

Latent Variable Models:

Factor Analyzers and Mixture of Factor Analyzers

Seungjin Choi

Department of Computer Science and Engineering
Pohang University of Science and Technology
77 Cheongam-ro, Nam-gu, Pohang 37673, Korea
seungjin@postech.ac.kr

Outline

- ▶ Why latent variable models?
- ▶ Principal component analysis (PCA)
- ▶ Maximum likelihood factor analysis
- ▶ Probabilistic PCA
- ▶ Mixture of factor analyzers

Why Latent Variable Models?: An Example



Gaussian $p(x_1, \dots, x_D)$ requires D independent parameters for mean vector and $\frac{D(D+1)}{2}$ independent parameters for covariance matrix, $\frac{D(D+3)}{2}$ parameters in total.

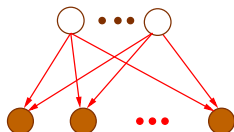
The number of independent parameters grows with D^2 .



Marginal independence assumes:

$$p(x_1, \dots, x_D) = \prod_{i=1}^D p(x_i)$$

which requires just $2D$ free parameters.



Conditional independence assumes:

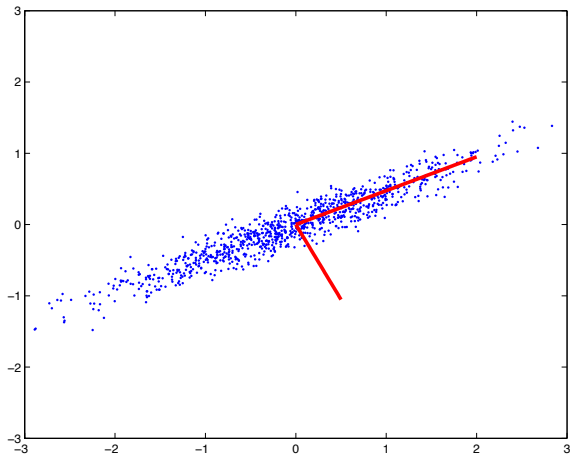
$$p(x_1, \dots, x_D | s_1, \dots, s_K) = \prod_{i=1}^D p(x_i | s_1, \dots, s_K).$$

In the case of linear models, the number of independent parameters grows with D (actually, need $(DK + 2D)$).

Principal Component Analysis (PCA)

- ▶ **Principal component analysis (PCA)** is a well-established technique for **dimension reduction**. Its applications include data compression, image processing, data visualization, exploratory data analysis, pattern recognition, and time series prediction.
- ▶ The most common derivation of PCA is in terms of **an orthogonal projection ($\mathbf{W}^\top \mathbf{W} = \mathbf{I}$) which maximizes the variance** in the projected space. Given a set of m -dimensional observation vector, $\{\mathbf{x}_t\}$, the PCA aims at finding a orthogonal linear projection $\mathbf{y} = \mathbf{W}^\top \mathbf{x}$ such that the variance of $\mathbf{y} \in \mathbb{R}^d$ ($d < D$) is maximized. The i th element of \mathbf{y} is called *i th principal component*.
- ▶ Alternatively, PCA provides **an orthogonal linear projection which minimize the squared reconstruction error $\sum_{t=1}^N \|\mathbf{x}_t - \mathbf{W}\mathbf{W}^\top \mathbf{x}_t\|^2$** . Thus PCA is the optimal linear encoding in MS sense.
- ▶ It was shown that \mathbf{w}_i (i th column vector of \mathbf{W}) corresponds to the normalized eigenvector associated with the i th largest eigenvalue of the covariance matrix $\mathbf{C}_x = \mathbb{E}\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top\}$.

PCA: An Example



Vanilla PCA

- ▶ Data centering: $\mathbf{X} \leftarrow \mathbf{X} \left(\mathbf{I} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \right)$, where $\mathbf{1}_N \in \mathbb{R}^N$ is the vector of all one's.
- ▶ Find rank- d approximation of the sample covariance matrix $\mathbf{C} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{D \times D}$:

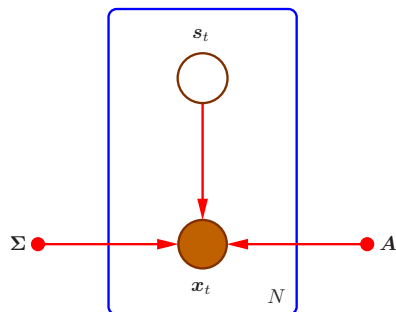
$$\mathbf{C} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{U}^\top,$$

where $\mathbf{U} \in \mathbb{R}^{D \times d}$ contains d leading eigenvectors in its columns and $\mathbf{\Sigma} \in \mathbb{R}^{d \times d}$ is diagonal matrix in which diagonal entries are eigenvalues in descending order.

- ▶ **Principal components** are computed as

$$\mathbf{Y} = \mathbf{U}^\top \mathbf{X}.$$

Factor Analysis: Graphical Model



- ▶ Gaussian latent variables:

$$p(\mathbf{s}_t) = \mathcal{N}(\mathbf{s}_t | 0, \mathbf{I}).$$

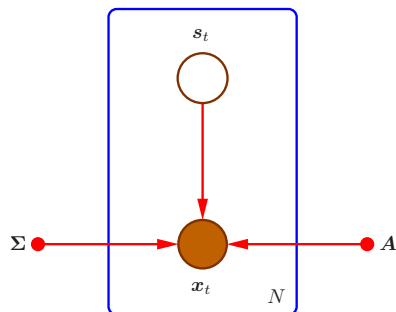
- ▶ The distribution over observed variables conditioned on latent variables is

$$p(\mathbf{x}_t | \mathbf{s}_t) = \mathcal{N}(\mathbf{x}_t | \mathbf{A}\mathbf{s}_t, \mathbf{\Sigma}),$$

where $\mathbf{\Sigma}$ is a diagonal matrix.

- ▶ Finds a **lower-dimensional projection** of high-dimensional data that captures the **covariance structure** of the data.

Learning Factor Analysis



- ▶ Compute maximum likelihood estimates of parameters

$$\{A, \Sigma\}.$$

- ▶ Compute the posterior distribution over latent variables

$$p(s_t | x_t)$$

which yields lower-dimensional projection.

- ▶ We will do this again using expectation maximization.

Factor Analysis

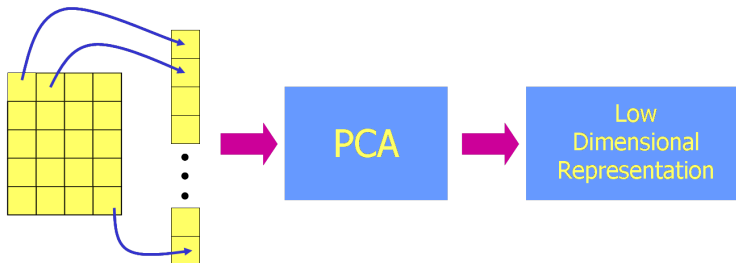
- ▶ In maximum likelihood factor analysis (FL), an D -dimensional real-valued data \mathbf{x} is modeled using an d -dimensional vector of real-valued **factors** \mathbf{s} . In general, d is much smaller than D .
- ▶ The linear generative model is given by

$$\mathbf{x} = \mathbf{A}\mathbf{s} + \epsilon,$$

where \mathbf{A} is known as the **factor loading matrix**. The factors \mathbf{s} are assumed to be $\mathcal{N}(0, \mathbf{I})$ distributed. The D -dimensional vector ϵ is $\mathcal{N}(0, \mathbf{\Sigma})$ distributed, where $\mathbf{\Sigma}$ is a diagonal matrix.

- ▶ The observed variables \mathbf{x} are independent given the factors. According to this model, \mathbf{x} is normally distributed with zero mean and covariance $\mathbf{A}\mathbf{A}^T + \mathbf{\Sigma}$.
- ▶ The goal of factor analysis is to find the \mathbf{A} and $\mathbf{\Sigma}$ that **best model the covariance structure of \mathbf{x}** .
- ▶ There is an indeterminacy in factor analysis, since a factor loading matrix $\mathbf{A}' = \mathbf{A}\mathbf{Q}$, where \mathbf{Q} is any **orthogonal matrix**, produces the same covariance of \mathbf{x} , using hidden factors given by $\mathbf{s}' = \mathbf{Q}^T \mathbf{s}$. (**Grassmann manifold**)

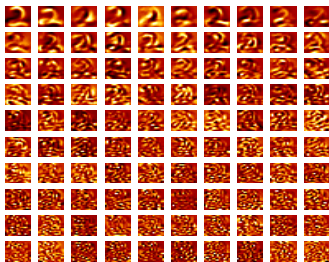
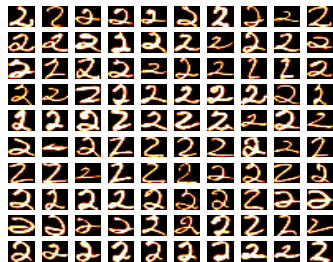
Eigenface



$$\text{Face Image} = a_1 \times \text{Eigenface}_1 + a_2 \times \text{Eigenface}_2 + a_3 \times \text{Eigenface}_3 + a_4 \times \text{Eigenface}_4 + \dots + a_n \times \text{Eigenface}_n$$

The equation shows a grayscale image of a man's face on the left, followed by an equals sign. To the right of the equals sign is a series of terms: a_1 multiplied by a grayscale image of a face profile, plus a_2 multiplied by a grayscale image of a face with a different expression, plus a_3 multiplied by another grayscale face image, plus a_4 multiplied by a grayscale face image, followed by three dots and a_n multiplied by a grayscale face image.

Eigendigits



Remarks on Factor Analysis

- ▶ The d factors play the same role as the principal components in PCA. They are informative projections of the data.
- ▶ The columns of \mathbf{A} span the space of the first d principal components.
- ▶ To compute the corresponding eigenvector and eigenvalues explicitly, the data can be projected into this d -dimensional subspace and an ordered orthogonal basis for the covariance in the subspace can be considered.
- ▶ For the case of isotropic Gaussian noise, the ML FA boils down to probabilistic PCA.
- ▶ In the case of zero noise limit, it is simplified to EM-PCA.

EM Algorithm for Maximum Likelihood Factor Analysis

- ▶ **E-step:** We compute the **posterior distribution over latent variables**, $p(\mathbf{s}|\mathbf{x})$ and calculate the **expected complete-data log-likelihood** $\mathbb{E}\mathcal{L}_c$ where the expectation is taken with respect to the posterior $p(\mathbf{s}|\mathbf{x})$.

$$\begin{aligned}\mathbb{E}\mathcal{L}_c &= \mathbb{E} \log p(\mathbf{X}, \mathbf{S} | \mathbf{A}, \mathbf{\Sigma}) \\ &= \mathbb{E} \sum_{t=1}^N \log p(\mathbf{x}_t, \mathbf{s}_t | \mathbf{A}, \mathbf{\Sigma}).\end{aligned}$$

- ▶ **M-step:** Re-estimate parameters \mathbf{A} and $\mathbf{\Sigma}$ which maximize the expected complete-data log-likelihood:

$$\begin{aligned}\mathbf{A} &\leftarrow \arg \max_{\mathbf{A}} \mathbb{E}\mathcal{L}_c, \\ \mathbf{\Sigma} &\leftarrow \arg \max_{\mathbf{\Sigma}} \mathbb{E}\mathcal{L}_c.\end{aligned}$$

Gaussian Identities

A D -dimensional Gaussian density for \mathbf{x} is

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{x} - \boldsymbol{\mu})\right\}.$$

Define the augmented vector $\mathbf{z} = [\mathbf{x}^\top, \mathbf{s}^\top]^\top$ which is jointly normal, i.e.,

$$\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \mathbf{s} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix}\right).$$

The marginal densities are $\mathbf{x} \sim \mathcal{N}(\mathbf{a}, \mathbf{A})$, $\mathbf{s} \sim \mathcal{N}(\mathbf{b}, \mathbf{B})$.

The conditional distributions are:

$$\begin{aligned} p(\mathbf{x} | \mathbf{s}) &= \mathcal{N}\left(\mathbf{a} + \mathbf{C}\mathbf{B}^{-1}(\mathbf{s} - \mathbf{b}), \mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top\right) \\ p(\mathbf{s} | \mathbf{x}) &= \mathcal{N}\left(\mathbf{b} + \mathbf{C}^\top\mathbf{A}^{-1}(\mathbf{x} - \mathbf{a}), \mathbf{B} - \mathbf{C}^\top\mathbf{A}^{-1}\mathbf{C}\right). \end{aligned}$$

E-Step for ML-FA

It follows from Gaussian identities that the posterior distribution over latent variables is calculated as

$$p(\mathbf{s}|\mathbf{x}) = \mathcal{N}(\mathbf{s}|\Phi\mathbf{x}, \mathbf{I} - \Phi\mathbf{A}),$$

where $\Phi = \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top + \Sigma)^{-1}$.

Since $p(\mathbf{s}|\mathbf{x})$ is Gaussian, we only need the following sufficient statistics to evaluate the expected complete-data log-likelihood:

$$\begin{aligned}\mathbb{E}\{\mathbf{s}_t|\mathbf{x}_t\} &= \Phi\mathbf{x}_t, \\ \mathbb{E}\{\mathbf{s}_t\mathbf{s}_t^\top|\mathbf{x}_t\} &= \text{var}(\mathbf{s}_t|\mathbf{x}_t) + \mathbb{E}\{\mathbf{s}_t|\mathbf{x}_t\} [\mathbb{E}\{\mathbf{s}_t|\mathbf{x}_t\}]^\top \\ &= \mathbf{I} - \Phi\mathbf{A} + \Phi\mathbf{x}_t\mathbf{x}_t^\top\Phi^\top.\end{aligned}$$

The complete-data log-likelihood is given by

$$\mathcal{L}_c = \log p(\mathbf{X}, \mathbf{S} | \mathbf{A}, \mathbf{\Sigma}) = \sum_{t=1}^N \log p(\mathbf{x}_t | \mathbf{s}_t) + \sum_{t=1}^N \log p(\mathbf{s}_t).$$

The expected complete-data log-likelihood is given by

$$\begin{aligned} \mathbb{E} \mathcal{L}_c &= \sum_{t=1}^N \mathbb{E} \log p(\mathbf{x}_t | \mathbf{s}_t) + \sum_{t=1}^N \mathbb{E} \log p(\mathbf{s}_t) \\ &\propto -\frac{N}{2} \log |\mathbf{\Sigma}| - \frac{1}{2} \sum_{t=1}^N \left\{ \mathbf{x}_t^\top \mathbf{\Sigma}^{-1} \mathbf{x}_t - 2 [\mathbb{E}\{\mathbf{s}_t | \mathbf{x}_t\}]^\top \mathbf{A}^\top \mathbf{\Sigma}^{-1} \mathbf{x}_t \right. \\ &\quad \left. + \text{tr} \left(\mathbf{A}^\top \mathbf{\Sigma}^{-1} \mathbf{A} \mathbb{E}\{\mathbf{s}_t \mathbf{s}_t^\top | \mathbf{x}_t\} \right) \right\}. \end{aligned}$$

M-Step for ML-FA

M-step involves re-estimating parameters \mathbf{A} and $\mathbf{\Sigma}$ which maximize $\mathbb{E}\mathcal{L}_c$:

$$\begin{aligned}\mathbf{A} &\leftarrow \arg \max_{\mathbf{A}} \mathbb{E}\mathcal{L}_c, \\ \mathbf{\Sigma} &\leftarrow \arg \max_{\mathbf{\Sigma}} \mathbb{E}\mathcal{L}_c.\end{aligned}$$

Solve $\frac{\partial \mathbb{E}\mathcal{L}_c}{\partial \mathbf{A}} = 0$ and $\frac{\partial \mathbb{E}\mathcal{L}_c}{\partial \mathbf{\Sigma}^{-1}} = 0$ for \mathbf{A} and $\mathbf{\Sigma}$, respectively, to obtain

$$\begin{aligned}\mathbf{A}^{\text{new}} &= \left(\sum_{t=1}^N \mathbf{x}_t [\mathbb{E}\{\mathbf{s}_t | \mathbf{x}_t\}]^{\top} \right) \left(\sum_{t=1}^N \mathbb{E}\{\mathbf{s}_t \mathbf{s}_t^{\top} | \mathbf{x}_t\} \right)^{-1}, \\ \mathbf{\Sigma}^{\text{new}} &= \frac{1}{N} \text{diag} \left\{ \sum_{t=1}^N \mathbf{x}_t \mathbf{x}_t^{\top} - \mathbf{A}^{\text{new}} \sum_{t=1}^N \mathbb{E}\{\mathbf{s}_t | \mathbf{x}_t\} \mathbf{x}_t^{\top} \right\}.\end{aligned}$$

Algorithm Outline: EM for ML-FA

- ▶ **E-step:** Compute sufficient statistics:

$$\begin{aligned}\mathbb{E}\{\mathbf{s}_t|\mathbf{x}_t\} &= \mathbf{\Phi}\mathbf{x}_t, \\ \mathbb{E}\{\mathbf{s}_t\mathbf{s}_t^\top|\mathbf{x}_t\} &= \mathbf{I} - \mathbf{\Phi}\mathbf{A} + \mathbf{\Phi}\mathbf{x}_t\mathbf{x}_t^\top\mathbf{\Phi}^\top,\end{aligned}$$

where $\mathbf{\Phi} = \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top + \mathbf{\Sigma})^{-1}$.

- ▶ **M-step:** Re-estimate parameters:

$$\begin{aligned}\mathbf{A}^{\text{new}} &= \left(\sum_{t=1}^N \mathbf{x}_t [\mathbb{E}\{\mathbf{s}_t|\mathbf{x}_t\}]^\top \right) \left(\sum_{t=1}^N \mathbb{E}\{\mathbf{s}_t\mathbf{s}_t^\top|\mathbf{x}_t\} \right)^{-1}, \\ \mathbf{\Sigma}^{\text{new}} &= \frac{1}{N} \text{diag} \left\{ \sum_{t=1}^N \mathbf{x}_t\mathbf{x}_t^\top - \mathbf{A}^{\text{new}} \sum_{t=1}^N \mathbb{E}\{\mathbf{s}_t|\mathbf{x}_t\}\mathbf{x}_t^\top \right\}.\end{aligned}$$

Probabilistic PCA (PPCA)

Tipping and Bishop showed that the subspace defined by the columns of the factor loading matrix \mathbf{A} corresponds to principal subspace when $\Sigma = \sigma^2 \mathbf{I}$ (isotropic Gaussian).

- ▶ **E-step:** Compute sufficient statistics:

$$\begin{aligned}\mathbb{E}\{\mathbf{s}_t | \mathbf{x}_t\} &= \mathbf{M}^{-1} \mathbf{A}^\top \mathbf{x}_t, \\ \mathbb{E}\{\mathbf{s}_t \mathbf{s}_t^\top | \mathbf{x}_t\} &= \sigma^2 \mathbf{M}^{-1} + \mathbb{E}\{\mathbf{s}_t | \mathbf{x}_t\} \mathbb{E}\{\mathbf{s}_t^\top | \mathbf{x}_t^\top\},\end{aligned}$$

where $\mathbf{M} = \mathbf{A}^\top \mathbf{A} + \sigma^2 \mathbf{I}$.

- ▶ **M-step:** Re-estimate parameters \mathbf{A} and σ^2 :

$$\begin{aligned}\mathbf{A} &\leftarrow \left[\sum_{t=1}^N \mathbf{x}_t \mathbb{E}\{\mathbf{s}_t^\top | \mathbf{x}_t^\top\} \right] \left[\sum_{t=1}^N \mathbb{E}\{\mathbf{s}_t \mathbf{s}_t^\top | \mathbf{x}_t\} \right]^{-1}, \\ \sigma^2 &\leftarrow \frac{1}{ND} \sum_{t=1}^N \left\{ \|\mathbf{x}_t\|^2 - 2 \mathbb{E}\{\mathbf{s}_t^\top | \mathbf{x}_t^\top\} \mathbf{A}^\top \mathbf{x}_t \right. \\ &\quad \left. + \text{tr}(\mathbf{E}\{\mathbf{s}_t \mathbf{s}_t^\top | \mathbf{x}_t\}) \mathbf{A}^\top \mathbf{A} \right\}.\end{aligned}$$

EM-PCA: Limiting Case

- ▶ In general, the subspace defined by the columns of the factor loading matrix \mathbf{A} does not correspond to the principal subspace of the data. However, for the case of isotropic Gaussian noise, i.e., $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$, FA behaves like PSA.
- ▶ In the case of **zero noise limit** ($\sigma^2 \rightarrow 0$), the inference reduces to simple least square projection:

$$\begin{aligned}\Phi &= \lim_{\sigma^2 \rightarrow 0} \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top + \sigma^2 \mathbf{I})^{-1} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top, \\ p(\mathbf{s}|\mathbf{x}) &= \mathcal{N}\left(\left(\mathbf{A}^\top \mathbf{A}\right)^{-1} \mathbf{A}^\top \mathbf{x}, 0\right) = \delta\left(\mathbf{s} - \left(\mathbf{A}^\top \mathbf{A}\right)^{-1} \mathbf{A}^\top \mathbf{x}\right).\end{aligned}$$

For $\mathbf{A} \in \mathbb{R}^{D \times K}$ with rank K , left multiplication by $\mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1}$ is exactly equivalent to left multiplication by $(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$.

Algorithm Outline: EM-PCA

Define

$$\begin{aligned}\mathbf{S} &= [\mathbf{s}_1, \dots, \mathbf{s}_N] \in \mathbb{R}^{d \times N}, \\ \mathbf{X} &= [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}.\end{aligned}$$

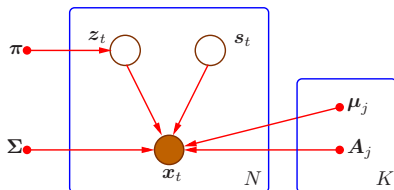
- ▶ E-step (LS projection):

$$\mathbf{S} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{X}.$$

- ▶ M-step (Wiener filter):

$$\mathbf{A} = \mathbf{X} \mathbf{S}^\top (\mathbf{S} \mathbf{S}^\top)^{-1}.$$

Mixture of Factor Analyzers: Graphical Model



- ▶ Gaussian+Multinomial latent variables:

$$p(\mathbf{s}_t) = \mathcal{N}(\mathbf{s}_t | \mathbf{0}, \mathbf{I}),$$

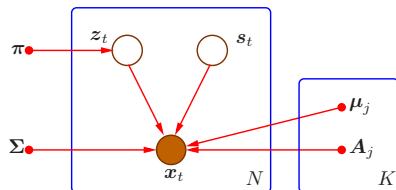
$$p(\mathbf{z}_t) = \prod_{j=1}^K \pi_j^{z_{jt}}.$$

- ▶ The distribution over observed variables conditioned on latent variables is

$$\begin{aligned} p(\mathbf{x}_t | \mathbf{s}_t, \mathbf{z}_t) &= \prod_{j=1}^K \mathcal{N}(\mathbf{x}_t | \mathbf{A}_j \mathbf{s}_t + \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)^{z_{jt}}. \end{aligned}$$

- ▶ Performs **clustering** and **dimensionality reduction** simultaneously.

Learning Mixture of Factor Analyzers



- ▶ Compute maximum likelihood estimates of parameters

$$\{\pi, \{A_j\}, \{\mu_j\}, \{\Sigma\}\}.$$

- ▶ Compute the posterior distribution over latent variables

$$p(s_t | x_t, z_{jt} = 1)$$

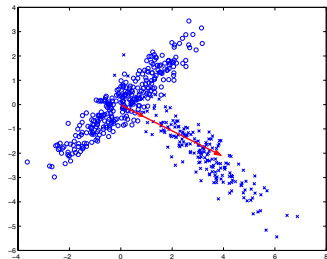
which yields lower-dimensional projection in cluster j .

- ▶ We will do this again using expectation maximization.

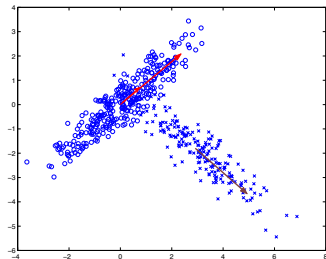
Mixtures of Factor Analyzers (MFA)

- ▶ Clustering and dimensionality reduction are two fundamental problems in unsupervised learning.
 - ▶ The goal of clustering is to group data points by similarity between their features.
 - ▶ The goal of dimensionality reduction is to group (or compress) features that are highly correlated.
- ▶ Mixtures of Factor Analyzers = Dimensionality Reduction (FA) + Clustering (Gaussian Mixture Models).
- ▶ Might be the best method for local dimensionality reduction because it incorporate with probabilistic model.
- ▶ Good applications: character/face recognition

MFA: An Example



(a)



(b)

Figure: (a) PCA; (b) Mixture of factor analyzers.

MFA: Complete-Data Log-Likelihood

- ▶ The complete-data log-likelihood is given by

$$\begin{aligned}\mathcal{L}_c &= \sum_{t=1}^N \log p(\mathbf{x}_t, \mathbf{s}_t, \mathbf{z}_t) \\ &= \sum_{t=1}^N [\log p(\mathbf{x}_t | \mathbf{s}_t, \mathbf{z}_t) + \log p(\mathbf{s}_t) + \log p(\mathbf{z}_t)],\end{aligned}$$

where ω_j represents $z_{jt} = 1$ and $\pi_j = p(\omega_j) = p(z_{jt} = 1)$, and

$$\begin{aligned}p(\mathbf{x}_t | \mathbf{s}_t, \mathbf{z}_t) &= \prod_{j=1}^K \mathcal{N}(\mathbf{x}_t | \mathbf{A}_j \mathbf{s}_t + \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)^{z_{jt}}, \\ p(\mathbf{s}_t) &= \mathcal{N}(\mathbf{s}_t | 0, \mathbf{I}), \quad p(\mathbf{z}_t) = \prod_{j=1}^K \pi_j^{z_{jt}}.\end{aligned}$$

- ▶ The posterior responsibility of the mixture j for generating data point \mathbf{x}_t , is defined by $\langle z_{jt} \rangle = R_{jt} = p(\omega_j | \mathbf{x}_t)$.
- ▶ The calculation of the expected complete-data log-likelihood, involves $\langle z_{jt} \mathbf{s}_t \rangle = R_{jt} \langle \mathbf{s}_t | \omega_j, \mathbf{x}_t \rangle$ and $\langle z_{jt} \mathbf{s}_t \mathbf{s}_t^\top \rangle = R_{jt} \langle \mathbf{s}_t \mathbf{s}_t^\top | \omega_j, \mathbf{x}_t \rangle$.

EM-MFA

► E-Step:

$$\begin{aligned}R_{jt} &= p(\omega_j | \mathbf{x}_t) \propto \pi_j \mathcal{N}(\boldsymbol{\mu}_j, \mathbf{A}_j \mathbf{A}_j^\top + \boldsymbol{\Sigma}), \\ \langle \mathbf{s}_t | \mathbf{x}_t, \omega_j \rangle &= \boldsymbol{\Phi}_j(\mathbf{x}_t - \boldsymbol{\mu}_j), \\ \langle \mathbf{s}_t \mathbf{s}_t^\top | \mathbf{x}_t, \omega_j \rangle &= \mathbf{I} - \boldsymbol{\Phi}_j \mathbf{A}_j + \boldsymbol{\Phi}_j(\mathbf{x}_t - \boldsymbol{\mu}_j)(\mathbf{x}_t - \boldsymbol{\mu}_j)^\top \boldsymbol{\Phi}_j^\top.\end{aligned}$$

► M-Step:

$$\mathbf{A}_j \leftarrow \left[\sum_t R_{jt} (\mathbf{x}_t - \boldsymbol{\mu}_j) \langle \mathbf{s}_t | \mathbf{x}_t, \omega_j \rangle \right] \left[\sum_t R_{jt} \langle \mathbf{s}_t \mathbf{s}_t^\top | \mathbf{x}_t, \omega_j \rangle \right]^{-1},$$

$$\boldsymbol{\mu}_j \leftarrow \frac{\sum_t R_{jt} (\mathbf{x}_t - \mathbf{A}_j \langle \mathbf{s}_t | \mathbf{x}_t, \omega_j \rangle)}{\sum_t R_{jt}},$$

$$\boldsymbol{\Sigma} \leftarrow \frac{1}{N} \text{diag} \left\{ \sum_t \sum_j R_{jt} \left[(\mathbf{x}_t - \boldsymbol{\mu}_j)(\mathbf{x}_t - \boldsymbol{\mu}_j)^\top - \mathbf{A}_j \langle \mathbf{s}_t | \mathbf{x}_t, \omega_j \rangle (\mathbf{x}_t - \boldsymbol{\mu}_j)^\top \right] \right\},$$

$$\pi_j \leftarrow \frac{1}{N} \sum_t R_{jt}.$$

Suggested Further Readings

1. C. Bishop (1999), "Latent variable models," In M. I. Jordan Ed., Learning in Graphical Models.
2. Z. Ghahramani and G. Hinton (1996), "The EM algorithm for mixtures of factor analyzers," University of Toronto Technical Report CRG-TR-96-1.
3. S. Roweis (1997), "EM algorithms for PCA and SPCA," NIPS-1997.
4. M. Tipping and C. Bishop (1999), "Probabilistic principal component analysis," Journal of the Royal Statistical Society, Series B.
5. M. Tipping and C. Bishop (1999), "Mixtures of probabilistic principal component analysers," Neural Computation.