

Exponential Families

David M. Blei

1 Introduction

- We discuss the **exponential family**, a very flexible family of distributions.
- Most distributions that you have heard of are in the exponential family.
 - Bernoulli, Gaussian, Multinomial, Dirichlet, Gamma, Poisson, Beta

2 Set-up

- An exponential family distribution has the following form,

$$p(x | \eta) = h(x) \exp\{\eta^\top t(x) - a(\eta)\} \quad (1)$$

- The different parts of this equation are
 - The natural parameter η
 - The sufficient statistic $t(x)$
 - The underlying measure $h(x)$, e.g., counting measure or Lebesgue measure
 - The log normalizer $a(\eta)$,

$$a(\eta) = \log \int h(x) \exp\{\eta^\top t(x)\}. \quad (2)$$

Here we integrate the unnormalized density over the sample space. This ensures that the density integrates to one.

- The statistic $t(x)$ is called *sufficient* because the likelihood for η only depends on x through $t(x)$.
- The exponential family has fundamental connections to the world of graphical models. For our purposes, we'll use exponential families as components in directed graphical models, e.g., in the mixtures of Gaussians.

3 The Gaussian distribution

- As a running example, consider the Gaussian distribution.
- The familiar form of the univariate Gaussian is

$$p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\} \quad (3)$$

- We put it in exponential family form by expanding the square

$$p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}} \exp \left\{ \frac{\mu}{\sigma^2} x - \frac{1}{2\sigma^2} x^2 - \frac{1}{2\sigma^2} \mu^2 - \log \sigma \right\} \quad (4)$$

- We see that

$$\eta = \langle \mu/\sigma^2, -1/2\sigma^2 \rangle \quad (5)$$

$$t(x) = \langle x, x^2 \rangle \quad (6)$$

$$a(\eta) = \mu^2/2\sigma^2 + \log \sigma \quad (7)$$

$$= -\eta_1^2/4\eta_2 - (1/2) \log(-2\eta_2) \quad (8)$$

$$h(x) = 1/\sqrt{2\pi} \quad (9)$$

- If you are new to this, work it out for others on the list.

4 Moments

- The derivatives of the log normalizer gives the moments of the sufficient statistics,

$$\frac{d}{d\eta} a(\eta) = \frac{d}{d\eta} \left(\log \int \exp \{ \eta^\top t(x) \} h(x) dx \right) \quad (10)$$

$$= \frac{\int t(x) \exp \{ \eta^\top t(x) \} h(x) dx}{\int \exp \{ \eta^\top t(x) \} h(x) dx} \quad (11)$$

$$= \int t(x) \exp \{ \eta^\top t(x) - a(\eta) \} h(x) dx \quad (12)$$

$$= \mathbb{E} [t(X)] \quad (13)$$

- The next derivatives are higher moments. The second derivative is the variance, etc.
- Let's go back to the Gaussian example.

– The derivative with respect to η_1 is

$$\frac{da(\eta)}{d\eta_1} = -\frac{\eta_1}{2\eta_2} \quad (14)$$

$$= \mu \quad (15)$$

$$= E[X] \quad (16)$$

– The derivative with respect to η_2 is

$$\frac{da(\eta)}{d\eta_2} = \frac{\eta_1^2}{4\eta_2^2} - \frac{1}{2\eta_2} \quad (17)$$

$$= \sigma^2 + \mu^2 \quad (18)$$

$$= E[X^2] \quad (19)$$

– This means that the variance is

$$\text{Var}(X) = E[X^2] - E[X]^2 \quad (20)$$

$$= -\frac{1}{2\eta_2} \quad (21)$$

- In a **minimal exponential family**, the components of the sufficient statistics $t(x)$ are linearly independent.
- In a minimal exponential family, the mean $\mu := E[t(X)]$ is another parameterization of the distribution. That is, there is a 1-1 mapping between η and μ .
 - The function $a(\eta)$ is convex. (It is log-sum-exponential.)
 - Thus there is a 1-1 mapping between its argument and its derivative.
 - Thus there is a 1-1 mapping between η and $E[t(X)]$.
- Side note: the MLE of an exponential family matches the mean parameters with the empirical statistics of the data.
 - Assume $x_{1:n}$ are from an exponential family.
 - Find $\hat{\eta}$ that maximizes the likelihood of x .
 - This is the η such that $E[t(X)] = (1/n) \sum_i t(x_i)$.

5 Conjugacy

- Consider the following set up:

$$\eta \sim F(\cdot | \lambda) \quad (22)$$

$$x_i \sim G(\cdot | \eta) \quad \text{for } i \in \{1, \dots, n\}. \quad (23)$$

- This is a classical Bayesian data analysis setting. And, this is used as a component in more complicated models, e.g., in hierarchical models.
- The posterior distribution of η given the data $x_{1:n}$ is

$$p(\eta | x_{1:n}, \lambda) \propto F(\eta | \lambda) \prod_{i=1}^n G(x_i | \eta). \quad (24)$$

When this distribution is in the same family as F , i.e., when its parameters are part of the parameter-space defined by λ , then we say that F and G make a **conjugate pair**.

- For example,
 - A Gaussian likelihood with fixed variance, and a Gaussian prior on the mean
 - A multinomial likelihood and a Dirichlet prior on the probabilities
 - A Bernoulli likelihood and a beta prior on the bias
 - A Poisson likelihood and a gamma prior on the rate

In all these settings, the conditional distribution of the parameter given the data is in the same family as the prior.

- Suppose the data come from an exponential family. Every exponential family has a conjugate prior (in theory),

$$p(x_i | \eta) = h_\ell(x) \exp\{\eta^\top t(x_i) - a_\ell(\eta)\} \quad (25)$$

$$p(\eta | \lambda) = h_c(\eta) \exp\{\lambda_1^\top \eta + \lambda_2^\top (-a_\ell(\eta)) - a_c(\lambda)\}. \quad (26)$$

- The natural parameter $\lambda = \langle \lambda_1, \lambda_2 \rangle$ has dimension $\dim(\eta) + 1$.
- The sufficient statistics are $\langle \eta, -a(\eta) \rangle$.
- The other terms depend on the form of the exponential family. For example, when η are multinomial parameters then the other terms help define a Dirichlet.

- Let's compute the posterior,

$$p(\eta | x_{1:n}, \lambda) \propto p(\eta | \lambda) \prod_{i=1}^n p(x_i | \eta) \quad (27)$$

$$= h(\eta) \exp\{\lambda_1^\top \eta + \lambda_2(-a(\eta)) - a_c(\lambda)\} \quad (28)$$

$$\cdot \left(\prod_{i=1}^n h(x_i) \right) \exp\{\eta^\top \sum_{i=1}^n t(x_i) - na_x(\eta)\} \quad (29)$$

$$\propto h(\eta) \exp\{(\lambda_1 + \sum t(x_i))^\top \eta + (\lambda_2 + n)(-a(\eta))\}. \quad (30)$$

This is the same exponential family as the prior, with parameters

$$\hat{\lambda}_1 = \lambda_1 + \sum_{i=1}^n t(x_i) \quad (31)$$

$$\hat{\lambda}_2 = \lambda_2 + n. \quad (32)$$

6 Example: Data from a unit variance Gaussian

- Suppose the data x_i come from a unit variance Gaussian

$$p(x | \mu) = \frac{1}{\sqrt{2\pi}} \exp\{-(x - \mu)^2/2\}. \quad (33)$$

- This is a simpler exponential family than the previous Gaussian

$$p(x | \mu) = \frac{\exp\{-x^2/2\}}{\sqrt{2\pi}} \exp\{\mu x - \mu^2/2\}. \quad (34)$$

In this case

$$\eta = \mu \quad (35)$$

$$t(x) = x \quad (36)$$

$$h(x) = \frac{\exp\{-x^2/2\}}{\sqrt{2\pi}} \quad (37)$$

$$a(\eta) = \mu^2/2 = \eta^2/2. \quad (38)$$

- We are interested in the conjugate prior. (State the end result on the next page.)

- Consider a model with an unknown mean. What is the conjugate prior? It is

$$p(\eta | \lambda) = h(\eta) \exp\{\lambda_1 \eta + \lambda_2(-\eta^2/2) - a_c(\lambda)\} \quad (39)$$

- Set $\lambda_1^* = \lambda_1$ and $\lambda_2^* = -\lambda_2/2$. This means the sufficient statistics are $\langle \eta, \eta^2 \rangle$.

- This is a **Gaussian distribution**. We now know the conjugate prior.

- Now consider the posterior,

$$\hat{\lambda}_1 = \lambda_1 + \sum_{i=1}^n x_i \quad (40)$$

$$\hat{\lambda}_2 = \lambda_2 + n \quad (41)$$

$$\hat{\lambda}_2^* = \frac{-(\lambda_2 + n)}{2}. \quad (42)$$

- Let's map this back to traditional Gaussian parameters.

- The mean is

$$E[\mu | x_{1:n}, \lambda] = \frac{\lambda_1 + \sum_{i=1}^n x_i}{\lambda_2 + n} \quad (43)$$

- The variance is

$$\text{Var}(\mu | x_{1:n}, \lambda) = \frac{1}{\lambda_2 + n} \quad (44)$$

- Finally, for closure, let's parameterize everything in the mean parameterization.

- Consider a prior mean and prior variance $\{\mu_0, \sigma_0^2\}$.

- We know that

$$\lambda_1 = \mu_0 / \sigma_0^2 \quad (45)$$

$$\lambda_2 = -1/2\sigma_0^2 \quad (46)$$

$$\lambda_2^* = 1/\sigma_0^2. \quad (47)$$

The expression λ_2^* is also called the **precision**.

- So the posterior mean is

$$E[\mu | x_{1:n}, \mu_0, \sigma_0^2] = \frac{\mu_0 / \sigma_0^2 + \sum_{i=1}^n x_i}{1/\sigma_0^2 + n} \quad (48)$$

- The posterior variance is

$$\text{Var}(\mu | x_{1:n}, \mu_0, \sigma_0^2) = \frac{1}{1/\sigma_0^2 + n} \quad (49)$$

- Intuitively, when we haven't seen any data then our estimate of the mean is the *prior mean*. As we see more data, our estimate of the mean moves towards the *sample mean*. Before seeing data, our “confidence” about the estimate is the prior variance. As we see more data, the confidence decreases.