

COS513: FOUNDATIONS OF PROBABILISTIC MODELS

LECTURE 15

ALEX LORBERT, GUNGOR POLATKAN

1. DIMENSIONALITY REDUCTION

The goal of dimensionality reduction is to compute a reduced representation of our data. The benefits of such a reduction include visualization of data, storage of data, and the *possible* extraction of systematic structures. In general, if we have a p -dimensional vector (x_1, x_2, \dots, x_p) we wish to find a way to represent this vector with q -dimensions as $(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_q)$ with $p > q$. In this lecture we assume only real valued vectors.

2. PRINCIPAL COMPONENT ANALYSIS (PCA)

The main idea of PCA is to project our data to a lower dimensional manifold. For example, if $p = 2$ and our data “seem” linear ($q = 1$) then we wish to project the data points onto a “suitable” line (see Figure 1). This projection is not without cost since our data do not really live on a line. In PCA our free parameter is the selection of q .

There are at least three ways to think about our lower dimensional subspace:

- (1) We can maximize the variance of the projection along \mathbb{R}^q [1]. In the previous example, a selection of a horizontal line results in the projected data points being “squashed”.
- (2) We can minimize the reconstruction error, i.e. the distance between the the original data and the projected data [2]. [Note: this is not the same as regression where we minimize the RSS].
- (3) We can view PCA via an MLE of a parameter in a latent variable model [3].

3. THE MULTIVARIATE GAUSSIAN DISTRIBUTION

The probability density function of a Gaussian random vector $X \in \mathbb{R}^p$ is

$$p(x|\mu, \Sigma) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

1

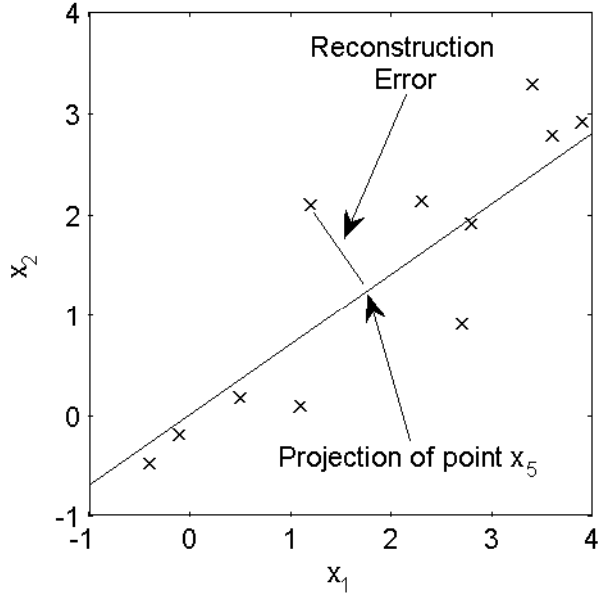


FIGURE 1. Example of dimensionality reduction with $p = 2$ and $q = 1$.

with mean $\mu \in \mathbb{R}^p$ and symmetric positive definite covariance matrix $\Sigma \in \mathbb{R}^{p \times p}$. If we let X_i denote the i^{th} element of X and σ_{ij} denote the ij^{th} element of Σ then we have the following relationships:

$$\begin{aligned} \mu_i &= \mathbb{E}[X_i] && \text{(mean)} \\ \sigma_{ij} &= \mathbb{E}[X_i X_j] - \mathbb{E}[X_i] \mathbb{E}[X_j] && \text{(covariance)} \\ \sigma_{ii} &= \mathbb{E}[X_i^2] - \mathbb{E}[X_i]^2 && \text{(variance)} \end{aligned}$$

Letting $f(x) = -\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)$ defines contours of equal probability (see Figure 2). When Σ is diagonal the elements of X are uncorrelated implying statistical independence in the case of Gaussian random vectors.

3.1. MLE of the Multivariate Gaussian. Let $X_1, \dots, X_N \in \mathbb{R}^p$ denote iid Gaussian random vectors. The MLE is given by

$$\left(\hat{\mu}, \hat{\Sigma} \right) = \arg \max_{(\mu, \Sigma)} \sum_{n=1}^N \log p(x_n | \mu, \Sigma)$$

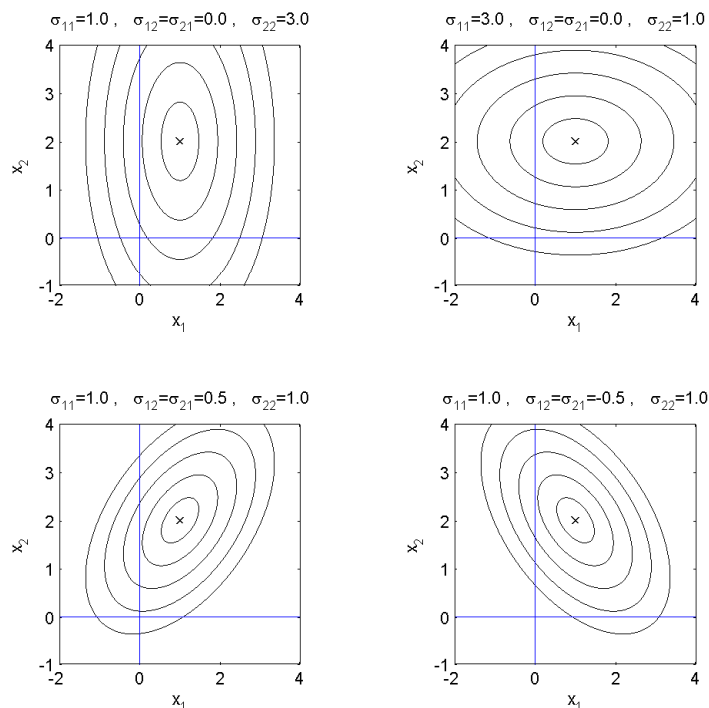


FIGURE 2. Some examples of multivariate Gaussian equiprobable contours in 2 dimensions with $\mu = [1 \ 2]^T$

and its solution is given by

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\hat{\Sigma} = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu})(x_n - \hat{\mu})^T$$

3.2. Subvectors of Multivariate Gaussian Random Vectors. For p -dimensional Gaussian random vector $X = \langle X_1, X_2, \dots, X_p \rangle$ we can write $X = \langle \tilde{X}_1, \tilde{X}_2 \rangle$ where $\tilde{X}_1 = \langle X_1, \dots, X_k \rangle$ and $\tilde{X}_2 = \langle X_{k+1}, \dots, X_p \rangle$. It follows that \tilde{X}_1 and \tilde{X}_2 are jointly normal and for X having mean μ and covariance matrix Σ , we partition our parameters as follows:

$$\mu = \langle \mu_1, \mu_2 \rangle \quad \text{and} \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

with

$$\mu_i = \mathbb{E}[\tilde{X}_i], \quad \text{and} \quad \Sigma_{ij} = \mathbb{E}[\tilde{X}_i \tilde{X}_j^T].$$

Using the chain rule we obtain $p(\tilde{x}_1, \tilde{x}_2) = p(\tilde{x}_2)p(\tilde{x}_1|\tilde{x}_2)$. To obtain the marginal of \tilde{X}_2 , denoted X_m , we marginalize over \tilde{X}_1 and obtain the Gaussian random vector with mean $\mu_m = \mu_2$ and covariance matrix $\Sigma_m = \Sigma_{22}$. The conditional variable $\tilde{X}_1|\tilde{X}_2$, denoted X_c , is also Gaussian with parameters

$$\begin{aligned}\mu_c &= \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\tilde{x}_2 - \mu_2) \\ \Sigma_c &= \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} \quad (\text{Schur complement})\end{aligned}$$

4. PCA AND FACTOR ANALYSIS (FA)

We begin with the latent variable graphical model (see Figure 3): the latent random variable $z_n \sim \mathcal{N}(\mathbf{0}, I_q)$, observed variable $x_n \sim \mathcal{N}(\mu + \Lambda z_n, \Psi)$, and with parameters $\Lambda \in \mathbb{R}^{p \times q}$ and positive definite diagonal matrix $\Psi \in \mathbb{R}^{q \times q}$ (I_k is the $k \times k$ identity matrix). Without loss of generality, we can assume $\mu = \mathbf{0}$ since we can always center the observed data.

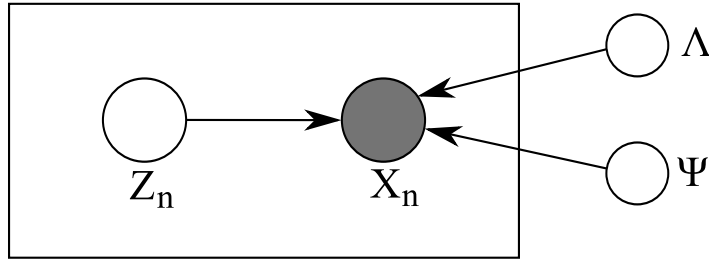


FIGURE 3. Latent variable graphical model

We can view this as a generative process as also seen by samples at Figure 4: for each $n = 1, \dots, N$

- (1) Select a random point on the q -manifold with distribution $\mathcal{N}(\mathbf{0}, I_q)$
 $\implies z_n$
- (2) Use Λ to map this random point to $\mathbb{R}^p \implies \Lambda z_n$
- (3) Select a random point in \mathbb{R}^p with distribution $\mathcal{N}(\Lambda z_n, \Psi) \implies x_n$

The difference between PCA and FA is the structure of Ψ :

$$\begin{array}{ll} \text{PCA} & \Psi = \text{diag}(\sigma^2 \mathbf{1}_p) \\ \text{FA} & \Psi = \text{diag}(\langle \sigma_1^2, \sigma_2^2, \dots, \sigma_p^2 \rangle) \end{array}$$

where $\mathbf{1}_k$ is the k -dimensional vector of all ones. The solution to PCA is exact and involves selecting the eigenvectors of $[x_1|x_2|\dots|x_n] \times [x_1|x_2|\dots|x_n]^T$ corresponding to the largest p eigenvalues (in magnitude). FA, on the other hand, does not have an explicit solution and so we rely on the EM algorithm. The graphical model above has a strong resemblance to the regression model (z_n resemble the covariates and x_n resemble the response). The

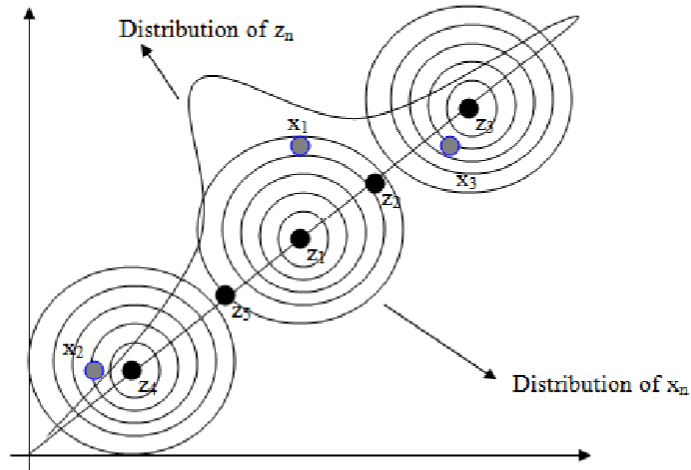


FIGURE 4. A few samples from the generative process given at Figure 3

steps of the EM algorithm are

$$\begin{array}{ll} \text{E step} & z|x \\ \text{M step} & \hat{\Lambda} \end{array}$$

where

$$\hat{\Lambda}^{(t+1)} = \left(\sum_{n=1}^N \mathbb{E} [z_n z_n^T | x] \right)^{-1} \left(\sum_{n=1}^N \mathbb{E} [z_n | x_n]^T x_n \right)$$

which resembles the normal equations.

REFERENCES

- [1] H. Hotelling, “Analysis of a Complex of Statistical Variables into Principal Components,” *J. Educational Psychology*, vol. 24, pp. 417–441, 1933.
- [2] K. Pearson, “On lines and planes of closest fit to systems of points in space,” *Philosophical Magazine*, vol. 2, no. 6, pp. 559–572, 1901.
- [3] C. M. Bishop and C. K. I. Williams, “Em optimization of latent-variable density models,” in *Advances in Neural Information Processing Systems 8*, pp. 465–471, MIT Press, 1996.