1(2\mu_1 - 1). \quad (9.14)$$

The gray area in Figure 9.2(b) shows the intersection of the three-dimensional marginal polytope  $\mathbb{M}(K_2)$  from panel (a) with the hyperplane  $\mu_1 = \mu_2$ . The intersection of the semidefinite constraint set  $\mathbb{S}_1(K_2)$  with this same hyperplane is characterized by the interval inclusion  $\mu_1 \in [0, 1]$  and the two inequalities in equation (9.14). Note that the semidefinite constraint set is an outer bound on  $\mathbb{M}(K_2)$ , but that it includes points that are clearly not valid marginals. For instance, it can be verified that  $(\mu_1, \mu_2, \mu_{12}) = (\frac{1}{4}, \frac{1}{4}, -\frac{1}{8})$  corresponds to a positive semidefinite  $M_1[\mu]$ , but this vector certainly does not belong to  $\mathbb{M}(K_2)$ .

In this case, if we move up one more step to the semidefinite outer bound  $\mathbb{S}_2(K_2)$ , then Proposition 14 guarantees that the description



**Fig. 9.2.** (a) (b) Nature of the semidefinite outer bound  $\mathbb{S}_1$  on the marginal polytope  $\mathbb{M}(K_2)$  for a pair  $(x_1, x_2) \in \{0, 1\}^2$ . The gray area shows the cross-section of the binary marginal polytope  $\mathbb{M}(K_2)$  corresponding to intersection with the hyperplane  $\mu_1 = \mu_2$ . The intersection of  $\mathbb{S}_1$  with this same hyperplane is defined by the inclusion  $\mu_1 \in [0, 1]$ , the linear constraint  $\mu_{12} \leq \mu_1$ , and the quadratic constraint  $\mu_{12} \geq 2\mu_1^2 - \mu_1$ . Consequently, there are points belonging to  $\mathbb{S}_1$  that lie strictly outside  $\mathbb{M}(K_2)$ .

should be exact. To verify this fact, note that  $\mathbb{S}_2(K_2)$  is based on imposing positive semidefiniteness of the moment matrix  $M_2[\mu]$ , as represented by the gray shaded region in Figure 9.1(b). Positivity of the diagonal element  $(4, 4)$  gives the constraint  $\mu_{12} \geq 0$ . Positivity of the  $(3, 4)$  subminor, combined with the constraint  $\mu_{12} \geq 0$ , leads to  $\mu_2 - \mu_{12} \geq 0$ . By symmetry, the  $(2, 4)$  subminor gives  $\mu_1 - \mu_{12} \geq 0$ . Finally, calculating the determinant of  $M_2[\mu]$  yields

$$\det M_2[\mu] = \mu_{12} [\mu_1 - \mu_{12}] [\mu_1 - \mu_{12}] [1 + \mu_{12} - \mu_1 - \mu_2]. \quad (9.15)$$

The constraint  $\det M_2[\mu] \geq 0$ , in conjunction with the previous constraints, implies the inequality  $1 + \mu_{12} - \mu_1 - \mu_2 \geq 0$ . In fact, the quantities  $\{\mu_{12}, \mu_1 - \mu_{12}, \mu_2 - \mu_{12}, 1 + \mu_{12} - \mu_1 - \mu_2\}$  are the eigenvalues of  $M_2[\mu]$ , so positive semidefiniteness of  $M_2[\mu]$  is equivalent to non-negativity of these four quantities. The positive semidefiniteness of  $\mathbb{S}_2(K_2)$  thus recovers the four inequalities (9.12) that define  $\mathbb{M}(K_2)$ , thereby providing an explicit confirmation of Proposition 14. ♣

### 9.3 Link to LP relaxations and graphical structure

Recall from our previous discussion that the complexity of a given marginal polytope  $\mathbb{M}(G)$  depends very strongly on the structure of the (hyper)graph  $G$ . One consequence of the junction tree theorem is that marginal polytopes associated with hypertrees are straightforward to characterize (see Propositions 1 and 13). This simplicity is also apparent in the context of semidefinite characterizations. In fact, the LP relaxations discussed in Section 8.5, which are tight for hypertrees of appropriate widths, can be obtained by imposing semidefinite constraints on particular minors of the overall moment matrix  $M_m[\mu]$ , as we illustrate here.

Given a hypergraph  $G = (V, E)$  and the associated subset of multi-indices  $\mathcal{I}(G)$ , let  $\mathbb{M}(G)$  be the  $|\mathcal{I}(G)|$ -dimensional marginal polytope (over the sufficient statistics  $\{x^\alpha \mid \alpha \in \mathcal{I}(G)\}$ ). For each hyperedge  $h \in E$ , let  $\mathcal{I}(h)$  be the subset of multi-indices with non-zero components only in the elements of  $h$ , and let  $\mathbb{M}(h)$  be the  $|\mathcal{I}(h)|$ -dimensional marginal polytope associated with the multinomials  $\{x^\alpha \mid \alpha \in \mathcal{I}(h)\}$ . Letting the maximal hyperedge have cardinality  $t + 1$ , we can then rewrite the hypertree-based LP relaxation from Section 8.5 in the form

$$\mathbb{L}_t(G) = \bigcap_{h \in E} \left\{ \mu \in \mathbb{R}^{|\mathcal{I}(G)|} \mid \Pi_h(\mu) \in \mathbb{M}(\{h\}) \right\}, \quad (9.16)$$

where  $\Pi_h : \mathbb{R}^{|\mathcal{I}(G)|} \rightarrow \mathbb{R}^{|\mathcal{I}(h)|}$  projects down to multi-indices in the set  $\mathcal{I}(h)$ .

In general, imposing semidefiniteness on moment matrices produces constraint sets that are convex but non-polyhedral (i.e., with curved boundaries; see Figure 9.2(b) for an illustration). In certain cases, however, a semidefinite constraint actually reduces to a set of linear inequalities, so that the associated constraint set is a polytope. We have already seen one instance of this phenomenon in Proposition 14, where by a basis transformation (between the indicator functions  $\mathbb{1}_j(x_s)$  and the monomials  $x_s^{\alpha_s}$ ), we showed that a semidefinite constraint on the full matrix  $M_m[\mu]$  is equivalent to a large set of linear inequalities.

In similar fashion, it turns out that the LP relaxation  $\mathbb{L}_t(G)$  can be described in terms of semidefiniteness constraints on a subset of minors from the full moment matrix  $M_m[\mu]$ . In particular, for each hyperedge

1	$\mu_1$	$\mu_2$	$\mu_3$	$\mu_{12}$	$\mu_{23}$	$\mu_{13}$	$\mu_{123}$
$\mu_1$	$\mu_1$	$\mu_{12}$	$\mu_{13}$	$\mu_{12}$	$\mu_{123}$	$\mu_{13}$	$\mu_{123}$
$\mu_2$	$\mu_{12}$	$\mu_2$	$\mu_{23}$	$\mu_{12}$	$\mu_{23}$	$\mu_{123}$	$\mu_{123}$
$\mu_3$	$\mu_{13}$	$\mu_{23}$	$\mu_3$	$\mu_{123}$	$\mu_{23}$	$\mu_{13}$	$\mu_{123}$
$\mu_{12}$	$\mu_{12}$	$\mu_{12}$	$\mu_{123}$	$\mu_{12}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$
$\mu_{23}$	$\mu_{123}$	$\mu_{23}$	$\mu_{23}$	$\mu_{123}$	$\mu_{23}$	$\mu_{123}$	$\mu_{123}$
$\mu_{13}$	$\mu_{13}$	$\mu_{123}$	$\mu_{13}$	$\mu_{123}$	$\mu_{123}$	$\mu_{13}$	$\mu_{123}$
$\mu_{123}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$

(a)

1	$\mu_1$	$\mu_2$	$\mu_3$	$\mu_{12}$	$\mu_{23}$	$\mu_{13}$	$\mu_{123}$
$\mu_1$	$\mu_1$	$\mu_{12}$	$\mu_{13}$	$\mu_{12}$	$\mu_{123}$	$\mu_{13}$	$\mu_{123}$
$\mu_2$	$\mu_{12}$	$\mu_2$	$\mu_{23}$	$\mu_{12}$	$\mu_{23}$	$\mu_{123}$	$\mu_{123}$
$\mu_3$	$\mu_{13}$	$\mu_{23}$	$\mu_3$	$\mu_{123}$	$\mu_{23}$	$\mu_{13}$	$\mu_{123}$
$\mu_{12}$	$\mu_{12}$	$\mu_{12}$	$\mu_{123}$	$\mu_{12}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$
$\mu_{23}$	$\mu_{123}$	$\mu_{23}$	$\mu_{23}$	$\mu_{123}$	$\mu_{23}$	$\mu_{123}$	$\mu_{123}$
$\mu_{13}$	$\mu_{13}$	$\mu_{123}$	$\mu_{13}$	$\mu_{123}$	$\mu_{123}$	$\mu_{13}$	$\mu_{123}$
$\mu_{123}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$	$\mu_{123}$

(b)

**Fig. 9.3.** Shaded regions correspond to the  $4 \times 4$  minors  $M_{\{12\}}[\mu]$  in panel (a), and  $M_{\{13\}}[\mu]$  in panel (b) are constrained by the Sherali-Adams relaxation at order 2. Also constrained is the minor  $M_{\{23\}}[\mu]$  (not shown).

$h \in E$ , let  $M_{\{h\}}[\mu]$  denote the  $|\mathcal{I}(h)| \times |\mathcal{I}(h)|$  minor of  $M_m[\mu]$ , and define the semidefinite constraint set

$$\mathbb{S}(\{h\}) = \left\{ \mu \in \mathbb{R}^{|\mathcal{I}^m|} \mid M_{\{h\}}[\mu] \succeq 0 \right\}. \tag{9.17}$$

As a special case of Proposition 14 for  $m' = |h|$ , we conclude that the projection of this constraint set down to the mean parameters indexed by  $\mathcal{I}(h)$ —namely, the set  $\Pi_h(\mathbb{S}(\{h\}))$ —is equal to the local marginal polytope  $\mathbb{M}(\{h\})$ . We have thus established the following:

**Proposition 15.** *Given a hypergraph  $G$  with maximal hyperedge size  $|h| = t + 1$ , an equivalent representation of the polytope  $\mathbb{L}_t(G)$  is in terms of the family of hyperedge-based semidefinite constraints:*

$$\mathbb{L}_t(G) = \bigcap_{h \in E} \mathbb{S}(\{h\}), \tag{9.18}$$



obtained by imposing semidefiniteness on certain minors of the full moment matrix  $M_m[\mu]$ . This relaxation is tight when  $G$  is a hypertree of width  $t$ .

Figure 9.3 provides an illustration of the particular minors that are constrained by the relaxation  $\mathbb{L}_t(G)$ , in the special case of the single cycle  $K_3$  on  $m = 3$  nodes with binary variables (see Figure 8.1). Each panel in Figure 9.3 shows the full 8 times 8 moment matrix  $M_3[\mu]$  associated with this graph; the shaded portions in each panel demarcate the minors  $M_{\{12\}}[\mu]$  and  $M_{\{13\}}[\mu]$  constrained in the first-order LP relaxation  $\mathbb{L}(G)$ . In addition, this LP relaxation also constrains the minor  $M_{\{23\}}[\mu]$ , not shown here.

Thus, the framework of moment matrices allows us to understand both the hierarchy of LP relaxations defined by the hypertree-based polytopes  $\mathbb{L}_t(G)$ , as well as the SDP hierarchy based on the non-polyhedral sets  $\mathbb{S}_t(G)$ . It is worth noting that at least for general graphs with  $m \geq 3$  vertices, this hierarchy of LP and SDP relaxations are mutually incomparable, in that neither one dominates the other at a fixed level of the hierarchy. We illustrate this mutual incomparability in the following example:

**Example 45 (Mutual incomparability of LP/SDP relaxations).**

For the single cycle  $K_3$  on  $m = 3$  vertices, neither the first-order LP relaxation  $\mathbb{L}_1(G)$  nor the first-order semidefinite relaxation  $\mathbb{S}_1(K_3)$  are exact. At the same time, these two relaxations are mutually incomparable, in that neither dominates the other. In one direction, Example 43 provides an instance of a pseudomarginal vector  $\tau$  that belongs to  $\mathbb{L}_1(K_3)$ , but violates the semidefinite constraint defining  $\mathbb{S}_1(K_3)$ . In the other direction, we need to construct a pseudomarginal vector  $\tau$  that satisfies the semidefinite constraint  $M_1[\tau] \succeq 0$ , but violates at least one of the linear inequalities defining  $\mathbb{L}_1(K_3)$ . Consider the pseudomarginal  $\tau$  with moment matrix  $M_1[\tau]$  of the form

$$M_1[\tau] = \begin{bmatrix} 1 & \tau_1 & \tau_2 & \tau_3 \\ \tau_1 & \tau_1 & \tau_{12} & \tau_{13} \\ \tau_2 & \tau_{12} & \tau_2 & \tau_{23} \\ \tau_3 & \tau_{13} & \tau_{23} & \tau_3 \end{bmatrix} = \begin{bmatrix} 1 & 0.75 & 0.75 & 0.75 \\ 0.75 & 0.75 & 0.5 & 0.5 \\ 0.75 & 0.5 & 0.75 & 0.5 \\ 0.75 & 0.5 & 0.5 & 0.75 \end{bmatrix}$$

Calculation shows that  $M_1[\tau]$  is positive definite, yet the constraint  $1 + \tau_{12} - \tau_1 - \tau_2 \geq 0$  from the conditions (9.12), necessary for membership in  $\mathbb{L}_1(K_3)$ , is not satisfied. Hence, the constructed vector  $\tau$  belongs to  $\mathbb{S}_1(K_3)$  but does not belong to  $\mathbb{L}_1(K_3)$ . ♣

Of course, for hypergraphs of low treewidth, the LP relaxation is definitely advantageous, in that it is tight once the order of relaxation  $t$  hits the hypergraph treewidth—that is, the equality  $\mathbb{L}_t G = \mathbb{M}(G)$  holds for any hypergraph of treewidth  $t$ , as shown in Proposition 13. In contrast, the semidefinite relaxation  $\mathbb{S}_t(G)$  is not tight for a general hypergraph of width  $t$ ; as concrete examples, see Example 44 for the failure of  $\mathbb{S}_1(G)$  for the graph  $G = K_2$ , which has treewidth  $t = 1$ , and Example 45 for the failure of  $\mathbb{S}_2(G)$  for the graph  $G = K_3$ , which has treewidth  $t = 2$ . However, if the order of semidefinite relaxation is raised to *one order higher* than the treewidth, then the semidefinite relaxation is tight:

**Proposition 16.** *For a random vector  $X \in \{0, 1, \dots, r - 1\}^m$  Markov with respect to any hypergraph  $G$ , we always have*

$$\mathbb{S}_{t+1}(G) \subseteq \mathbb{L}_t(G),$$

where the inclusion is strict unless  $G$  is a hypertree of width  $t$ . For any such hypertree, the equality  $\mathbb{S}_{t+1}(G) = \mathbb{M}(G)$  holds.

*Proof.* Each maximal hyperedge  $h$  in a hypergraph  $G$  with treewidth  $t$  has cardinality  $|h| = t + 1$ . Note that the moment matrix  $M_{t+1}[\mu]$  constrained by the relaxation  $\mathbb{S}_{t+1}(G)$  includes as minors the  $|\mathcal{I}(h)| \times |\mathcal{I}(h)|$  matrix  $M_{\{h\}}[\mu]$ , for every hyperedge  $h \in E$ . This fact implies that  $M_{t+1}(G) \subseteq \mathbb{S}(\{h\})$  for each hyperedge  $h \in E$ , where the set  $\mathbb{S}(\{h\})$  was defined in equation (9.17). Hence, the asserted inclusion follows from Proposition 15. For hypergraphs of treewidth  $t$ , the equality  $\mathbb{S}_{t+1}(G) = \mathbb{M}(G)$  follows from this inclusion, and the equivalence between  $\mathbb{L}_t(G)$  and  $\mathbb{M}(G)$  for hypergraphs of width  $t$ , as asserted in Proposition 13. □

## 9.4 Second-order cone relaxations

In addition to the LP and SDP relaxations discussed thus far, another class of constraints on marginal polytopes are based on the second-order cone (see Appendix A.2). The resulting second-order cone programming (SOCP) relaxations of the mode-finding problem (as well as more general non-convex optimization problems) have been studied by various researchers [143, 131]. In the framework of moment matrices, second-order cone constraints can be understood as a particular weakening of a positive semidefiniteness constraint, as we describe here.

The semidefinite relaxations that we have developed are based on enforcing that certain moment matrices  $\Sigma = \mathbb{E}[YY^T]$  be positive semidefinite, where the random vector  $Y$  was an appropriately chosen function of  $(X_1, \dots, X_m)$ . As we have noted, the semidefinite constraint  $\Sigma \succeq 0$  is equivalent to an infinite number of linear constraints on the elements of  $\Sigma$ . Recall that a matrix  $\Sigma$  is positive semidefinite if and only if its Frobenius inner product  $\langle\langle \Sigma, \Lambda \rangle\rangle := \text{trace}(\Sigma\Lambda)$  with any other positive semidefinite matrix  $\Lambda \in \mathcal{S}_+^d$  is non-negative [111]; this fact corresponds to the *self-duality* of the positive semidefinite cone (e.g., [36]). By the singular value decomposition [111], any  $\Lambda \in \mathcal{S}_+^d$  can be written as  $\Lambda = UU^T$  for some matrix  $U \in \mathbb{R}^{d \times k}$  with  $k \leq d$ . Using this notation, we can rewrite the constraint  $\langle\langle \Lambda, \Sigma \rangle\rangle \geq 0$  as

$$\langle\langle \Lambda, \mathbb{E}[YY^T] \rangle\rangle \geq \|U^T \mathbb{E}[Y]\|^2. \quad (9.19)$$

For any fixed  $\Lambda$ , the inequality (9.19) corresponds to a second-order cone constraint on the elements of the second-order moment matrix

$$M_1[\mu] := \mathbb{E} \left[ \begin{bmatrix} 1 \\ Y \end{bmatrix} \begin{bmatrix} 1 & Y \end{bmatrix} \right]. \quad (9.20)$$

By the self-duality condition, imposing the constraint (9.19) for *all* matrices  $\Lambda \in \mathcal{S}_+^d$  is equivalent to enforcing that the matrix  $\Sigma = \mathbb{E}[YY^T] - \mathbb{E}[Y]\mathbb{E}[Y^T]$  is positive semidefinite. Using the Schur complement formula [111], a little calculation shows that the condition  $\Sigma \succeq 0$  is equivalent to the condition  $M_1[\mu] \succeq 0$ .

The SOCP approach is based on imposing the constraint (9.19) only for a subset of positive semidefinite matrices  $\Lambda$ , so that SOCPs are

weaker than the associated SDP relaxation. The benefit of this weakening is that many SOCPs can be solved with lower computational complexity than semidefinite programs [36]. Kim and Kojima [131] study the use of SOC constraints in deriving relaxations for fairly general classes of non-convex quadratic programming problems. Kumar et al. [143] study the use of SOCPs for mode-finding in pairwise Markov random fields, in particular using matrices  $\Lambda$  that are locally defined on the edges of the graph, as well as additional linear inequalities, such as the cycle inequalities in the binary case (see Example 42). In later work, Kumar et al. [142] showed that one form of their SOCP relaxation is equivalent to a form of the quadratic programming (QP) relaxation proposed by Lafferty and Ravikumar [194]. Kumar et al. [142] also provide a cautionary message, by demonstrating that certain classes of SOCP constraints fail to improve upon the first-order tree-based LP relaxation (8.5). An example of such redundant SOCP constraints are those in which the matrix  $\Lambda$  has non-zero entries only in pairs of elements, corresponding to a single edge, or more generally, is associated with a sub-tree of the original graph on which the MRF is defined. This result can be established using moment matrices by the following reasoning: any SOC constraint (9.19) that constrains only elements associated with some sub-tree of the graph is redundant with the first-order LP constraints (8.5), since the LP constraints ensure that any moment vector  $\mu$  is consistent on any sub-tree embedded within the graph. Of course, there are higher-order parallels to this statement when dealing with high-order LP relaxations, guaranteed to be tight on hypertrees up to a given treewidth (e.g., see Proposition 13.) Kumar et al. [142] also give other cycle-structured SOC constraints that are redundant in certain cases (in particular, for certain settings of the parameters that define the mode-finding problem).

# 10

---

## Discussion

---

The core of this paper is a general set of variational principles for the problems of computing marginal probabilities and modes, applicable to multivariate statistical models in the exponential family. A fundamental object underlying these optimization problems is the set of realizable mean parameters associated with the exponential family; indeed, the structure of this set largely determines whether or not the associated variational problem can be solved exactly in an efficient manner. Moreover, a large class of well-known algorithms for both exact and approximate inference—including mean field methods, the sum-product and max-product algorithms, as well as generalizations thereof—can be derived and understood as methods for solving various forms (either exact or approximate) of these variational problems. The variational perspective also suggests convex relaxations of the exact principle, which in turn lead to new algorithms for approximate inference.

Many of the algorithms described in this paper are already essential tools in various practical applications (e.g., the sum-product algorithm in error-correcting coding). While such empirical successes underscore the promise of variational approaches, a variety of theoretical questions remain to be addressed. One important direction to pursue is

obtaining *a priori* guarantees on the accuracy of a given variational method for a particular subclass of problems. For instance, it remains to be seen whether techniques used to obtain performance guarantees for relaxations of combinatorial optimization problems can be adapted to analyze other types of inference problems (e.g., computing approximate marginal distributions). Another major area with various open issues is the application of variational methods to parameter estimation. Although mean field methods are already widely used for parameter estimation in directed graphical models, open questions include how to exploit more powerful variational methods, and also how to deal with undirected graphical models. Variational methods that provide upper bounds on the cumulant function are likely to be useful for parameter estimation in the undirected setting.

We have focused in this article on regular exponential families where the parameters lie in an open linear space. This class covers a broad class of graphical models, particularly undirected graphical models, where clique potentials often factor into products over parameters and sufficient statistics. However, there are also examples in which nonlinear constraints are imposed on the parameters in a graphical model; this often arises in particular in the directed graphical model setting. Such constraints require the general machinery of *curved exponential families* [71]. While there have been specialized examples of the application of variational methods for such families [119], there does not yet exist a general treatment of variational methods for curved exponential families.

Another direction that requires further exploration is the study of variational inference for distributions outside of the exponential family. In particular, it is of interest to develop variational methods for the stochastic processes that underlie nonparametric Bayesian modeling. Again, special cases of variational inference have been presented for such models—see in particular the work of Blei and Jordan [28] on variational inference for Dirichlet processes—but there is as of yet no general framework.

Finally, it should be emphasized that the variational approach provides a set of techniques that are complementary to Monte Carlo methods. One interesting program of research, then, is to characterize the

classes of problems for which variational methods (or conversely, Monte Carlo methods) are best suited, and moreover to develop a theoretical characterization of the trade-offs in complexity versus accuracy inherent to each method.

### **Acknowledgements**

A large number of people contributed to the gestation of this paper, and it is a pleasure to acknowledge them here. The intellectual contributions and support of Alan Willsky and Tommi Jaakkola were particularly significant in the development of the ideas presented here. In addition, we thank the following individuals for their comments and insights along the way: Constantine Caramanis, Laurent El Ghaoui, Jon Feldman, G. David Forney Jr., David Karger, John Lafferty, Adrian Lewis, Jon McAuliffe, Michal Rosen-Zvi, Lawrence Saul, Nathan Srebro, Sekhar Tatikonda, Romain Thibaux, Yee Whye Teh, Lieven Vandenberghe, Yair Weiss and Jonathan Yedidia.

# A

---

## Background material

---

In this appendix, we provide some basic background on graph theory, and convex sets and functions.

### A.1 Background on graphs and hypergraphs

Here we collect some basic definitions and results on graphs and hypergraphs; see the standard books [15, 31, 32] for further background. A graph  $G = (V, E)$  consists of a set  $V = \{1, 2, \dots, m\}$  of vertices, and a set  $E \subset V \times V$  of edges. By definition, a graph does not contain self-loops (i.e.,  $(s, s) \notin E$  for all vertices  $s \in V$ ), nor does it contain multiple copies of the same edge. (These features are permitted only in the extended notion of a multigraph.) For a *directed graph*, the ordering of the edge matters, meaning that  $(s, t)$  is distinct from  $(t, s)$ , whereas for an *undirected graph*, the quantities  $(s, t)$  and  $(t, s)$  refer to the same edge. A *subgraph*  $F$  of a given graph  $G = (V, E)$  is a graph  $(V(F), E(F))$  such that  $V(F) \subseteq V$  and  $E(F) \subseteq E$ . Given a subset  $S \subseteq V$  of the vertex set of a graph  $G = (V, E)$ , the associated *vertex-induced subgraph* is the subgraph  $F[S] := (S, E(S))$  with edge-set  $E(S) = \{(s, t) \in E \mid (s, t) \in E\}$ . Given a subset of edges



$F \subseteq E$ , the associated *edge-induced subgraph* is  $F(F) = (V(F), F)$ , where  $V(F) = \{s \in V \mid (s, t) \in F\}$ . A *clique* of a graph is a vertex-induced subgraph  $F[S]$  that is completely connected (i.e., all vertices  $s, t \in S$  are joined by an edge  $(s, t) \in E$ ).

A *path*  $P$  is a graph  $P = (V(P), E(P))$  with vertex set  $V(P) := \{v_0, v_1, v_2, \dots, v_k\}$ , and edge set

$$E(P) := \{(v_0, v_1), (v_1, v_2), \dots, (v_{k-1}, v_k)\}.$$

We say that the path  $P$  joins vertex  $v_0$  to vertex  $v_k$ . Of particular interest are paths that are subgraphs of a given graph  $G = (V, E)$ , meaning that  $V(P) \subseteq V$  and  $E(P) \subseteq E$ . A *cycle* is a graph  $C = (V(C), E(C))$  with  $V(C) = \{v_0, v_1, \dots, v_k\}$  with  $k \geq 2$  and edge set  $E(C) = \{(v_0, v_1), (v_1, v_2), \dots, (v_{k-1}, v_k), (v_k, v_0)\}$ . An undirected graph is *acyclic* if it contains no cycles. A graph is *bipartite* if its vertex set can be partitioned as a disjoint union  $V = V_A \uplus V_B$  (where  $\uplus$  denotes disjoint union), such that  $(s, t) \in E$  implies that  $s \in V_A$  and  $t \in V_B$  (or vice versa).

A *chord* of a cycle  $C$  on the vertex set  $V(C) = \{v_0, v_1, \dots, v_k\}$  is an edge  $(v_i, v_j)$  that is *not* part of the cycle edge set  $E(C) = \{(v_0, v_1), (v_1, v_2), \dots, (v_{k-1}, v_k), (v_k, v_0)\}$ . A cycle  $C$  of length 4 within greater contained with a larger graph  $G$  is *chordless* if the graph contains no chords for the cycle. A graph is *triangulated* if it contains no cycles of length four or greater that are chordless.

A *connected component* of a graph is a subset of vertices  $S \subseteq V$  such that for all  $s, t \in S$ , there exists a path contained within the graph  $G$  joining  $s$  to  $t$ . We say that a graph  $G$  is *singly connected*, or simply *connected*, if it consists of a single connected component. A *tree* is an acyclic graph with a single connected component; it can be shown by induction that any tree on  $m$  vertices must have  $m - 1$  edges. More generally, a *forest* is an acyclic graph consisting of one or more connected components.

A graph is *triangulated*

A *hypergraph* is  $G = (V, E)$  is a natural generalization of a graph: it consists of a vertex set  $V = \{1, 2, \dots, m\}$ , and a set  $E$  of *hyperedges*, which each hyperedge  $h \in E$  is a particular subset of  $V$ . (Thus, an ordinary graph is the special case in which  $|h| = 2$  for each hyperedge.)

The *factor graph* associated with a hypergraph  $G$  is a bipartite graph  $(V', E')$  with vertex set  $V'$  corresponding to the union  $V \cup E$  of the hyperedge vertex set  $V$ , and the set of hyperedges  $E$ , and an edge set  $E'$  that includes the edge  $(s, h)$ , where  $s \in V$  and  $h \in E$ , if and only if the hyperedge  $h$  includes vertex  $s$ .

## A.2 Basics of convex sets and functions

Here we collect some basic definitions and results on convex sets and functions. Rockafellar [198] is a standard reference on convex analysis; see also the books by Hiriart-Urruty and Lemaréchal [109, 110], Boyd and Vandenberghe [36], and Bertsekas [21].

### A.2.1 Convex sets and cones

A set  $C \subseteq \mathbb{R}^d$  is *convex* if for all  $x, y \in C$  and  $\alpha \in [0, 1]$ , the set  $C$  also contains the point  $\alpha x + (1 - \alpha)y$ . Equivalently, the set  $C$  must contain the entire line segment joining any two of its elements. Note that convexity is preserved by various operations on sets:

- the *intersection*  $\bigcap_{j \in \mathcal{J}} C_j$  of a family  $\{C_j\}_{j \in \mathcal{J}}$  of convex sets remains convex.
- the *Cartesian product*  $C_1 \otimes C_2 \otimes \dots \otimes C_k$  of a family of convex sets is convex.
- letting  $f(x) = Ax + b$  be an affine mapping, then the image  $f(C)$  of a convex set  $C$  is also convex; and
- the *closure*  $\bar{C}$  and *interior*  $C^\circ$  of a convex set  $C$  are also convex sets.

Note that the union of convex sets is not convex, in general.

A *cone*  $K$  is a set such for any  $x \in K$ , the ray  $\{\lambda x \mid \lambda > 0\}$  also belongs to  $K$ . A *convex cone* is a cone that is also convex. (To appreciate the distinction, the union of two distinct rays is a cone, but not convex in general.) Important examples of convex cones include:

- any subspace  $\{x \in \mathbb{R}^d \mid Ax = 0\}$  for some matrix  $A \in \mathbb{R}^{m \times d}$ .

- the non-negative orthant

$$\{x \in \mathbb{R}^d \mid x_1 \geq 0, x_2 \geq 0, \dots, x_d \geq 0\}.$$

- the *conical hull* of a collection of vectors  $\{x_1, \dots, x_n\}$ :

$$\{y \in \mathbb{R}^d \mid y = \sum_{i=1}^n \lambda_i x_i \quad \text{with } \lambda_i \geq 0, i = 1, \dots, n\}. \quad (\text{A.1})$$

- the second-order cone  $\{(x, t) \in \mathbb{R}^{d-1} \times \mathbb{R} \mid \|x\|_2 \leq t\}$ ;
- the cone  $\mathcal{S}_+^d$  of symmetric positive semidefinite matrices

$$\mathcal{S}_+^d = \left\{ X \in \mathbb{R}^{d \times d} \mid X = X^T, X \succeq 0 \right\}. \quad (\text{A.2})$$

### A.2.2 Convex and affine hulls

A *linear combination* of elements  $x_1, x_2, \dots, x_k$  from a set  $S$  is a sum  $\sum_{i=1}^k \alpha_i x_i$ , for arbitrary scalars  $\alpha_i \in \mathbb{R}$ . An *affine combination* is a linear combination with the further restriction that  $\sum_{i=1}^k \alpha_i = 1$ , whereas a *convex combination* is a linear combination with the further restrictions  $\sum_{i=1}^k \alpha_i = 1$  and  $\alpha_i \geq 0$  for all  $i = 1, \dots, k$ . The *affine hull* of a set  $S$ , denoted  $\text{aff}(S)$ , is the smallest set that contains all affine combinations. Similarly, the *convex hull* of a set  $S$ , denoted  $\text{conv}(S)$ , is the smallest set that contains all its convex combinations. Note that  $\text{conv}(S)$  is a convex set, by definition.

### A.2.3 Affine hulls and relative interior

For any  $\epsilon > 0$  and vector  $z \in \mathbb{R}^d$ , define the Euclidean ball

$$\mathbb{B}_\epsilon(z) := \{y \in \mathbb{R}^d \mid \|y - z\| < \epsilon\}. \quad (\text{A.3})$$

The *interior* of a convex set  $C \subseteq \mathbb{R}^d$ , denoted by  $C^\circ$ , is given by

$$C^\circ := \{z \in C \mid \exists \epsilon > 0 \text{ s.t. } \mathbb{B}_\epsilon(z) \subset C\}. \quad (\text{A.4})$$

The relative interior is defined similarly, except that the interior is taken with respect to the affine hull of  $C$ , denoted  $\text{aff } C$ . More formally, the *relative interior* of  $C$ , denoted  $\text{ri}(C)$ , is given by

$$\text{ri}(C) := \{z \in C \mid \exists \epsilon > 0 \text{ s.t. } \mathbb{B}_\epsilon(z) \cap \text{aff}(C) \subset C\}. \quad (\text{A.5})$$

To illustrate the distinction, note that the interior of the convex set  $[0, 1]$ , when viewed as a subset of  $\mathbb{R}^2$ , is empty. In contrast, the affine hull of  $[0, 1]$  is the real line, so that the relative interior is the open interval  $(0, 1)$ .

A key property of any convex set  $C$  is that its relative interior is always non-empty [198]. A convex set  $C \subseteq \mathbb{R}^d$  is *full-dimensional* if its affine hull is equal to  $\mathbb{R}^d$ . For instance, the interval  $[0, 1]$ , when viewed as a subset of  $\mathbb{R}^2$ , is not full-dimensional, since its affine hull is only the real line. For a full-dimensional convex set, the notion of interior and relative interior coincide.

#### A.2.4 Polyhedra and their representations

A *polyhedron*  $P$  is a set that can be represented as the intersection of a finite number of half-spaces—that is, if  $\{(a_j, b_j) \in \mathbb{R}^d \times \mathbb{R}, j \in \mathcal{J}\}$  is a collection of halfspaces, then

$$P = \{x \in \mathbb{R}^d \mid \langle a_j, x \rangle \leq b_j \quad \forall j \in \mathcal{J}\} \quad (\text{A.6})$$

We refer to a bounded polyhedron as a *polytope*. Given a polyhedron  $P$ , we say that  $x \in P$  is an *extreme point* if it is not possible to write  $x = \lambda y + (1 - \lambda)z$  for some  $\lambda \in (0, 1)$  and  $y, z \in P$ . We say that  $x$  is a *vertex* if there exists some  $c \in \mathbb{R}^d$  such that  $\langle c, x \rangle > \langle c, y \rangle$  for all  $y \in P$  with  $y \neq x$ . For a polyhedron,  $x$  is a vertex if and only if it is an extreme point. A consequence of the Minkowski-Weyl theorem is that any non-empty polytope can be written as the convex hull of its extreme points. This convex hull representation is dual to the half-space representation. Conversely, the convex hull of any finite collection of vectors is a polytope, and so has a half-space representation (A.6).

#### A.2.5 Convex functions

It is convenient to allow convex functions to take the value  $+\infty$ , particularly in performing dual calculations. More formally, an *extended real-valued function*  $f$  on  $\mathbb{R}^d$  takes values in the extended reals  $\mathbb{R}_* := \mathbb{R} \cup \{+\infty\}$ . A function  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  is *convex* if for all  $x, y \in \mathbb{R}^d$  and  $\alpha \in [0, 1]$ , we have

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y). \quad (\text{A.7})$$

The function is *strictly convex* if the inequality (A.7) holds strictly for all  $x \neq y$  and  $\alpha \in (0, 1)$ .

The *domain* of an extended real-valued convex function  $f$  is the set

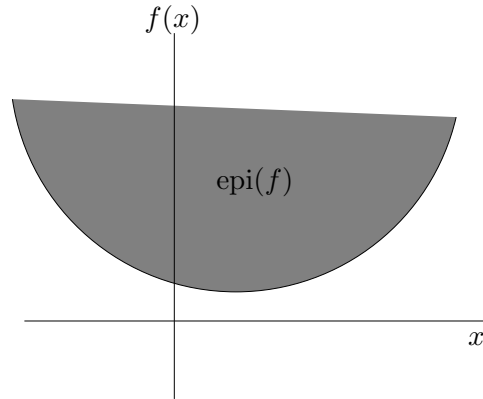
$$\text{dom}(f) := \left\{ x \in \mathbb{R}^d \mid f(x) < +\infty \right\}. \quad (\text{A.8})$$

Note that the domain is always a convex subset of  $\mathbb{R}^d$ . Throughout, we restrict our attention to *proper* convex functions, meaning that  $\text{dom}(f)$  is non-empty.

A key consequence of the definition (A.7) is *Jensen's inequality*, which says that for any convex combination  $\sum_{i=1}^k \alpha_i x_i$  of points  $\{x_i\}$  in the domain of  $f$ , we have  $f\left(\sum_{i=1}^k \alpha_i x_i\right) \leq \sum_{i=1}^k \alpha_i f(x_i)$ .

Convexity of functions is preserved by various operations. In particular, given a collection of convex functions  $\{f_j \mid j \in \mathcal{J}\}$ :

- any linear combination  $\sum_{j \in \mathcal{J}} \alpha_j f_j$  with  $\alpha_i \geq 0$  is a convex function;
- the pointwise supremum  $f(x) := \sup_{j \in \mathcal{J}} f_j(x)$  is also a convex function.



**Fig. A.1.** The epigraph of a convex function is a subset of  $\mathbb{R}^d \times \mathbb{R}$ , given by  $\text{epi}(f) = \{(x, t) \in \mathbb{R}^d \times \mathbb{R} \mid f(x) \leq t\}$ . For a convex function, this set is always convex.

The convexity of a function can also be defined in terms of its *epi-*

*graph*, namely the subset of  $\mathbb{R}^d \times \mathbb{R}$  given by

$$\text{epi}(f) := \left\{ (x, t) \in \mathbb{R}^d \times \mathbb{R} \mid f(x) \leq t \right\}, \quad (\text{A.9})$$

as illustrated in Figure A.1. Many properties of convex functions can be stated in terms of properties of this epigraph. For instance, it is a straightforward exercise to show that the function  $f$  is convex if and only if its epigraph is a convex subset of  $\mathbb{R}^d \times \mathbb{R}$ . A convex function is *lower semi-continuous* if  $\lim_{y \rightarrow x} f(y) \geq f(x)$  for all  $x$  in its domain. It can be shown that a convex function  $f$  is lower semi-continuous if and only if its epigraph is a closed set; in this case,  $f$  is said to be a *closed convex* function.

### A.2.6 Conjugate duality

Given a convex function  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  with non-empty domain, its conjugate dual is a new extended real-valued function  $f^* : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ , defined as

$$\begin{aligned} f^*(y) &:= \sup_{x \in \mathbb{R}^d} \{ \langle y, x \rangle - f(x) \} \\ &= \sup_{x \in \text{dom}(f)} \{ \langle y, x \rangle - f(x) \}, \end{aligned} \quad (\text{A.10})$$

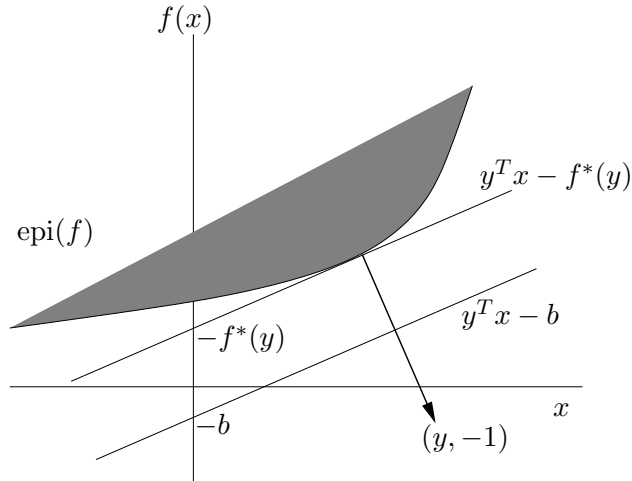
where the second equality follows by definition of the domain of  $f$ . Note that  $f^*$  is always a convex function, since it is the pointwise supremum of a family of convex functions. Geometrically, this operation can be interpreted as computing the intercept of the supporting hyperplane to  $\text{epi}(f)$  with normal vector  $(y, -1) \in \mathbb{R}^d \times \mathbb{R}$ , as illustrated in Figure A.2.

It is also possible to take the conjugate dual of this dual function thereby yielding a new function known as the biconjugate

$$(f^*)^*(z) = \sup_{y \in \mathbb{R}^d} \{ \langle z, y \rangle - f^*(y) \}. \quad (\text{A.11})$$

An important property of conjugacy is that whenever  $f$  is a closed convex function, then the biconjugate  $f^{**}$  is equal to the conjugate  $f^*$ .

For closed, convex, and differentiable functions  $f$  that satisfy additional technical conditions, the conjugate pair  $(f, f^*)$  induce a one-to-



**Fig. A.2.** Interpretation of the conjugate dual function in terms of supporting hyperplanes to the epigraph. The negative value of the dual function  $-f^*(y)$  corresponds to the intercept of the supporting hyperplane to  $\text{epi}(f)$  with normal vector  $(y, -1) \in \mathbb{R}^d \times \mathbb{R}$ .

one correspondence between  $\text{dom}(f)$  and  $\text{dom}(f^*)$ , based on the gradient mappings  $\nabla f$  and  $\nabla f^*$ . In particular, we say that a convex function  $f$  is *essentially smooth* if it has a non-empty domain  $C = \text{dom}(f)$ ,  $f$  is differentiable throughout  $C^\circ$ , and  $\lim_{i \rightarrow \infty} \nabla f(x^i) = +\infty$  for any sequence  $\{x^i\}$  contained in  $C$ , and converging to a boundary point of  $C$ . Note that a convex function with domain  $\mathbb{R}^d$  is always essentially smooth, since its domain has no boundary. This property is also referred to as *steepness* in some statistical treatments [40].

When  $f$  is strictly convex, differentiable and essentially smooth, then defining  $C = \text{dom}(f)$ , then we say that  $(f, C^\circ)$  is a convex function of *Legendre type*. Letting  $D = \text{dom}(f^*)$ , it can be shown that  $(f, C^\circ)$  is a convex function of Legendre type if and only if  $(f^*, D^\circ)$  is a convex function of Legendre type. In this case, the gradient mapping  $\nabla f$  is one-to-one from  $C^\circ$  onto the open convex set  $D^\circ$ , continuous in both directions, and moreover  $\nabla f^* = (\nabla f)^{-1}$ . See Section 26 of Rockafellar [198] for further details on these types of Legendre correspondences between  $f$  and its dual  $f^*$ .

# B

---

## Proofs and auxiliary results: Exponential families and duality

---

In this section, we provide the proofs of Theorems 1 and 2 from Section 3, as well as some auxiliary results of independent interest.

### B.1 Proof of Theorem 1

We prove the result first for a minimal representation, and then discuss its extension to the overcomplete case. Recall from Appendix A.2.3 that a convex subset of  $\mathbb{R}^d$  is full-dimensional if its affine hull is equal to  $\mathbb{R}^d$ . We first note that  $\mathcal{M}$  is a full-dimensional convex set if and only if the exponential family representation is minimal. Indeed, the representation is *not* minimal if and only if there exists some vector  $a \in \mathbb{R}^d$  and constant  $b \in \mathbb{R}$  such that  $\langle a, \phi(x) \rangle = b$  holds  $\nu$ -a.e.. By definition of  $\mathcal{M}$ , this equality holds if and only if  $\langle a, \mu \rangle = b$  for all  $\mu \in \mathcal{M}$ , which is equivalent to  $\mathcal{M}$  not being full-dimensional.

Thus, if we restrict attention to minimal exponential families for the time-being, the set  $\mathcal{M}$  is full-dimensional. Our proof makes use of the following properties of a full-dimensional convex set [see 109, 198]: (a) its interior  $\mathcal{M}^\circ$  is non-empty, and the interior of the closure  $[\overline{\mathcal{M}}]^\circ$  is equal to the interior  $\mathcal{M}^\circ$ ; and (b) the interior  $\mathcal{M}^\circ$  contains the zero-



vector 0 if and only if for all non-zero  $\gamma \in \mathbb{R}^d$ , there exists some  $\mu \in \mathcal{M}$  with  $\langle \gamma, \mu \rangle > 0$ .

$\nabla A(\Omega) \subseteq \mathcal{M}^\circ$ : By shifting the potential  $\phi$  by a constant vector if necessary, it suffices to consider the case  $0 \in \nabla A(\Omega)$ . Let  $\theta^0 \in \Omega$  be the associated canonical parameter satisfying  $\nabla A(\theta^0) = 0$ . We prove that for all non-zero directions  $\gamma \in \mathbb{R}^d$ , there is some  $\mu \in \mathcal{M}$  such that  $\langle \gamma, \mu \rangle > 0$ , which implies  $0 \in \mathcal{M}^\circ$  by property (b).

For any  $\gamma \in \mathbb{R}^d$ , the openness of  $\Omega$  ensures the existence of some  $\delta > 0$  such that  $(\theta^0 + \delta\gamma) \in \Omega$ . Using the strict convexity and differentiability of  $A$  on  $\Omega$  and the fact that  $\nabla A(\theta^0) = 0$  by assumption, there holds  $A(\theta^0 + \delta\gamma) > A(\theta^0) + \langle \nabla A(\theta^0), \delta\gamma \rangle = A(\theta^0)$ . Similarly, defining  $\mu^\delta := \nabla A(\theta^0 + \delta\gamma)$ , we can write  $A(\theta^0) > A(\theta^0 + \delta\gamma) + \langle \mu^\delta, -\delta\gamma \rangle$ . These two inequalities in conjunction imply that

$$\delta \langle \mu^\delta, \gamma \rangle > A(\theta^0 + \delta\gamma) - A(\theta^0) > 0.$$

Since  $\mu^\delta \in \nabla A(\Omega) \subseteq \mathcal{M}$  and  $\gamma \in \mathbb{R}^d$  was arbitrary, this establishes that  $0 \in \mathcal{M}^\circ$ .

$\mathcal{M}^\circ \subseteq \nabla A(\Omega)$ : As in the preceding argument, we may take  $0 \in \mathcal{M}^\circ$  without loss of generality. Then, we must establish the existence of  $\theta \in \Omega$  such that  $\nabla A(\theta) = 0$ . By convexity, it is equivalent to show that  $\inf_{\theta \in \Omega} A(\theta)$  is attained. To establish the attainment of this infimum, we prove that  $A$  has no directions of recession, meaning that  $\lim_{n \rightarrow +\infty} A(\theta^n) = +\infty$  for all sequences  $\{\theta^n\}$  such that  $\|\theta^n\| \rightarrow +\infty$ .

For an arbitrary non-zero direction  $\gamma \in \mathbb{R}^d$  and  $\epsilon > 0$ , consider the half space  $H_{\gamma, \epsilon} := \{x \in \mathcal{X}^m \mid \langle \gamma, \phi(x) \rangle \geq \epsilon\}$ . Since  $0 \in \mathcal{M}^\circ$ , this half-space must have positive measure under  $\nu$  for all sufficiently small  $\epsilon > 0$ . Otherwise, the inequality  $\langle \gamma, \phi(x) \rangle \leq 0$  would hold  $\nu$ -a.e., which implies that  $\langle \gamma, \mu \rangle \leq 0$  for all  $\mu \in \overline{\mathcal{M}}$ . By the convexity of  $\mathcal{M}$ , this inequality would imply that  $0 \notin [\overline{\mathcal{M}}]^\circ = \mathcal{M}^\circ$ , which contradicts our starting assumption.

For an arbitrary  $\theta^0 \in \Omega$ , we now write

$$\begin{aligned} A(\theta^0 + t\gamma) &\geq \log \int_{H_{\gamma,\epsilon}} \exp \{ \langle \theta^0 + t\gamma, \phi(x) \rangle \} \nu(dx) \\ &\geq t\epsilon + \underbrace{\log \int_{H_{\gamma,\epsilon}} \exp \{ \langle \theta^0, \phi(x) \rangle \} \nu(dx)}_{C(\theta^0)}. \end{aligned}$$

Note that we must have  $C(\theta^0) > -\infty$ , because

$$\exp \{ \langle \theta^0, \phi(x) \rangle \} > 0 \quad \text{for all } x \in \mathcal{X}^m,$$

and  $\nu(H_{\gamma,\epsilon}) > 0$ . Hence, we conclude that  $\lim_{t \rightarrow +\infty} A(\theta^0 + t\gamma) = +\infty$  for all directions  $\gamma \in \mathbb{R}^d$ , showing that  $A$  has no directions of recession.

Extension to overcomplete case: For any overcomplete representation  $\phi$ , let  $\varphi$  be a set of potential functions in an equivalent minimal representation. In particular, a collection  $\varphi$  can be specified by eliminating elements of  $\phi$  until no affine dependencies remain. Let  $\nabla A_\varphi$  and  $\nabla A_\phi$  be the respective mean parameter mappings associated with  $\varphi$  and  $\phi$ , with the sets  $\mathcal{M}_\varphi$  and  $\mathcal{M}_\phi$  similarly defined. By the result just established,  $\nabla A_\varphi$  is onto the interior of  $\mathcal{M}_\varphi$ . By construction of  $\varphi$ , each member in the relative interior of  $\mathcal{M}_\phi$  is associated with a unique element in the interior of  $\mathcal{M}_\varphi$ . We conclude that the mean parameter mapping  $\nabla A_\phi$  is onto the relative interior of  $\mathcal{M}_\phi$ .

## B.2 Proof of Theorem 2

(a) Case (i)  $\mu \in \mathcal{M}^\circ$ : In this case, Theorem 1 guarantees that the inverse image  $(\nabla A)^{-1}(\mu)$  is non-empty. Any point in this inverse image attains the supremum in equation (3.42). In a minimal representation, there is only one optimizing point, whereas there is an affine subset for an overcomplete representation. Nonetheless, for any  $\theta(\mu) \in (\nabla A)^{-1}(\mu)$ , the value of the optimum is:  $A^*(\mu) = \langle \theta(\mu), \mu \rangle - A(\theta(\mu))$ . By definition (3.2) of the entropy, we have

$$\begin{aligned} -H(p_{\theta(\mu)}) &= \mathbb{E}_{\theta(\mu)}[\langle \theta(\mu), \phi(X) \rangle - A(\theta(\mu))] \\ &= \langle \theta(\mu), \mu \rangle - A(\theta(\mu)), \end{aligned}$$

where the final equality uses linearity of expectation, and the fact that  $\mathbb{E}_{\theta(\mu)}[\phi(X)] = \mu$ .

Case (ii)  $\mu \notin \overline{\mathcal{M}}$ : Let  $\text{dom } A^* = \{\mu \in \mathbb{R}^d \mid A^*(\mu) < +\infty\}$  denote the effective domain of  $A^*$ . With this notation, we must prove that  $\overline{\mathcal{M}} \supseteq \text{dom } A^*$ . In order to do so, we use Corollary 26.4.1 of Rockafellar [198] which asserts that if  $A$  is essentially smooth and lower semi-continuous, then we have the inclusions

$$[\text{dom } A^*]^\circ \subseteq \nabla A(\Omega) \subseteq \text{dom } A^*,$$

or equivalently

$$[\text{dom } A^*]^\circ \subseteq \mathcal{M}^\circ \subseteq \text{dom } A^*,$$

since  $\nabla A(\Omega) = \mathcal{M}^\circ$  from Theorem 1. Since both  $\mathcal{M}$  and  $\text{dom } A^*$  are convex sets, taking closures in these inclusions yields that  $\overline{\text{dom } A^*} = \overline{[\mathcal{M}^\circ]} = \overline{\mathcal{M}}$ , where the second equality follows by the convexity of  $\mathcal{M}$ . Therefore, by definition of the effective domain,  $A^*(\mu) = +\infty$  for any  $\mu \notin \overline{\mathcal{M}}$ .

It remains to verify that  $A$  is both lower semi-continuous, and essentially smooth. (See Appendix A.2.5 for the definition.) Recall from Proposition 2 that  $A$  is differentiable; hence, it is lower semi-continuous on its domain. To establish that  $A$  is essentially smooth, let  $\theta^b$  be a boundary point, and let  $\theta^0 \in \Omega$  be arbitrary. By Since the set  $\Omega$  is convex and open, the line  $\theta^t := t\theta^b + (1-t)\theta^0$  is contained in  $\Omega$  for all  $t \in [0, 1)$  (see Theorem 6.1 of Rockafellar [198]). Using the differentiability of  $A$  on  $\Omega$  and its convexity (see Proposition 2(b)), for any  $t < 1$ , we can write

$$A(\theta^0) \geq A(\theta^t) + \langle \nabla A(\theta^t), \theta^0 - \theta^t \rangle.$$

Re-arranging and applying the Cauchy-Schwartz inequality to the inner product term yields that

$$A(\theta^t) - A(\theta^0) \leq \|\theta^t - \theta^0\| \|\nabla A(\theta^t)\|.$$

Now as  $t \rightarrow 1^-$ , the left-hand side tends to infinity by the lower semi-continuity of  $A$ , and the regularity of the exponential family. Consequently, the right-hand side must also tend to infinity; since  $\|\theta^t - \theta^0\|$

is bounded, we conclude that  $\|\nabla A(\theta^t)\| \rightarrow +\infty$ , which shows that  $A$  is essentially smooth.

Case (iii)  $\mu \in \overline{\mathcal{M}} \setminus \mathcal{M}^\circ$ : Since  $A^*$  is defined as a conjugate function, it is lower semi-continuous. Therefore, the value of  $A^*(\mu)$  for any boundary point  $\mu \in \overline{\mathcal{M}} \setminus \mathcal{M}^\circ$  is determined by the limit over a sequence approaching  $\mu$  from inside  $\mathcal{M}^\circ$ , as claimed.

(b) From Proposition 2,  $A$  is lower semi-continuous, which ensures that  $(A^*)^* = A$ , and part (a) shows that  $\overline{\text{dom } A^*} = \overline{\mathcal{M}}$ . Consequently, we can write

$$\begin{aligned} A(\theta) &= \sup_{\mu \in \overline{\mathcal{M}}} \{\langle \theta, \mu \rangle - A^*(\mu)\} \\ &= \sup_{\mu \in \mathcal{M}} \{\langle \theta, \mu \rangle - A^*(\mu)\}, \end{aligned} \tag{B.1}$$

since it is inconsequential whether we take the supremum over  $\mathcal{M}$  or its closure  $\overline{\mathcal{M}}$ .

(c) For a minimal representation, Proposition 3 and Theorem 1 guarantee that the gradient mapping  $\nabla A$  is a bijection between  $\Theta$  and  $\mathcal{M}^\circ$ . On this basis, it follows that the gradient mapping  $\nabla A^*$  also exists and is bijective [198], whence the supremum (B.1) is attained at a unique point whenever  $\theta \in \Omega$ . The analogous statement for an over-complete representation can be proved via reduction to a minimal representation.

### B.3 General properties of $\mathcal{M}$ and $A^*$

In this section, we state and prove some auxiliary results of independent interest on general properties of the dual function  $A^*$ , and the associated set  $\mathcal{M}$  of valid mean parameters.

#### B.3.1 Properties of $\mathcal{M}$

From its definition, it is clear that  $\mathcal{M}$  is always a convex set. Other more specific properties of  $\mathcal{M}$  turn out to be determined by the properties of the exponential family. Recall from Appendix A.2.3 that a convex set

$\mathcal{M} \subseteq \mathbb{R}^d$  is full-dimensional if its affine hull is equal to  $\mathbb{R}^d$ . With this notion, we have the following:

**Proposition 17.** *The set  $\mathcal{M}$  has the following properties:*

- (a)  $\mathcal{M}$  is full-dimensional if and only if the exponential family is minimal.
- (b)  $\mathcal{M}$  is bounded if and only if  $\Theta = \mathbb{R}^d$  and  $A$  is globally Lipschitz on  $\mathbb{R}^d$ .

*Proof.* (a) The representation is *not* minimal if and only if there exists some vector  $a \in \mathbb{R}^d$  and constant  $b \in \mathbb{R}$  such that  $\langle a, \phi(x) \rangle = b$  holds  $\nu$ -a.e. By definition of  $\mathcal{M}$ , this equality holds if and only if  $\langle a, \mu \rangle = b$  for all  $\mu \in \mathcal{M}$ , which is equivalent to  $\mathcal{M}$  not being full-dimensional.

(b) The recession function  $A_\infty$  is the support function  $\sup_{\mu \in \mathcal{M}} \langle \mu, \theta \rangle$ . Therefore, the set  $\mathcal{M}$  is bounded if and only if  $A_\infty(\theta)$  is finite for all  $\theta \in \mathbb{R}^d$ . The recession function  $A_\infty$  is finite-valued if and only if  $A$  is Lipschitz and hence finite on all of  $\mathbb{R}^d$  (see Proposition 3.2.7 in Hiriart-Urruty and Lemaréchal [109]).  $\square$

The necessity of the condition  $\Omega = \mathbb{R}^d$  for  $\mathcal{M}$  to be bounded is clear from the boundary behavior of  $\nabla A$  given in Proposition 2. However, the additional global Lipschitz condition is also necessary, as demonstrated by the Poisson family (see Table 3.2). In this case, we have  $\Theta = \mathbb{R}$  yet the set of mean parameters  $\mathcal{M} = (0, +\infty)$  is unbounded. This unboundedness occurs because the function  $A(\theta) = \exp(\theta)$ , while finite on  $\mathbb{R}$ , is not globally Lipschitz.

### B.3.2 Properties of $A^*$

Despite the fact that  $A^*$  is not given in closed form, a number of properties can be inferred from its variational definition (3.42). For instance, an immediate consequence is that  $A^*$  is always convex. More specific properties of  $A^*$  depend on the nature of the exponential family. Recall from Appendix A.2.6 the definition of an essentially smooth convex function.

**Proposition 18.** *The dual function  $A^*$  is always convex and lower semi-continuous. Moreover, in a minimal and regular exponential family:*

- (a)  $A^*$  is differentiable on  $\mathcal{M}^\circ$ , and  $\nabla A^*(\mu) = (\nabla A)^{-1}(\mu)$ .
- (b)  $A^*$  is strictly convex and essentially smooth.

*Proof.* The convexity and lower semi-continuity follow because  $A^*$  is the supremum of collection of functions linear in  $\mu$ . Since both  $A$  and  $A^*$  are lower semi-continuous, the dual  $A^*$  has these properties if and only if  $A$  is strictly convex and essentially smooth (see our discussion of Legendre duality in Appendix A.2.6, and Theorem 26.3 in Rockafellar [198]). For a minimal representation,  $A$  is strictly convex by Proposition 2, and from the proof of Theorem 2(a), it is also essentially smooth, so that the stated result follows.  $\square$

This result is analogous to Proposition 2, in that the conditions stated ensure the essential smoothness of the dual function  $A^*$  in a minimal representation. The boundary behavior of  $\nabla A^*$  can be verified explicitly for the examples shown in Table 3.2, for which we have closed form expressions for  $A^*$ . For instance, in the Bernoulli case, we have  $\mathcal{M} = [0, 1]$  and  $|\nabla A^*(\mu)| = |\log[(1 - \mu)/\mu]|$ , which tends to infinity as  $\mu \rightarrow 0^+$  or  $\mu \rightarrow 1^-$ . Similarly, in the Poisson case, we have  $\mathcal{M} = (0, +\infty)$  and  $|\nabla A^*(\mu)| = |\log \mu|$ , which tends to infinity as  $\mu$  tends to the boundary point 0.

#### B.4 Proof of Theorem 3(b)

For any MRF defined by a tree  $T$ , both components of the Bethe variational problem (4.16) are exact:

- (a) by Proposition 4, the set  $\mathbb{L}(T)$  is equivalent to the marginal polytope  $\mathbb{M}(T)$ , and
- (b) the Bethe entropy  $H_{Bethe}(\cdot)$  is equivalent to the negative dual function  $-A^*(\cdot)$ .

Consequently, for the tree-structured problem, the Bethe variational principle is equivalent to a special case of the conjugate duality between

$(A, A^*)$  from Theorem 2; consequently, the value of the optimized Bethe problem is equal to the cumulant function  $A(\theta)$  as claimed. Finally, Theorem 2(c) implies that the optimum  $\tau^*$  corresponds to the exact marginal distributions of the tree-structured MRF. Strict convexity implies that the solution is unique.

### B.5 Proof of Theorem 5

We begin by proving assertion (8.4a). Let  $\mathcal{P}$  be the space of all densities  $p$ , taken with respect to the base measure  $\nu$  associated with the given exponential family. On the one hand, for any  $p \in \mathcal{P}$ , we have  $\int \langle \theta, \phi(x) \rangle p(x) \nu(dx) \leq \max_{x \in \mathcal{X}^m} \langle \theta, \phi(x) \rangle$ , whence

$$\sup_{p \in \mathcal{P}} \int \langle \theta, \phi(x) \rangle p(x) \nu(dx) \leq \max_{x \in \mathcal{X}^m} \langle \theta, \phi(x) \rangle. \quad (\text{B.2})$$

Since the support of  $\nu$  is  $\mathcal{X}^m$ , equality is achieved in the inequality (B.2) by taking a sequence  $p^n$  converging to a delta function  $\delta_{x^*}(x)$ , where  $x^* \in \arg \max_x \langle \theta, \phi(x) \rangle$ . Finally, by linearity of expectation and the definition of  $\mathcal{M}$ , we have  $\sup_{p \in \mathcal{P}} \int \langle \theta, \phi(x) \rangle p(x) \nu(dx) = \sup_{\mu \in \mathcal{M}} \langle \theta, \mu \rangle$ , which establishes the claim (8.4a).

We now turn to the claim (8.4b). By Proposition 2, the function  $A$  is lower semi-continuous. Therefore, for all  $\theta \in \Omega$ , the quantity  $\lim_{\beta \rightarrow +\infty} A(\beta\theta)/\beta$  is equivalent to the recession function of  $A$ , which we denote by  $A_\infty$  (see Corollary 8.5.2 of Rockafellar [198]). Hence, it suffices to establish that  $A_\infty(\theta)$  is equal to  $\sup_{\mu \in \mathcal{M}} \langle \theta, \mu \rangle$ . Using the lower semi-continuity of  $A$  and Theorem 13.3 of Rockafellar [198], the recession function of  $A$  corresponds to the support function of the effective domain of its dual. By Theorem 2, we have  $\overline{\text{dom } A^*} = \overline{\mathcal{M}}$ , whence  $A_\infty(\theta) = \sup_{\mu \in \overline{\mathcal{M}}} \langle \theta, \mu \rangle$ . Finally, the supremum is not affected by taking the closure.

# C

---

## Variational principles for multivariate Gaussians

---

In this appendix, we describe how the general variational representation (3.45) specializes to the multivariate Gaussian case. For a Gaussian Markov random field, as described in Example 9, computing the mapping  $(\theta, \Theta) \mapsto (\mu, \Sigma)$  amounts to computing the mean vector  $\mu \in \mathbb{R}^d$ , and the covariance matrix  $\Sigma - \mu\mu^T$ . These computations are typically carried out through the so-called normal equations [123]. In this section, we derive such normal equations for Gaussian inference from the principle (3.45).

### C.1 Gaussian with known covariance

We first consider the case of  $X = (X_1, \dots, X_m)$  with unknown mean  $\mu \in \mathbb{R}^m$ , but *known covariance* matrix  $\Lambda$ , which we assume to be strictly positive definite (hence invertible). Under this assumption, we may alternatively use the precision matrix  $P$ , corresponding to the inverse covariance  $\Lambda^{-1}$ . This collection of models can be written as a  $m$ -dimensional exponential family with respect to the base measure  $\nu(dx) = \exp(-\frac{1}{2}x^T Px)dx$ , with densities of the form  $p_\theta(x) = \exp\{\sum_{i=1}^m \theta_i x_i - A(\theta)\}$ . With this notation, the well-known *normal*



equations for computing the Gaussian mean vector  $\mu = \mathbb{E}[X]$  are given by  $\mu = P^{-1}\theta$ .

Let us now see how these normal equations are a consequence of the variational principle (3.45). We first compute the cumulant function for the specified exponential family as follows

$$\begin{aligned} A(\theta) &= \log \int_{\mathbb{R}^m} \exp\left\{\sum_{i=1}^m \theta_i x_i\right\} \exp\left(-\frac{1}{2}x^T P x\right) dx \\ &= \frac{1}{2}\theta^T P^{-1}\theta \\ &= \frac{1}{2}\theta^T \Lambda \theta, \end{aligned}$$

after some algebra.<sup>1</sup> Consequently, we have  $A(\theta) < +\infty$  for all  $\theta \in \mathbb{R}^m$ , meaning that  $\Omega = \mathbb{R}^m$ . Moreover, the mean parameter  $\mu = \mathbb{E}[X]$  takes values in all of  $\mathbb{R}^m$ , so that  $\mathcal{M} = \mathbb{R}^m$ . Next, a straightforward calculation yields that the conjugate dual of  $A$  is given by  $A^*(\mu) = \frac{1}{2}\mu^T P \mu$ . Consequently, the variational principle (3.45) takes the form

$$A(\theta) = \sup_{\mu \in \mathbb{R}^m} \left\{ \langle \theta, \mu \rangle - \frac{1}{2}\mu^T P \mu \right\},$$

which is an unconstrained quadratic program. Taking derivatives to find the optimum yields that  $\mu(\theta) = P^{-1}\theta$ ; thus, the normal equations for Gaussian inference are a consequence of the variational principle (3.45) specialized to this particular exponential family. Moreover, by the Hammersley-Clifford theorem [103, 23, 99], the precision matrix  $P$  has the same sparsity pattern as the graph adjacency matrix (i.e., for all  $i \neq j$ ,  $P_{ij} \neq 0$  implies that  $(i, j) \in E$ ). Consequently, the matrix-inverse-vector problem  $P^{-1}\theta$  can often be solved very quickly, by specialized methods that exploit the graph structure of  $P$  (see, for instance, Golub and van Loan [96]). Perhaps the best-known example is the case of a tri-diagonal matrix  $P$ , in which case the original Gaussian vector  $X$  is Markov with respect to a chain-structured graph, with edges  $(i, i + 1)$  for  $i = 1, \dots, m - 1$ . More generally, the Kalman filter

<sup>1</sup>Note that since  $[\nabla^2 A(\theta)]_{ij} = \text{cov}(X_i, X_j)$ , this calculation is consistent with Proposition 2 (see equation (3.41b)).

on trees [257] can be viewed as a fast algorithm for solving the matrix-inverse-vector system of normal equations. For graphs without cycles, there are also local iterative methods for solving the normal equations. In particular, as we discuss in Example 28 in Section 5, Gaussian mean field theory leads to the Gauss-Jacobi updates for solving the normal equations.

### C.2 Gaussian with unknown covariance

We now turn to the more general case of a Gaussian with unknown covariance (see Example 5). Its mean parameterization is specified by the vector  $\mu = \mathbb{E}[X]$  and the second-order moment matrix  $\Sigma = \mathbb{E}[XX^T]$ . The set  $\overline{\mathcal{M}}$  of valid mean parameters is completely characterized by the positive semidefiniteness (PSD) constraint  $\Sigma - \mu\mu^T \succeq 0$  (see Example 9). Turning to the form of the dual function, given a set of valid mean parameters  $(\mu, \Sigma) \in \mathcal{M}^\circ$ , we can associate them with a Gaussian random vector  $X$  with mean  $\mu$  and covariance  $\Sigma - \mu\mu^T \succ 0$ . The entropy of such a multivariate Gaussian takes the form [54]

$$h(X) = -A^*(\mu, \Sigma) = \frac{1}{2} \log \det [\Sigma - \mu\mu^T] + \frac{m}{2} \log 2\pi e.$$

By the Schur complement formula [111], we may rewrite the negative dual function as

$$-A^*(\mu, \Sigma) = \frac{1}{2} \log \det \begin{bmatrix} 1 & \mu^T \\ \mu & \Sigma \end{bmatrix} + \frac{m}{2} \log 2\pi e. \quad (\text{C.1})$$

Since the log-determinant function is strictly concave on the cone of positive semidefinite matrices [36], this representation demonstrates the convexity of  $-(A^*)$  in an explicit way.

These two ingredients allow us to write the variational principle (3.45) as specialized to the multivariate Gaussian with unknown covariance as follows:

$$A(\theta, \Theta) = \sup_{(\mu, \Sigma), \Sigma - \mu\mu^T \succ 0} \left\{ \langle \theta, \mu \rangle + \frac{1}{2} \langle \langle \Theta, \Sigma \rangle \rangle + \frac{1}{2} \log \det \begin{bmatrix} 1 & \mu^T \\ \mu & \Sigma \end{bmatrix} + \frac{1}{2} \log 2\pi e \right\}. \quad (\text{C.2})$$

This problem is an instance of a well-known class of concave maximization problems, known (for obvious reasons) as a log-determinant problem [230]. In general, log-determinant problems can be solved efficiently using standard methods for convex programming, such as interior point methods. However, this specific log-determinant problem (C.2) actually has a closed-form solution: in particular, it can be verified by a Lagrangian reformulation that the optimal pair  $(\mu^*, \Sigma^*)$  is specified by the relations

$$\Sigma^* - \mu^*(\mu^*)^T = -[\Theta]^{-1}, \quad \text{and} \quad \mu^* = -[\Theta]^{-1} \theta. \quad (\text{C.3})$$

Note that these relations, which emerge as the optimality conditions for the log-determinant problem (C.2), are actually familiar statements. Recall from our exponential representation of the Gaussian density (3.12) that  $-\Theta$  is the precision matrix. Thus, the first equality in equation (C.3) confirms that the covariance matrix is the inverse of the precision matrix, whereas the second equality corresponds to the normal equations for the mean  $\mu$  of a Gaussian. Thus, as a special case of the general variational principle (3.45), we have re-derived the familiar equations for Gaussian inference.

### C.3 Gaussian mode computation

In Section C.2, we have seen that for a multivariate Gaussian with unknown covariance, the general variational principle (3.45) corresponds to a log-determinant optimization problem. In this section, we show that the mode-finding problem for a Gaussian is a semidefinite program (SDP), another well-known class of convex optimization problems [229, 230]. Consistent with our development in Section 8.1, this SDP arises as the zero-temperature limit of the log-determinant problem. Recall that for the multivariate Gaussian, the mean parameters are the mean vector  $\mu = \mathbb{E}[X]$  and second-moment matrix  $\Sigma = \mathbb{E}[XX^T]$ , and the set  $\overline{\mathcal{M}}$  of valid mean parameters is characterized by the semidefinite constraint  $\Sigma - \mu\mu^T \succeq 0$ , or equivalently (via the Schur complement formula [111]) by a linear matrix inequality

$$\Sigma - \mu\mu^T \succeq 0 \iff \begin{bmatrix} 1 & \mu^T \\ \mu & \Sigma \end{bmatrix} \succeq 0. \quad (\text{C.4})$$

Consequently, by applying Theorem 5, we conclude that

$$\max_{x \in \mathcal{X}} \left[ \langle \theta, x \rangle + \frac{1}{2} \langle \langle \Theta, xx^T \rangle \rangle \right] = \max_{(\mu, \Sigma), \Sigma - \mu\mu^T \succeq 0} \left[ \langle \theta, \mu \rangle + \frac{1}{2} \langle \langle \Theta, \Sigma \rangle \rangle \right].$$

The problem on the left-hand side is simply a quadratic program, with optimal solution corresponding to the mode  $x^* = -[\Theta]^{-1}\theta$  of the Gaussian density. The problem on the right-hand side is a semidefinite program [229], as it involves a linear objective function subject to the linear matrix inequality (LMI) constraint (C.4).

In order to demonstrate how the optimum of this semidefinite program (SDP) recovers the Gaussian mode  $x^*$ , we begin with the fact [111] that a matrix  $B$  is positive semidefinite if and only if  $\langle \langle B, C \rangle \rangle := \text{trace}(BC) \geq 0$  for all other positive semidefinite matrices  $C$ . (This property corresponds to the self-duality of the cone  $\mathcal{S}_+^m$ ). Applying this fact to the matrices  $B = \Sigma - \mu\mu^T \succeq 0$  and  $C = -\Theta \succ 0$  yields the inequality  $\langle \langle \Theta, \Sigma \rangle \rangle \leq \langle \langle \Theta, \mu\mu^T \rangle \rangle$ , which in turn implies that

$$\langle \theta, \mu \rangle + \frac{1}{2} \langle \langle \Theta, \Sigma \rangle \rangle \leq \langle \theta, \mu \rangle + \frac{1}{2} \langle \langle \Theta, \mu\mu^T \rangle \rangle. \quad (\text{C.5})$$

Observe that the right-hand side is simply a quadratic program in  $\mu \in \mathbb{R}^m$ , with its maximum attained at  $\mu^* = -[\Theta]^{-1}\theta$ . Consequently, if we maximize the left-hand side over the set  $\Sigma - \mu\mu^T \succeq 0$ , the bound will be achieved at  $\mu^*$  and  $\Sigma^* = \mu^*(\mu^*)^T$ . The interpretation is the optimal solution corresponds to a degenerate Gaussian, specified by mean parameters  $(\mu^*, \Sigma^*)$ , with a zero covariance matrix, so that all mass concentrates on its mean  $\mu^* = x^*$ .

In practice, of course, one would not compute the mode of a Gaussian problem by solving this semidefinite program, since it can be computed by more direct methods, including Kalman filtering [123], corresponding to the sum-product algorithm on trees, by numerical methods such as conjugate gradient [64], by the max-product/min-sum algorithm (see Section 8.3), both of which solve the quadratic program in an iterative manner, or by iterative algorithms based on tractable subgraphs [218, 158]. However, the SDP-based formulation (C.5) provides valuable perspective on the use of semidefinite relaxations for integer programming problems, as discussed in Section 8.5.

# D

---

## Clustering and augmented hypergraphs

---

In this appendix, we elaborate on various techniques to transform hypergraphs so as to apply Kikuchi and related cluster variational methods. The most natural strategy is to develop techniques for approximate inference based directly on the structure of the original hypergraph. The Bethe approximation is of this form, corresponding to the case when the original structure is an ordinary graph. Alternatively, it can be beneficial to build approximations based on an *augmented hypergraphs*  $G = (V, E)$  built on top of the original hypergraph. A natural way in which to construct such augmented hypergraphs is by clustering nodes so as to define new hyperedges; a variety of different techniques of this nature have been discussed in the literature [e.g., 263, 264, 184, 163].

### D.1 Covering augmented hypergraphs

For the purposes of this discussion, we focus on a subclass of augmented hypergraphs. We begin by requiring that the original hypergraph  $G'$  is *covered* by the augmented hypergraph, meaning that the hyperedge set  $E$  of the augmented hypergraph includes all hyperedges in  $E'$  (as well as the vertices of  $G'$ ). A desirable feature of this requirement is that

any Markov random field defined by  $G'$  can also be viewed as an MRF on a covering hypergraph  $G$ , simply by setting  $\theta_h = 0$  for all  $h \in E \setminus E'$ .

**Example 46 (Covering hypergraph).** To illustrate, suppose that the original hypergraph  $G'$  is simply an ordinary graph—namely, the  $3 \times 3$  grid shown in Figure D.1(a). As illustrated in panel (b), we cluster the nodes into groups of four, which is known as Kikuchi 4-plaque clustering in statistical physics [263, 264]. We then form the augmented hypergraph  $G$  shown in panel (c), with hyperedge set  $E := E' \cup \{(1245), (2356), (4578), (5689)\}$ . The darkness of the boxes in this diagram reflects the depth of the hyperedges in the poset diagram. This hypergraph covers the original (hyper)graph, since it includes as hyperedges all edges and vertices of the original  $3 \times 3$  grid. ♣

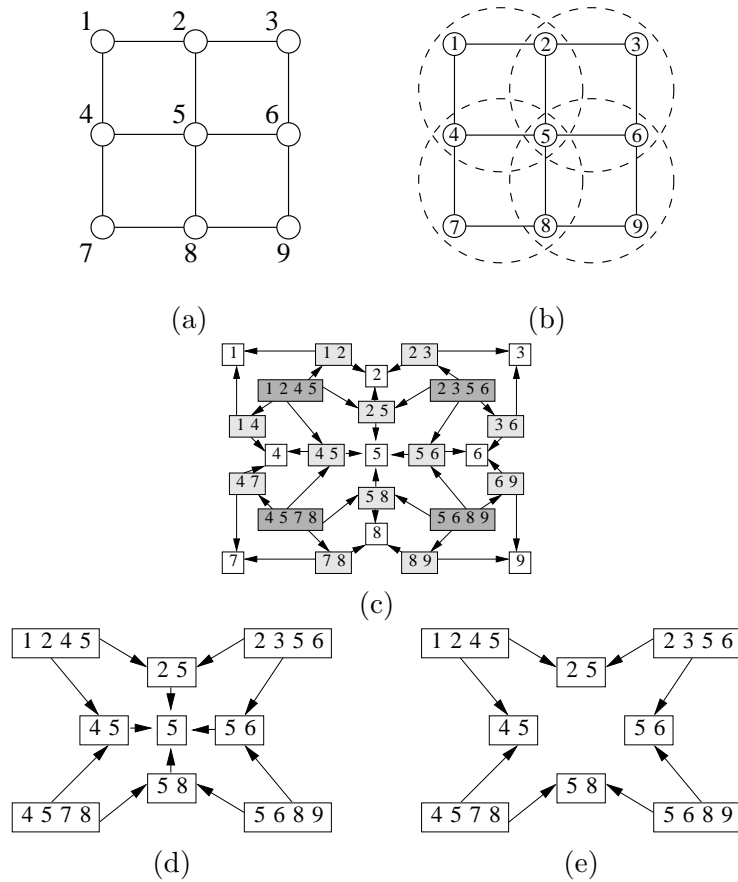
As emphasized by Yedidia et al. [264], it turns out to be important to ensure that every hyperedge (including vertices) in the original hypergraph  $G'$  is counted exactly once in the augmented hypergraph  $G$ . More specifically, for a given hyperedge  $h' \in E'$ , consider the set  $\mathcal{C}(h') := \{f \in E \mid f \supseteq h'\}$  of hyperedges in  $E$  that contain  $h'$ . For ease of reference, we restate the definition (4.44) of the overcounting numbers  $c(\cdot)$  associated with the hypergraph  $G$ . In particular, these overcounting numbers are defined in terms of the Möbius function associated with  $G$ , viewed as a poset, in the following way:

$$c(f) := \sum_{e \in \mathcal{A}^+(f)} \omega(f, e). \tag{D.1}$$

The *single counting criterion* requires that for all  $h' \in E'$  (including all single vertices), there holds

$$\sum_{f \in \mathcal{C}(h')} c(f) = 1. \tag{D.2}$$

**Example 47 (Single counting).** To illustrate the single counting criterion, we consider two additional hypergraphs that can be constructed from the  $3 \times 3$  grid of Figure D.1(a). The vertex set and edge set of the grid form the original hyperedge set  $E'$ . The hypergraph in panel (d) is constructed by the Kikuchi method described by Yedidia



**Fig. D.1.** Constructing new hypergraphs via clustering and the single counting criterion. (a) Original (hyper)graph  $G'$  is a  $3 \times 3$  grid. Its hyperedge set  $E'$  consists of the union of the vertex set with the (ordinary) edge set. (b) Nodes are clustered into groups of four. (c) A covering hypergraph  $G$  formed by adjoining these 4-clusters to the original hyperedge set  $E'$ . Darkness of the boxes indicates depth of the hyperedges in the poset representation. (d) An augmented hypergraph constructed by the Kikuchi method. (e) A third augmented hypergraph that fails the single counting criterion for (5).

et al. [263, 264]. In this construction, we include the four clusters, all of their pairwise intersections, and all the intersections of intersections (only (5) in this case). The hypergraph in panel (e) includes only the

hyperedges of size four and two; that is, it omits the hyperedge (5).

Let us focus first on the hypergraph (e), and understand why it violates the single counting criterion for hyperedge (5). Viewed as a poset, all of the maximal hyperedges (of size four) in this hypergraph have a counting number of  $c(h) = \omega(h, h) = 1$ . Any hyperedge  $f$  of size two has two parents, each with an overcounting number of 1, so that  $c(f) = 1 - (1 + 1) = -1$ . The hyperedge (5) is a member of the original hyperedge set  $E'$  (of the  $3 \times 3$  grid), but not of the augmented hypergraph. It is included in all the hyperedges, so that  $\mathcal{C}(5) = E$  and  $\sum_{h \in \mathcal{C}(5)} c(h) = 0$ . Thus, the single criterion condition fails to hold for hypergraph (e). In contrast, it can be verified that for the hypergraphs in panels (c) and (d), the single counting condition holds for all hyperedges  $h' \in E'$ .

There is another interesting fact about hypergraphs (c) and (d). If we eliminate from hypergraph (c) all hyperedges that have zero overcounting numbers, the result is hypergraph (d). To understand this reduction, consider for instance the hyperedge (14) which appears in (c) but not in (d). Since it has only one parent (which is a maximal hyperedge), we have  $c(14) = 0$ . In a similar fashion, we see that  $c(12) = 0$ . These two equalities together imply that  $c(1) = 0$ , so that we can eliminate hyperedges (12), (14) and (1) from hypergraph (c). By applying a similar argument to the remaining hyperedges, we can fully reduce hypergraph (c) to hypergraph (d). ♣

It turns out that if the augmented hypergraph  $G$  covers the original hypergraph  $G'$ , then the single counting criterion is always satisfied. Implicit in this definition of covering is that the hyperedge set  $E'$  of the original hypergraph includes the vertex set, so that equation (D.2) should hold for the vertices. The proof is quite straightforward: we begin by observing that under the covering condition, the set  $\mathcal{C}(h)$  is equal to  $\mathcal{A}^+(h)$  in the augmented hypergraph  $G$ . We then invoke the following result:

**Lemma 3** (Single counting). *For any  $h \in E$ , the associated overcounting numbers satisfy the identity  $\sum_{e \in \mathcal{A}^+(h)} c(e) = 1$ , which can be written equivalently as  $c(h) = 1 - \sum_{e \in \mathcal{A}(h)} c(e)$ .*



*Proof.* From the definition of  $c(h)$ , we have the identity:

$$\sum_{h \in \mathcal{A}^+(g)} c(h) = \sum_{h \in \mathcal{A}^+(g)} \sum_{f \in \mathcal{A}^+(h)} \omega(h, f). \quad (\text{D.3})$$

Considering the double sum on the right-hand side, we see that for a fixed  $d \in \mathcal{A}^+(g)$ , there is a term  $\omega(d, e)$  for each  $e$  such that  $g \subseteq e \subseteq d$ . Using this observation, we can write

$$\begin{aligned} \sum_{h \in \mathcal{A}^+(g)} \sum_{f \in \mathcal{A}^+(h)} \omega(h, f) &= \sum_{d \in \mathcal{A}^+(g)} \sum_{\{e | g \subseteq e \subseteq d\}} \omega(e, d) \\ &\stackrel{(a)}{=} \sum_{d \in \mathcal{A}^+(g)} \delta(d, g) \\ &\stackrel{(b)}{=} 1. \end{aligned}$$

Here equality (a) follows from the definition of the Möbius function (see Appendix E.1), and  $\delta(d, g)$  is the Kronecker delta function, from which equality (b) follows.  $\square$

Thus, the construction that we have described, in which the hyperedges (including all vertices) of the original hypergraph  $G'$  are covered by  $G$  and the partial ordering is set inclusion, ensures that the single counting criterion is always satisfied. We emphasize that there is a broad range of other possible constructions [e.g., 263, 264, 184, 163].

## D.2 Specification of compatibility functions

Our final task is to specify how to assign the compatibility functions associated with the original hypergraph  $G' = (V, E')$  with the hyperedges of the augmented hypergraph  $G = (V, E)$ . It is convenient to use the notation  $\psi'_g(x_g) := \exp\{\theta_g(x_g)\}$  for the compatibility functions of the original hypergraph, corresponding to terms in the product (4.45). We can extend this definition to all hyperedges in  $E$  by setting  $\psi'_h(x_h) \equiv 1$  for any hyperedge  $h \in E \setminus E'$ . For each hyperedge  $h \in E$ , we then define a new compatibility function  $\psi_h$  as follows:

$$\psi_h(x_h) := \psi'_h(x_h) \prod_{g \in \mathcal{S}(h)} \psi'_g(x_g), \quad (\text{D.4})$$

where  $\mathcal{S}(h) := \{g \in E' \setminus E \mid g \subset h\}$  is the set of hyperedges in  $E' \setminus E$  that are subsets of  $h$ . To illustrate this definition, consider the Kikuchi construction of Figure D.1(d), which is an augmented hypergraph for the  $3 \times 3$  grid in Figure D.1(a). For the hyperedge (25), we have  $\mathcal{S}(25) = \{(2)\}$ , so that  $\psi_{25} = \psi'_{25}\psi'_2$ . On the other hand, for the hyperedge (1245), we have  $\psi'_{1245} \equiv 1$  (since (1245) appears in  $E$  but not in  $E'$ ), and  $\mathcal{S}(1245) = \{(1), (12), (14)\}$ . Accordingly, equation (D.4) yields  $\psi_{1245} = \psi'_1\psi'_{12}\psi'_{14}$ . More generally, using the definition (D.4), it is straightforward to verify that the equivalence  $\prod_{h \in E} \psi_h(x_h) = \prod_{g \in E'} \psi'_g(x_g)$  holds, so that we have preserved the structure of the original MRF.

# E

---

## Miscellaneous results

---

This appendix collects together various miscellaneous results on graphical models and posets.

### E.1 Möbius inversion

This appendix provides a brief overview of the Möbius function associated with a partially-ordered set (poset); see Stanley [216] for a thorough treatment. The *zeta function*  $\zeta(g, h)$  of a poset is defined as:

$$\zeta(g, h) = \begin{cases} 1 & \text{if } g \subseteq h \\ 0 & \text{otherwise} \end{cases} \quad (\text{E.1})$$

The Möbius function  $\omega$  arises as the multiplicative inverse of this zeta function. It is defined in a recursive fashion, by first specifying  $\omega(g, g) = 1$  for all  $g$ . Once  $\omega(g, f)$  has been defined for all  $f$  such that  $g \subseteq f \subset h$ , we then define:

$$\omega(g, h) = - \sum_{\{f \mid g \subseteq f \subset h\}} \omega(g, f). \quad (\text{E.2})$$

With this definition, it can be seen that  $\omega$  and  $\zeta$  are multiplicative inverses, in the sense that

$$\sum_{\{f \mid g \subseteq f \subseteq h\}} \omega(g, f) \zeta(f, h) = \delta(g, h),$$

where  $\delta(g, h)$  is the Kronecker delta.

**Lemma 4** (Möbius inversion formula). *Let  $\Upsilon(h)$  be a real-valued function defined for  $h$  in a poset. Define a new real-valued function  $\Omega$  via:*

$$\Omega(h) = \sum_{g \in \mathcal{D}^+(h)} \Upsilon(g), \quad (\text{E.3})$$

where  $\mathcal{D}^+(h) := \{g \mid g \subseteq h\}$  is the set of descendants of  $h$ . Then we have the relation

$$\Upsilon(h) = \sum_{g \in \mathcal{D}^+(h)} \omega(g, h) \Omega(g), \quad (\text{E.4})$$

where  $\omega$  is the associated Möbius function.

## E.2 Conversion to a pairwise Markov random field

In this appendix, we describe how any Markov random field with discrete random variables can be converted to an equivalent pairwise form (i.e., with interactions only between pairs of variables). To illustrate the general principle, it suffices to show how to convert a compatibility function  $\psi_{123}$  defined on a triplet  $\{x_1, x_2, x_3\}$  of random variables into a pairwise form. To do so, we introduce an auxiliary node  $A$ , and associate with it random variable  $z$  that takes values in the Cartesian product space  $\mathcal{X}_1 \times \mathcal{X}_2 \times \mathcal{X}_3$ . In this way, each configuration of  $z$  can be identified with a triplet  $(z_1, z_2, z_3)$ . For each  $s \in \{1, 2, 3\}$ , we define a pairwise compatibility function  $\psi_{As}$ , corresponding to the interaction between  $z$  and  $x_s$ , by  $\psi_{As}(z, x_s) := [\psi_{123}(z_1, z_2, z_3)]^{1/3} \mathbb{I}[z_s = x_s]$ . (The purpose of the  $1/3$  power is to incorporate  $\psi_{123}$  with the correct exponent.) With this definition, it is straightforward to verify that the equivalence

$$\psi_{123}(x_1, x_2, x_3) = \sum_z \prod_{s=1}^3 \psi_{As}(z, x_s)$$

holds, so that our augmented model faithfully captures the interaction among the triplet  $\{x_1, x_2, x_3\}$ .

## References

---

- [1] A. Agresti. *Categorical Data Analysis*. Wiley, New York, 2002.
- [2] S.M. Aji and R.J. McEliece. The generalized distributive law. *IEEE Transactions on Information Theory*, 46:325–343, 2000.
- [3] N. I. Akhiezer. *The Classical Moment Problem and Some Related Questions in Analysis*. Hafner Publishing Company, New York, 1966.
- [4] S. Amari. Differential geometry of curved exponential families—curvatures and information loss. *Annals of Statistics*, 10(2):357–385, 1982.
- [5] S. Amari and H. Nagaoka. *Methods of information geometry*. AMS, Providence, RI, 2000.
- [6] G. An. A note on the cluster variation method. *Journal of Statistical Physics*, 52(3):727–734, 1988.
- [7] O. Banerjee, L. El Ghaoui, and A. d’Aspremont. Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data. *Journal of Machine Learning Research*, 9:485–516, 2008.
- [8] F. Barahona and M. Groetschel. On the cycle polytope of a binary matroid. *J. Combin. Theory, Series B*, B 40:40–62, 1986.
- [9] D. Barber and W. Wiegerinck. Tractable variational structures for approximating graphical models. In *NIPS 11*. MIT Press, 1999.
- [10] O. E. Barndorff-Nielsen. *Information and Exponential Families*. Wiley, Chichester, UK, 1978.
- [11] R. J. Baxter. *Exactly Solved Models in Statistical Mechanics*. Academic Press, New York, 1982.
- [12] M. Bayati, C. Borgs, J. Chayes, and R. Zecchina. Belief-propagation for weighted b-matchings on arbitrary graphs and its relation to linear programs

- with integer solutions. Technical Report arxiv:0709.1190, Microsoft Research, September 2007.
- [13] M. Bayati, D. Shah, and M. Sharma. Maximum weight matching for max-product belief propagation. In *International Symposium on Information Theory*, Adelaide, Australia, 2005.
- [14] M. J. Beal. *Variational algorithms for approximate Bayesian inference*. PhD thesis, Gatsby Computational Neuroscience Unit, University College, London, 2003.
- [15] C. Berge. *The Theory of Graphs and its Applications*. Wiley, New York, 1964.
- [16] E. Berlekamp, R. McEliece, and H. van Tilborg. On the inherent intractability of certain coding problems. *IEEE Transactions on Information Theory*, 24:384–386, 1978.
- [17] U. Bertele and F. Brioschi. *Nonserial Dynamic Programming*. Academic Press, New York, 1972.
- [18] D. P. Bertsekas. *Dynamic Programming and Stochastic Control*, volume 1. Athena Scientific, Belmont, MA, 1995.
- [19] D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific, Belmont, MA, 1995.
- [20] D. P. Bertsekas. *Network Optimization: Continuous and Discrete Methods*. Athena Scientific, Belmont, MA, 1998.
- [21] D. P. Bertsekas. *Convex Analysis and Optimization*. Athena Scientific, Belmont, MA, 2003.
- [22] D. Bertsimas and J. N. Tsitsiklis. *Introduction to Linear Optimization*. Athena Scientific, Belmont, MA, 1997.
- [23] J. Besag. Spatial interaction and the statistical analysis of lattice systems. *Journal of the Royal Statistical Society, Series B*, 36:192–236, 1974.
- [24] J. Besag. Statistical analysis of non-lattice data. *The Statistician*, 24(3):179–195, 1975.
- [25] J. Besag. On the statistical analysis of dirty pictures. *Journal of the Royal Statistical Society, Series B*, 48(3):259–279, 1986.
- [26] J. Besag and P. J. Green. Spatial statistics and Bayesian computation. *Journal of the Royal Statistical Society, Series B*, 55(1):25–37, 1993.
- [27] H. A. Bethe. Statistics theory of superlattices. *Proc. Royal Soc. London, Series A*, 150(871):552–575, 1935.
- [28] D. M. Blei and M. I. Jordan. Variational inference for Dirichlet process mixtures. *Bayesian Analysis*, 1:121–144, 2005.
- [29] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- [30] H. Bodlaender. A tourist guide through treewidth. *Acta Cybernetica*, 11:1–21, 1993.
- [31] B. Bollobás. *Graph Theory: An Introductory Course*. Springer-Verlag, New York, 1979.
- [32] B. Bollobás. *Modern Graph Theory*. Springer-Verlag, New York, 1998.
- [33] E. Boros, Y. Crama, and P. L. Hammer. Upper bounds for quadratic 0-1 maximization. *Operations Research Letters*, 9:73–79, 1990.

- [34] E. Boros and P. L. Hammer. Pseudo-boolean optimization. *Discrete Applied Mathematics*, 123:155–225, 2002.
- [35] J. Borwein and A. Lewis. *Convex Analysis*. Springer-Verlag, New York, 1999.
- [36] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, UK, 2004.
- [37] X. Boyen and D. Koller. Tractable inference for complex stochastic processes. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, pages 33–42. Morgan Kaufmann, 1998.
- [38] A. Braunstein, M. Mézard, and R. Zecchina. Survey propagation: an algorithm for satisfiability. Technical report, 2003. arXiv:cs.CC/02122002 v2.
- [39] L. M. Bregman. The relaxation method for finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7:191–204, 1967.
- [40] L.D. Brown. *Fundamentals of Statistical Exponential Families*. Institute of Mathematical Statistics, Hayward, CA, 1986.
- [41] C. B. Burge and S. Karlin. Finding the genes in genomic "dna". *Current Opinion in Structural Biology*, 8:346–354, 1998.
- [42] Y. Censor and S. A. Zenios. *Parallel Optimization: Theory, Algorithms, and Applications*. Oxford University Press, Oxford, 1988.
- [43] M. Cetin, L. Chen, J. W. Fisher, A. T. Ihler, R. L. Moses, M. J. Wainwright, and A. S. Willsky. Distributed fusion in sensor networks. *IEEE Signal Processing Magazine*, 23:42–55, 2006.
- [44] D. Chandler. *Introduction to Modern Statistical Mechanics*. Oxford University Press, Oxford, 1987.
- [45] C. Chekuri, S. Khanna, J. Naor, and L. Zosin. A linear programming formulation and approximation algorithms for the metric labeling problem. *SIAM Journal on Discrete Mathematics*, 18(3):608–625, 2005.
- [46] L. Chen, M. J. Wainwright, M. Cetin, and A. Willsky. Multitarget-multisensor data association using the tree-reweighted max-product algorithm. In *SPIE Aerosense Conference*, Orlando, FL, April 2003.
- [47] M. Chertkov and V. Y. Chernyak. Loop calculus helps to improve belief propagation and linear programming decoding of LDPC codes. In *Proceedings of the Allerton Conference on Control, Communication and Computing*, Monticello, IL, September 2006.
- [48] M. Chertkov and V. Y. Chernyak. Loop series for discrete statistical models on graphs. *Journal of Statistics Mechanics*, page P06009, June 2006.
- [49] M. Chertkov and V. Y. Chernyak. Loop calculus helps to improve belief propagation and linear programming decodings of low density parity check codes. In *Proceedings of the Allerton Conference on Control, Communication and Computing*, 2007.
- [50] M. Chertkov and M. G. Stepanov. An efficient pseudo-codeword search algorithm for linear programming decoding of LDPC codes. Technical Report arxiv:cs.IT/0601113, Los Alamos National Laboratories, September 2006.
- [51] S. Chopra. On the spanning tree polyhedron. *Operations Research Letters*, 8:25–29, 1989.



- [52] S. Cook. The complexity of theorem-proving procedures. In *Proceedings of the Third Annual ACM Symposium on Theory of Computing*, pages 151–158, 1971.
- [53] T. H. Cormen, C. E. Leiserson, and R. L. Rivest. *Introduction to Algorithms*. MIT Press, Cambridge, MA, 1990.
- [54] T.M. Cover and J.A. Thomas. *Elements of Information Theory*. John Wiley and Sons, New York, 1991.
- [55] G. Cross and A. Jain. Markov random field texture models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 5:25–39, 1983.
- [56] I. Csiszár and G. Tusnády. Information geometry and alternating minimization procedures. *Statistics and Decisions, Supplemental Issue 1*, pages 205–237, 1984.
- [57] I. Csiszár. A geometric interpretation of Darroch and Ratcliff’s generalized iterative scaling. *Annals of Statistics*, 17(3):1409–1413, 1989.
- [58] J. N. Darroch and D. Ratcliff. Generalized iterative scaling for log-linear models. *Annals of Mathematical Statistics*, 43:1470–1480, 1972.
- [59] C. Daskalakis, A. G. Dimakis, R. M. Karp, and M. J. Wainwright. Probabilistic analysis of linear programming decoding. *IEEE Trans. Information Theory*, 54(8):3565–3578, 2008.
- [60] A. d’Aspremont, O. Banerjee, and L. El Ghaoui. First order methods for sparse covariance selection. *SIAM Journal on Matrix Analysis and its Applications*, 30(1):55–66, 2008.
- [61] J. Dauwels, H. A. Loeliger, P. Merkli, and M. Ostojic. On structured-summary propagation, LFSR synchronization, and low-complexity trellis decoding. In *Proceedings of the Allerton Conference on Control, Communication and Computing*, pages 459–467, Monticello, IL, 2003.
- [62] A. P. Dawid. Applications of a general propagation algorithm for probabilistic expert systems. *Statistics and Computing*, 2:25–36, 1992.
- [63] R. Dechter. *Constraint Processing*. Morgan Kaufmann, Palo Alto, CA, 2003.
- [64] J.W. Demmel. *Applied numerical linear algebra*. SIAM, Philadelphia, 1997.
- [65] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 39:1–38, 1977.
- [66] M. Deza and M. Laurent. *Geometry of Cuts and Metric Embeddings*. Springer-Verlag, New York, 1997.
- [67] A. G. Dimakis and M. J. Wainwright. Guessing facets: Improved LP decoding and polytope structure. In *International Symposium on Information Theory*, Seattle, Washington, 2006.
- [68] M. Dudik, S. J. Phillips, and R. E. Schapire. Maximum entropy density estimation with generalized regularization and an application to species distribution modeling. *Journal of Machine Learning Research*, 8:1217–1260, 2007.
- [69] R. Durbin, S. Eddy, A. Krogh, and G. Mitchison, editors. *Biological Sequence Analysis*. Cambridge University Press, Cambridge, 1998.
- [70] J. Edmonds. Matroids and the greedy algorithm. *Mathematical Programming*, 1:127–136, 1971.

- [71] B. Efron. The geometry of exponential families. *Annals of Statistics*, 6:362–376, 1978.
- [72] G. Elidan, I. McGraw, and D. Koller. Residual belief propagation: Informed scheduling for asynchronous message-passing. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*. Morgan Kaufmann, 2006.
- [73] J. Feldman, D. R. Karger, and M. J. Wainwright. Linear programming-based decoding of turbo-like codes and its relation to iterative approaches. In *Proceedings of the Allerton Conference on Control, Communication and Computing*, Monticello, IL, October 2002.
- [74] J. Feldman, T. Malkin, R. A. Servedio, C. Stein, and M. J. Wainwright. LP decoding corrects a constant fraction of errors. *IEEE Transactions on Information Theory*, 53(1):82–89, 2007.
- [75] J. Feldman, M. J. Wainwright, and D. R. Karger. Using linear programming to decode binary linear codes. *IEEE Transactions on Information Theory*, 51:954–972, 2005.
- [76] J. Felsenstein. Evolutionary trees from DNA sequences: a maximum likelihood approach. *Journal of Molecular Evolution*, 17:368–376, 1981.
- [77] S. E. Fienberg. Contingency tables and log-linear models: Basic results and new developments. *Journal of the American Statistical Association*, 95(450):643–647, 2000.
- [78] M. E. Fisher. On the dimer solution of planar Ising models. *Journal of Mathematical Physics*, 7:1776–1781, 1966.
- [79] G. D. Forney, Jr. The Viterbi algorithm. *Proceedings of the IEEE*, 61:268–277, 1973.
- [80] G. D. Forney, Jr., R. Koetter, F. R. Kschischang, and A. Reznick. On the effective weights of pseudocodewords for codes defined on graphs with cycles. In *Codes, Systems and Graphical Models*, pages 101–112. Springer, New York, 2001.
- [81] W. T. Freeman, E. C. Pasztor, and O. T. Carmichael. Learning low-level vision. *Intl. J. Computer Vision*, 40(1):25–47, 2000.
- [82] B. Frey, R. Koetter, and N. Petrovic. Very loopy belief propagation for unwrapping phase images. In *Advances in Neural Information Processing Systems*, volume 14, pages 737–743, Cambridge, MA, 2001. MIT Press.
- [83] B. J. Frey, R. Koetter, and A. Vardy. Signal-space characterization of iterative decoding. *IEEE Trans. Info. Theory*, 47:766–781, 2001.
- [84] R. G. Gallager. *Low-density parity check codes*. MIT Press, Cambridge, MA, 1963.
- [85] A. E. Gelfand and A. F. M. Smith. Sampling-based approaches to calculating marginal densities. *Journal of the American Statistical Association*, 85:398–409, 1990.
- [86] S. Geman and D. Geman. Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6:721–741, 1984.
- [87] H. O. Georgii. *Gibbs Measures and Phase Transitions*. De Gruyter, New York, 1988.

- [88] Z. Ghahramani and M. J. Beal. Propagation algorithms for variational Bayesian learning. In *Advances in Neural Information Processing Systems*, volume 13, pages 507–513, Cambridge, MA, 2001. MIT Press.
- [89] Z. Ghahramani and M. I. Jordan. Factorial hidden Markov models. *Machine Learning*, 29:245–273, 1997.
- [90] W. Gilks, S. Richardson, and D. Spiegelhalter, editors. *Markov Chain Monte Carlo in Practice*. Chapman and Hall, New York, 1996.
- [91] A. Globerson and T. Jaakkola. Approximate inference using planar graph decomposition. In *Advances in Neural Information Processing Systems*, volume 19, pages 473–480, Cambridge, MA, 2006. MIT Press.
- [92] A. Globerson and T. Jaakkola. Approximate inference using conditional entropy decompositions. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, Puerto Rico, 2007.
- [93] A. Globerson and T. Jaakkola. Convergent propagation algorithms via oriented trees. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*. AUAI Press, 2007.
- [94] A. Globerson and T. Jaakkola. Fixing max-product: Convergent message passing algorithms for MAP LP-relaxations. In *Advances in Neural Information Processing Systems*, volume 20, pages 553–560, Cambridge, MA, 2007. MIT Press.
- [95] M. X. Goemans and D. P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *Journal of the ACM*, 42:1115–1145, 1995.
- [96] G. Golub and C. Van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, 1996.
- [97] V. Gómez, J. M. Mooij, and H. J. Kappen. Truncating the loop series expansion for BP. *Journal of Machine Learning Research*, 8:1987–2016, 2007.
- [98] D. M. Greig, B. T. Porteous, and A. H. Seheuly. Exact maximum a posteriori estimation for binary images. *Journal of the Royal Statistical Society B*, 51:271–279, 1989.
- [99] G. R. Grimmett. A theorem about random fields. *Bulletin of the London Mathematical Society*, 5:81–84, 1973.
- [100] M. Groetschel and K. Truemper. Decomposition and optimization over cycles in binary matroids. *J. Combin. Theory, Series B*, B 46:306–337, 1989.
- [101] M. Grötschel, L. Lovász, and A. Schrijver. *Geometric Algorithms and Combinatorial Optimization*. Springer-Verlag, Berlin, 1993.
- [102] P.L. Hammer, P. Hansen, and B. Simeone. Roof duality, complementation, and persistency in quadratic 0-1 optimization. *Mathematical Programming*, 28:121–155, 1984.
- [103] J. M. Hammersley and P. Clifford. Markov fields on finite graphs and lattices. Unpublished, 1971.
- [104] M. Hassner and J. Sklansky. The use of Markov random fields as models of texture. *Computer Graphics and Image Processing*, 12:357–370, 1980.
- [105] T. Hazan and A. Shashua. Convergent message-passing algorithms for inference over general graphs with convex free energy. In *Uncertainty in Artificial Intelligence (UAI)*, volume 24, 2008.

- [106] T. Heskes. On the uniqueness of loopy belief propagation fixed points. *Neural Computation*, 16:2379–2413, 2004.
- [107] T. Heskes. Convexity arguments for efficient minimization of the Bethe and Kikuchi free energies. *Journal of Artificial Intelligence Research*, 26:153–190, 2006.
- [108] T. Heskes, K. Albers, and B. Kappen. Approximate inference and constrained optimization. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, pages 313–320, San Mateo, CA, 2003. Morgan Kaufmann.
- [109] J. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms*, volume 1. Springer-Verlag, New York, 1993.
- [110] J. Hiriart-Urruty and C. Lemaréchal. *Fundamentals of Convex Analysis*. Springer-Verlag, New York, 2001.
- [111] R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, Cambridge, UK, 1985.
- [112] J. Hu, H.-A. Loeliger, J. Dauwels, and F. Kschischang. A general computation rule for lossy summaries/messages with examples from equalization. In *Proceedings of the Allerton Conference on Control, Communication and Computing*, Monticello, IL, September 2005.
- [113] B. Huang and T. Jebara. Loopy belief propagation for bipartite maximum weight b-matching. In *AISTATS*, San Juan, Puerto Rico, March 2007.
- [114] X. Huang, A. Acero, H-W. Hon, and R. Reddy. *Spoken Language Processing*. Prentice Hall, New York, 2001.
- [115] A. Ihler, J. Fisher, and A. S. Willsky. Loopy belief propagation: Convergence and effects of message errors. *Journal of Machine Learning Research*, 6:905–936, 2005.
- [116] E. Ising. Beitrag zur theorie der ferromagnetismus. *Zeitschrift für Physik*, 31:253–258, 1925.
- [117] T. S. Jaakkola. Tutorial on variational approximation methods. In M. Opper and D. Saad, editors, *Advanced Mean Field Methods: Theory and Practice*, pages 129–160. MIT Press, Cambridge, MA, 2001.
- [118] T. S. Jaakkola and M. Jordan. Improving the mean field approximation via the use of mixture distributions. In M. Jordan, editor, *Learning in graphical models*, pages 105–161. MIT Press, 1999.
- [119] T. S. Jaakkola and M. I. Jordan. Variational probabilistic inference and the QMR-DT network. *Journal of Artificial Intelligence Research*, 10:291–322, 1999.
- [120] E T Jaynes. Information theory and statistical mechanics. *Phys. Rev.*, 106:620–630, 1957.
- [121] J. K. Johnson, D. M. Malioutov, and A. S. Willsky. Lagrangian relaxation for MAP estimation in graphical models. In *Proceedings of the Allerton Conference on Control, Communication and Computing*, Urbana-Champaign, IL, September 2007.
- [122] M. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. Saul. An introduction to variational methods for graphical models. In *Learning in graphical models*, pages 105–161. MIT Press, 1999.

- [123] T. Kailath, A. H. Sayed, and B. Hassibi. *Linear Estimation*. Prentice Hall, Englewood Cliffs, NJ, 2000.
- [124] H. Kappen and P. Rodriguez. Efficient learning in Boltzmann machines using linear response theory. *Neural Computation*, 10:1137–1156, 1998.
- [125] H. J. Kappen and W. Wiegerinck. Novel iteration schemes for the cluster variation method. In *Neural Information Processing Systems 14*, pages 415–422. MIT Press, Cambridge, Ma, 2002.
- [126] D. Karger and N. Srebro. Learning Markov networks: maximum bounded tree-width graphs. In *Symposium on Discrete Algorithms*, pages 392–401, 2001.
- [127] S. Karlin and W. Studden. *Chebyshev Systems, with Applications in Analysis and Statistics*. Interscience Publishers, New York, 1966.
- [128] R. Karp. Reducibility among combinatorial problems. In *Complexity of Computer Computations*, pages 85–103. Plenum Press, New York, 1972.
- [129] P. W. Kastelyn. Dimer statistics and phase transitions. *Journal of Mathematical Physics*, 4:287–293, 1963.
- [130] R. Kikuchi. The theory of cooperative phenomena. *Physical Review*, 81:988–1003, 1951.
- [131] S. Kim and M. Kojima. Second order cone programming relaxation of non-convex quadratic optimization problems. Technical report, Tokyo Institute of Technology, July 2000.
- [132] J. Kleinberg and E. Tardos. Approximation algorithms for classification problems with pairwise relationships: Metric labeling and Markov random fields. *Journal of the ACM*, 49:616–639, 2002.
- [133] R. Koetter and P. O. Vontobel. Graph-covers and iterative decoding of finite length codes. In *Proceedings of the 3rd International Symposium on Turbo Codes*, pages 75–82, Brest, France, 2003.
- [134] V. Kolmogorov. Convergent tree-reweighted message-passing for energy minimization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(10):1568–1583, 2006.
- [135] V. Kolmogorov and M. J. Wainwright. On optimality properties of tree-reweighted message-passing. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, pages 316–322. AUAI Press, 2005.
- [136] N. Komodakis, N. Paragios, and G. Tziritas. MRF optimization via dual decomposition: Message-passing revisited. In *International Conference on Computer Vision*, Rio de Janeiro, Brazil, 2007.
- [137] V. K. Koval and M. I. Schlesinger. Two-dimensional programming in image analysis problems. *USSR Academy of Science, Automatics and Telemechanics*, 8:149–168, 1976.
- [138] A. Krogh, B. Larsson, G. von Heijne, and E. L. L. Sonnhammer. Predicting transmembrane protein topology with a hidden Markov model: Application to complete genomes. *Journal of Molecular Biology*, 305(3):567–580, 2001.
- [139] F. R. Kschischang and B. J. Frey. Iterative decoding of compound codes by probability propagation in graphical models. *IEEE Selected Areas in Communications*, 16(2):219–230, 1998.

- [140] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger. Factor graphs and the sum-product algorithm. *IEEE Transactions on Information Theory*, 47(2):498–519, 2001.
- [141] A. Kulesza and F. Pereira. Structured learning with approximate inference. In *Advances in Neural Information Processing Systems*, volume 20, pages 785–792, Cambridge, MA, 2008. MIT Press.
- [142] P. Kumar, V. Kolmogorov, and P. H. S. Torr. An analysis of convex relaxations for MAP estimation. In *Advances in Neural Information Processing Systems*, volume 20, pages 1041–1048, Cambridge, MA, 2008. MIT Press.
- [143] P. Kumar, P. H. S. Torr, and A. Zisserman. Solving Markov random fields using second order cone programming. *IEEE Conference of Computer Vision and Pattern Recognition*, pages 1045–1052, 2006.
- [144] J. B. Lasserre. An explicit equivalent positive semidefinite program for non-linear 0-1 programs. *SIAM Journal on Optimization*, 12:756–769, 2001.
- [145] J. B. Lasserre. Global optimization with polynomials and the problem of moments. *SIAM Journal on Optimization*, 11(3):796–817, 2001.
- [146] M. Laurent. Semidefinite relaxations for Max-Cut. In *The Sharpest Cut: Festschrift in Honor of M. Padberg's 60th Birthday*, New York, 2002. MPS-SIAM Series in Optimization.
- [147] M. Laurent. A comparison of the Sherali-Adams, Lovász-Schrijver and Lasserre relaxations for 0-1 programming. *Mathematics of Operations Research*, 28:470–496, 2003.
- [148] S. L. Lauritzen. *Lectures on Contingency Tables*. Department of Mathematics, Aalborg University, 1989.
- [149] S. L. Lauritzen. Propagation of probabilities, means and variances in mixed graphical association models. *Journal of the American Statistical Association*, 87:1098–1108, 1992.
- [150] S. L. Lauritzen. *Graphical Models*. Oxford University Press, Oxford, 1996.
- [151] S. L. Lauritzen and D. J. Spiegelhalter. Local computations with probabilities on graphical structures and their application to expert systems. *Journal of the Royal Statistical Society, Series B*, 50:155–224, 1988.
- [152] M. A. R. Leisink and H. J. Kappen. Learning in higher order Boltzmann machines using linear response. *Neural Networks*, 13:329–335, 2000.
- [153] M. A. R. Leisink and H. J. Kappen. A tighter bound for graphical models. In *Advances in Neural Information Processing Systems*, volume 13, pages 266–272, Cambridge, MA, 2001. MIT Press.
- [154] H. A. Loeliger. An introduction to factor graphs. *IEEE Signal Processing Magazine*, 21:28–41, 2004.
- [155] L. Lovász. Submodular functions and convexity. In A. Bachem, M. Grötschel, and B. Korte, editors, *Mathematical Programming: The State of the Art*, pages 235–257. Springer-Verlag, New York, 1983.
- [156] L. Lovász and A. Schrijver. Cones of matrices, set-functions and 0-1 optimization. *SIAM Journal of Optimization*, 1:166–190, 1991.
- [157] M. Luby, M. Mitzenmacher, M. A. Shokrollahi, and D. Spielman. Improved low-density parity check codes using irregular graphs. *IEEE Transactions on Information Theory*, 47:585–598, 2001.

- [158] D. M. Malioutov, J. M. Johnson, and A. S. Willsky. Walk-sums and belief propagation in gaussian graphical models. *Journal of Machine Learning Research*, 7:2013–2064, 2006.
- [159] E. Maneva, E. Mossel, and M. J. Wainwright. A new look at survey propagation and its generalizations. *Journal of the ACM*, 54(4):2–41, 2007.
- [160] C. D. Manning and H. Schütze. *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge, MA, 1999.
- [161] S. Maybeck. *Stochastic Models, Estimation, and Control*. Academic Press, New York, 1982.
- [162] J. McAuliffe, L. Pachter, and M. I. Jordan. Multiple-sequence functional annotation and the generalized hidden Markov phylogeny. *Bioinformatics*, 20:1850–1860, 2004.
- [163] R. J. McEliece and M. Yildirim. Belief propagation on partially ordered sets. In D. Gilliam and J. Rosenthal, editors, *Mathematical Theory of Systems and Networks*. Institute for Mathematics and its Applications, Minneapolis, MN, 2002.
- [164] R.J. McEliece, D.J.C. McKay, and J.F. Cheng. Turbo decoding as an instance of Pearl’s belief propagation algorithm. *IEEE Journal on Selected Areas in Communications*, 16(2):140–152, February 1998.
- [165] T. Meltzer, C. Yanover, and Y. Weiss. Globally optimal solutions for energy minimization in stereo vision using reweighted belief propagation. In *International Conference on Computer Vision*, pages 428–435, Silver Springs, MD, 2005. IEEE Computer Society.
- [166] M. Mézard and A. Montanari. *Information, Physics and Computation*. Oxford University Press, Oxford, 2008.
- [167] M. Mézard, G. Parisi, and R. Zecchina. Analytic and algorithmic solution of random satisfiability problems. *Science*, 297, 812, 2002.
- [168] M. Mézard and R. Zecchina. Random K-satisfiability: From an analytic solution to an efficient algorithm. *Physical Review E*, 66:056126, 2002.
- [169] T. Minka. Expectation propagation and approximate Bayesian inference. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, pages 362–369. Morgan Kaufmann, 2001.
- [170] T. Minka. Power EP. Technical Report MSR-TR-2004-149, Microsoft Research, October 2004.
- [171] T. Minka and Y. Qi. Tree-structured approximations by expectation propagation. In *Advances in Neural Information Processing Systems*, volume 16, pages 193–200. Morgan Kaufmann, 2004.
- [172] T. P. Minka. *A family of algorithms for approximate Bayesian inference*. PhD thesis, MIT, January 2001.
- [173] C. Moallemi and B. van Roy. Convergence of the min-sum message-passing algorithm for quadratic optimization. Technical report, Stanford University, March 2006.
- [174] C. Moallemi and B. van Roy. Convergence of the min-sum algorithm for convex optimization. Technical report, Stanford University, May 2007.

- [175] J. M. Mooij and H. J. Kappen. Sufficient conditions for convergence of loopy belief propagation. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, pages 396–403. AUAI Press, 2005.
- [176] R. Neal and G. E. Hinton. A view of the EM algorithm that justifies incremental, sparse, and other variants. In *Learning in Graphical Models*. MIT Press, Cambridge, MA, 1999.
- [177] G. L. Nemhauser and L. A. Wolsey. *Integer and Combinatorial Optimization*. Wiley-Interscience, New York, 1999.
- [178] M. Opper and D. Saad. Adaptive TAP equations. In M. Opper and D. Saad, editors, *Advanced Mean Field Methods: Theory and Practice*, pages 85–98. MIT Press, Cambridge, MA, 2001.
- [179] M. Opper and O. Winther. Gaussian processes for classification: Mean field algorithms. *Neural Computation*, 12(11):2177–2204, 2000.
- [180] M. Opper and O. Winther. Tractable approximations for probabilistic models: The adaptive Thouless-Anderson-Palmer approach. *Physical Review Letters*, 64:3695, 2001.
- [181] M. Opper and O. Winther. Expectation-consistent approximate inference. *Journal of Machine Learning Research*, 6:2177–2204, 2005.
- [182] J. G. Oxley. *Matroid Theory*. Oxford University Press, Oxford, 1992.
- [183] M. Padberg. The Boolean quadric polytope: some characteristics, facets and relatives. *Mathematical Programming*, 45:139–172, 1989.
- [184] P. Pakzad and V. Anantharam. Iterative algorithms and free energy minimization. In *Annual Conference on Information Sciences and Systems*, Princeton, NJ, 2002.
- [185] P. Pakzad and V. Anantharam. Estimation and marginalization using Kikuchi approximation methods. *Neural Computation*, 17(8):1836–1873, 2005.
- [186] G. Parisi. *Statistical Field Theory*. Addison-Wesley, New York, 1988.
- [187] P. Parrilo. Semidefinite programming relaxations for semialgebraic problems. *Mathematical Programming, Series B*, 96:293–320, 2003.
- [188] J. Pearl. *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann, 1988.
- [189] J. S. Pedersen and J. Hein. Gene finding with a hidden Markov model of genome structure and evolution. *Bioinformatics*, 19:219–227, 2003.
- [190] P. Pevzner. *Computational Molecular Biology: An Algorithmic Approach*. MIT Press, Cambridge, MA, 2000.
- [191] T. Plefka. Convergence condition of the TAP equation for the infinite-ranged Ising model. *Journal of Physics A*, 15(6):1971–1978, 1982.
- [192] L. R. Rabiner and B. H. Juang. *Fundamentals of Speech Recognition*. Prentice Hall, Englewood Cliffs, N.J., 1993.
- [193] P. Ravikumar, A. Agarwal, and M. J. Wainwright. Message-passing for graph-structured linear programs: Proximal projections, convergence and rounding schemes. In *Intl. Conf. Machine Learning (ICML)*, pages 800–807, New York, 2008. ACM Press.
- [194] P. Ravikumar and J. Lafferty. Quadratic programming relaxations for metric labeling and markov random field map estimation. In *Intl. Conf. Machine Learning (ICML)*, pages 737–744, 2006.



- [195] T. Richardson and R. Urbanke. The capacity of low-density parity check codes under message-passing decoding. *IEEE Transactions on Information Theory*, 47:599–618, 2001.
- [196] T. Richardson and R. Urbanke. *Modern Coding Theory*. Cambridge University Press, Cambridge, UK, 2008.
- [197] B. D. Ripley. *Spatial Statistics*. Wiley, New York, 1981.
- [198] G. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, 1970.
- [199] T. G. Roosta, M. J. Wainwright, and S. S. Sastry. Convergence analysis of reweighted sum-product algorithms. *IEEE Transactions on Signal Processing*, 56(9):4293–4305, 2008.
- [200] P. Rusmevichientong and B. Van Roy. An analysis of turbo decoding with Gaussian densities. In *Advances in Neural Information Processing Systems*, volume 12, pages 575–581, Cambridge, MA, 2000. MIT Press.
- [201] S. Sanghavi, D. Malioutov, and A. Willsky. Linear programming analysis of loopy belief propagation for weighted matching. In *Advances in Neural Information Processing Systems*, volume 20, pages 1273–1280, Cambridge, MA, 2007. MIT Press.
- [202] S. Sanghavi, D. Shah, and A. Willsky. Message-passing for max-weight independent set. In *Neural Information Processing Systems*, Vancouver, Canada, December 2007.
- [203] L. K. Saul and M. I. Jordan. Boltzmann chains and hidden Markov models. In *Advances in Neural Information Processing Systems*, volume 7, pages 435–442. MIT Press, Cambridge, MA, 1995.
- [204] L. K. Saul and M. I. Jordan. Exploiting tractable substructures in intractable networks. In *Advances in Neural Information Processing Systems*, volume 8, pages 486–492, Cambridge, MA, 1996. MIT Press.
- [205] M. I. Schlesinger. Syntactic analysis of two-dimensional visual signals in noisy conditions. *Kibernetika*, 4:113–130, 1976.
- [206] A. Schrijver. *Theory of Linear and Integer Programming*. Wiley-Interscience Series in Discrete Mathematics, New York, 1989.
- [207] A. Schrijver. *Combinatorial Optimization: Polyhedra and Efficiency*. Springer-Verlag, New York, 2003.
- [208] M. Seeger. Expectation propagation for exponential families. Technical report, Max Planck Institute, Tuebingen, November 2005.
- [209] M. Seeger, F. Steinke, and K. Tsuda. Bayesian inference and optimal design in the sparse linear model. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, Puerto Rico, 2007.
- [210] P. D. Seymour. Matroids and multi-commodity flows. *European J. Combin.*, 2:257–290, 1981.
- [211] G. R. Shafer and P. P. Shenoy. Probability propagation. *Annals of Mathematics and Artificial Intelligence*, 2:327–352, 1990.
- [212] H. D. Sherali and W. P. Adams. A hierarchy of relaxations between the continuous and convex hull representations for zero-one programming problems. *SIAM Journal on Discrete Mathematics*, 3:411–430, 1990.

- [213] A. Siepel and D. Haussler. Combining phylogenetic and hidden Markov models in biosequence analysis. In *Proceedings of the Seventh Annual International Conference on Computational Biology*, pages 277–286, 2003.
- [214] D. Sontag and T. Jaakkola. New outer bounds on the marginal polytope. In *Advances in Neural Information Processing Systems*, volume 20, pages 1393–1400, Cambridge, MA, 2007. MIT Press.
- [215] T. P. Speed and H. T. Kiiveri. Gaussian Markov distributions over finite graphs. *Annals of Statistics*, 14(1):138–150, 1986.
- [216] R. P. Stanley. *Enumerative Combinatorics*, volume 1. Cambridge University Press, Cambridge, UK, 1997.
- [217] R. P. Stanley. *Enumerative Combinatorics*, volume 2. Cambridge University Press, Cambridge, UK, 1997.
- [218] E. B. Sudderth, M. J. Wainwright, and A. S. Willsky. Embedded trees: Estimation of Gaussian processes on graphs with cycles. *IEEE Transactions on Signal Processing*, 52(11):3136–3150, 2004.
- [219] E. B. Sudderth, M. J. Wainwright, and A. S. Willsky. Loop series and Bethe variational bounds for attractive graphical models. In *Advances in Neural Information Processing Systems*, volume 20, pages 1425–1432, Cambridge, MA, 2008. MIT Press.
- [220] C. Sutton and A. McCallum. Piecewise training of undirected models. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*. Morgan Kaufmann, July 2005.
- [221] C. Sutton and A. McCallum. Improved dynamic schedules for belief propagation. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*. Morgan Kaufmann, 2007.
- [222] R. Szeliski, R. Zabih, D. Scharstein, O. Veskler, V. Kolmogorov, A. Agarwala, M. Tappen, and C. Rother. Comparative study of energy minimization methods for Markov random fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30:1068–1080, 2008.
- [223] M. H. Taghavi and P. H. Siegel. Adaptive linear programming decoding. In *International Symposium on Information Theory*, Seattle, WA, July 2006.
- [224] K. Tanaka and T. Morita. Cluster variation method and image restoration problem. *Physics Letters A*, 203:122–128, 1995.
- [225] S. Tatikonda and M. I. Jordan. Loopy belief propagation and Gibbs measures. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, pages 493–500. Morgan Kaufmann, 2002.
- [226] Y. W. Teh and M. Welling. On improving the efficiency of the iterative proportional fitting procedure. In *Workshop on Artificial Intelligence and Statistics*, 2003.
- [227] A. Thomas, A. Gutin, V. Abkevich, and A. Bansal. Multilocus linkage analysis by blocked Gibbs sampling. *Statistics and Computing*, 10:259–269, 2000.
- [228] A. W. van der Vaart. *Asymptotic statistics*. Cambridge University Press, Cambridge, UK, 1998.
- [229] L. Vandenberghe and S. Boyd. Semidefinite programming. *SIAM Review*, 38:49–95, 1996.

- [230] L. Vandenberghe, S. Boyd, and S. Wu. Determinant maximization with linear matrix inequality constraints. *SIAM Journal on Matrix Analysis and Applications*, 19:499–533, 1998.
- [231] V. V. Vazirani. *Approximation Algorithms*. Springer-Verlag, New York, 2003.
- [232] S. Verdú and H. V. Poor. Abstract dynamic programming models under commutativity conditions. *SIAM Journal of Control and Optimization*, 25(4):990–1006, 1987.
- [233] P. O. Vontobel and R. Koetter. Lower bounds on the minimum pseudo-weight of linear codes. In *International Symposium on Information Theory*, Chicago, IL, June 2004.
- [234] P. O. Vontobel and R. Koetter. Towards low-complexity linear-programming decoding. In *Proceedings of the International Conference on Turbo Codes and Related Topics*, Munich, Germany, April 2006.
- [235] M. J. Wainwright. *Stochastic processes on graphs with cycles: geometric and variational approaches*. PhD thesis, MIT, May 2002.
- [236] M. J. Wainwright. Estimating the “wrong” graphical model: Benefits in the computation-limited regime. *Journal of Machine Learning Research*, 7:1829–1859, 2006.
- [237] M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky. Tree-based reparameterization framework for analysis of sum-product and related algorithms. *IEEE Transactions on Information Theory*, 49(5):1120–1146, 2003.
- [238] M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky. Tree-reweighted belief propagation algorithms and approximate ML estimation by pseudomoment matching. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, January 2003.
- [239] M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky. Tree consistency and bounds on the max-product algorithm and its generalizations. *Statistics and Computing*, 14:143–166, 2004.
- [240] M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky. Exact MAP estimates via agreement on (hyper)trees: Linear programming and message-passing. *IEEE Transactions on Information Theory*, 51(11):3697–3717, 2005.
- [241] M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky. A new class of upper bounds on the log partition function. *IEEE Transactions on Information Theory*, 51(7):2313–2335, 2005.
- [242] M. J. Wainwright and M. I. Jordan. Treewidth-based conditions for exactness of the Sherali-Adams and Lasserre relaxations. Technical Report 671, University of California, Berkeley, Department of Statistics, September 2004.
- [243] M. J. Wainwright and M. I. Jordan. Log-determinant relaxation for approximate inference in discrete Markov random fields. *IEEE Transactions on Signal Processing*, 54(6):2099–2109, 2006.
- [244] Y. Weiss. Correctness of local probability propagation in graphical models with loops. *Neural Computation*, 12:1–41, 2000.
- [245] Y. Weiss and W. T. Freeman. Correctness of belief propagation in Gaussian graphical models of arbitrary topology. In *Advances in Neural Information Processing Systems*, volume 12, pages 673–679, Cambridge, MA, 2000. MIT Press.

- [246] Y. Weiss, C. Yanover, and T. Meltzer. MAP estimation, linear programming, and belief propagation with convex free energies. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*. AUAI Press, 2007.
- [247] M. Welling, T. Minka, and Y. W. Teh. Structured region graphs: Morphing EP into GBP. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, pages 609–614. AUAI Press, 2005.
- [248] M. Welling and S. Parise. Bayesian random fields: The Bethe-Laplace approximation. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*. AUAI Press, 2006.
- [249] M. Welling and Y. W. Teh. Belief optimization: A stable alternative to loopy belief propagation. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, pages 554–561. Morgan Kaufmann, 2001.
- [250] M. Welling and Y. W. Teh. Linear response for approximate inference. In *Advances in Neural Information Processing Systems*, Cambridge, MA, 2003. MIT Press.
- [251] T. Werner. A linear programming approach to max-sum problem: A review. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(7):1165–1179, 2007.
- [252] N. Wiberg. *Codes and decoding on general graphs*. PhD thesis, University of Linköping, Sweden, 1996.
- [253] N. Wiberg, H. A. Loeliger, and R. Koetter. Codes and iterative decoding on general graphs. *European Transactions on Telecommunications*, 6:513–526, 1995.
- [254] W. Wiegnerinck. Variational approximations between mean field theory and the junction tree algorithm. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, pages 626–633. Morgan Kaufmann Publishers, 2000.
- [255] W. Wiegnerinck. Approximations with reweighted generalized belief propagation. In *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics*, pages 421–428, Barbardos, 2005.
- [256] W. Wiegnerinck and T. Heskes. Fractional belief propagation. In *Advances in Neural Information Processing Systems*, volume 12, pages 438–445, Cambridge, MA, 2002. MIT Press.
- [257] A. S. Willsky. Multiresolution Markov models for signal and image processing. *Proceedings of the IEEE*, 90(8):1396–1458, 2002.
- [258] J.W. Woods. Markov image modeling. *IEEE Transactions on Automatic Control*, 23:846–850, 1978.
- [259] N. Wu. *The maximum entropy method*. Springer, New York, 1997.
- [260] C. Yanover, T. Meltzer, and Y. Weiss. Linear programming relaxations and belief propagation: An empirical study. *Journal of Machine Learning Research*, 7:1887–1907, 2006.
- [261] C. Yanover, O. Schueler-Furman, and Y. Weiss. Minimizing and learning energy functions for side-chain prediction. In *RECOMB*, volume 1, pages 381–395, 2007.
- [262] J. S. Yedidia. An idiosyncratic journey beyond mean field theory. In M. Opper and D. Saad, editors, *Advanced Mean Field Methods: Theory and Practice*, pages 21–36. MIT Press, Cambridge, MA, 2001.

- [263] J. S. Yedidia, W. T. Freeman, and Y. Weiss. Generalized belief propagation. In *Advances in Neural Information Processing Systems*, volume 13, pages 689–695, Cambridge, MA, 2001. MIT Press.
- [264] J.S. Yedidia, W. T. Freeman, and Y. Weiss. Constructing free energy approximations and generalized belief propagation algorithms. *IEEE Transactions on Information Theory*, 51(7):2282–2312, 2005.
- [265] A. Yuille. CCCP algorithms to minimize the Bethe and Kikuchi free energies: Convergent alternatives to belief propagation. *Neural Computation*, 14:1691–1722, 2002.
- [266] G. M. Ziegler. *Lectures on Polytopes*. Springer-Verlag, New York, 1995.