

NOTES

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

495 Advanced Statistical Machine Learning and Pattern Recognition

Author:

Thomas Teh (CID: 0124 3008)

Date: February 10, 2017

1 Expectation Maximization

1.1 General Approach for Expectation Maximization

1.1.1 Notes

1. The goal of the Expectation-Maximization is to find maximum likelihood solutions for models having latent variables
2. The general concept is that since our knowledge of the latent variables in \mathbf{Z} is given by the posterior distribution $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$, we use the expectation of the latent variables instead of the actual values.
3. EM algorithm can be used to find MAP solutions models in which a prior is defined over the parameters.

1.1.2 Algorithm

Given a joint distribution $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ over observed variables \mathbf{X} and latent variables \mathbf{Z} , governed by parameter $\boldsymbol{\theta}$ the goal is to maximize the likelihood function $p(\mathbf{X}|\boldsymbol{\theta})$ w.r.t $\boldsymbol{\theta}$.

1. Choose an initial for the parameters $\boldsymbol{\theta}^{old}$.
2. E Step: Evaluate $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old})$
3. M Step: Evaluate $\boldsymbol{\theta}^{new}$ given by

$$\boldsymbol{\theta}^{new} = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{old})$$

where

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) = \sum_{\mathbf{z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$

4. Check for convergence of either the log-likelihood or the parameter values. If convergence is not satisfied

$$\boldsymbol{\theta}^{old} \leftarrow \boldsymbol{\theta}^{new}$$

and return to step 2

1.2 Gaussian Mixture Models

The Gaussian mixture distribution can be written as a linear superposition of Gaussians

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$



Figure 1: Gaussian Mixture Model.

1.2.1 Formulation of the Gaussian Mixture Models

Let z be a K -dimensional binary random variable with 1-of- K representation.

$$p(z_k = 1) = \pi_k \Rightarrow p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

Conditional probability of \mathbf{x} given a particular value for latent variable \mathbf{z} :

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$$

Using Bayes theorem and marginalize the latent variable \mathbf{z}

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Similarly, the posterior probability of \mathbf{z} is given by

$$\gamma(z_k) = \frac{p(\mathbf{x}|z_k = 1)p(z_k = 1)}{p(\mathbf{x})} = \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

1.2.2 Maximum Likelihood

The log-likelihood function is given by

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)$$

The maximum likelihood method will yield the following:

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})}$$

$$\Sigma_k = \frac{\sum_{n=1}^N \gamma(z_{nk})(\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top}{\sum_{n=1}^N \gamma(z_{nk})}$$

$$\pi_k = \frac{\sum_{n=1}^N \gamma(z_{nk})}{N}$$

Issues with maximum likelihood:

1. Presence of singularities: When we have at least two components in the mixture, one of them can have a finite variance and assign finite probability to all the data points, while the other component can shrink onto one specific data point and therefore contribute to an ever increasing additive value to the log likelihood.
2. Identifiability: Solutions may not be unique, hence it may be hard to interpret the parameter values discovered by a model.
3. The log likelihood equation is difficult to optimize over.

1.2.3 Expectation Maximization Formulation

Supposed that in addition to \mathbf{X} , we were also given the values of the latent variables \mathbf{Z} , the likelihood and the log likelihood function are given by

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \Sigma, \boldsymbol{\pi}) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k)^{z_{nk}}$$

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \Sigma, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k))$$

1. Expectation Step:

Taking the expectation on the log likelihood

$$\mathbb{E}_{p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \Sigma, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}_{p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})} [z_{nk}] (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k))$$

$$= G(\boldsymbol{\theta})$$

$$p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}) \propto \prod_{n=1}^N \prod_{k=1}^K [\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k)]^{z_{nk}}$$

$$\mathbb{E}_{p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})} [z_{nk}] = \frac{\sum_{z_{nk}} z_{nk} [\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k)]^{z_{nk}}}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \Sigma_j)}$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \Sigma_j)}$$

$$= \gamma(z_{nk})$$

2. **Maximization Step:** By taking the derivative of (θ) w.r.t θ and set them to 0, the parameters can be found to be

$$\begin{aligned}\mu_k &= \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} \\ \Sigma_k &= \frac{\sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^\top}{\sum_{n=1}^N \gamma(z_{nk})} \\ \pi_k &= \frac{\sum_{n=1}^N \gamma(z_{nk})}{N}\end{aligned}$$

1.2.4 Expectation Maximization Algorithm

Given a Gaussian mixture model, the goal is to maximize the likelihood functions w.r.t to the parameters:

1. Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k and evaluate the initial of the log likelihood.
2. E Step: Evaluate the responsibilities using the current parameter values

$$\gamma(z_k) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})} = \frac{\pi_k \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)}$$

3. M Step: Re-estimate the parameters using the current responsibilities

$$\begin{aligned}\mu_k^{new} &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \\ \Sigma_k^{new} &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^\top \\ \pi_k^{new} &= \frac{N_k}{N}\end{aligned}$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk})$$

4. Evaluate the log likelihood

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k) \right)$$

and check for convergence of either parameter of the log likelihood. If the convergence criterion is not satisfied, return to step 2

1.3 Bernoulli Mixture Models

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{i=1}^D \mu_i^{x_i} (1 - \mu_i)^{1-x_i}$$

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$$

$$\mathbb{V}[\mathbf{x}] = \text{diag}\{\mu_i(1 - \mu_i)\}$$

The mixture of the Bernoulli distributions is given by

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{k=1}^K \pi_k p(\mathbf{x}|\boldsymbol{\mu}_k)$$

$$p(\mathbf{x}|\boldsymbol{\mu}_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i}$$

$$\mathbb{E}[\mathbf{x}] = \sum_{k=1}^K \pi_k \boldsymbol{\mu}_k$$

$$\mathbb{V}[\mathbf{x}] = \sum_{k=1}^K \pi_k (\boldsymbol{\Sigma}_k + \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top) - \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^\top$$

$$\boldsymbol{\Sigma}_k = \text{diag}\{\mu_{ki}(1 - \mu_{ki})\}$$

Let \mathbf{z} be the one hot representation, we have

$$p(\mathbf{x}|\mathbf{z}, \boldsymbol{\mu}) = \prod_{k=1}^K p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k}$$

$$p(\mathbf{z}|\boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{z_k}$$

The log likelihood function is given by

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left\{ \ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki})] \right\}$$

The E step: we then take the expectation of the log likelihood above

$$\mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left\{ \ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki})] \right\}$$

$$\gamma(z_{nk}) = \frac{\sum_{z_{nk}} z_{nk} [\pi_k p(\mathbf{x}_n|\boldsymbol{\mu}_k)]^{z_{nk}}}{\sum_{z_{nj}} [\pi_j p(\mathbf{x}_n|\boldsymbol{\mu}_j)]^{z_{nj}}} = \frac{\pi_k p(\mathbf{x}_n|\boldsymbol{\mu}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n|\boldsymbol{\mu}_j)}$$

The M step: we then maximize the log likelihood wrt to the parameters

$$\begin{aligned}\mu_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \\ \pi_k &= \frac{N_k}{N} \\ N_k &= \sum_{n=1}^N \gamma(z_{nk})\end{aligned}$$

1.4 Convergence of the EM Algorithm

$$\begin{aligned}p(\mathbf{X}|\boldsymbol{\theta}) &= \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \Rightarrow \ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + KL(q \parallel p) \\ \mathcal{L}(q, \boldsymbol{\theta}) &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left[\frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right] \\ KL(q \parallel p) &= - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left[\frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right]\end{aligned}$$

In the E step, the lower bound $\mathcal{L}(q, \boldsymbol{\theta}^{old})$ is maximized with respect to $q(\mathbf{Z})$. The solution to this maximization problem can be done by setting $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old})$ such that $KL(q \parallel p) = 0$.

In the M step, the lower bound $\mathcal{L}(q, \boldsymbol{\theta}^{old})$ is maximized with respect to $\boldsymbol{\theta}$. However, the KL divergence will no longer be zero. Hence, the increase in the upper bound of the log likelihood function is greater than the increase in the lower bound.

By substituting $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old})$, we get the lower bound to be the following after the E step:

$$\begin{aligned}\mathcal{L}(q, \boldsymbol{\theta}) &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \\ &= \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) + const\end{aligned}$$

2 Probabilistic PCA

$$\begin{aligned} \mathbf{x} &= \mathbf{W}\mathbf{y} + \boldsymbol{\mu} + \mathbf{e} \\ \mathbf{e} &\sim \mathcal{N}(\mathbf{e}|\mathbf{0}, \sigma^2 \mathbf{I}) \\ \mathbf{y} &\sim \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{I}) \end{aligned}$$



Figure 2: Gaussian Mixture Model.

2.1 Maximum Likelihood

1. We are interested in the distribution of \mathbf{x} and the parameters $\boldsymbol{\theta}$. The marginal distribution of \mathbf{x} can be obtained by integrating \mathbf{y} out of the joint distribution.

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_N | \boldsymbol{\theta}) &= \int_{\mathbf{y}_i} p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{y}_1, \dots, \mathbf{y}_N | \boldsymbol{\theta}) d\mathbf{y}_1 \dots d\mathbf{y}_N \\ &= \int_{\mathbf{y}_i} \prod_{i=1}^N p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) p(\mathbf{y}_i) d\mathbf{y}_1 \dots d\mathbf{y}_N \\ &= \prod_{i=1}^N \int_{\mathbf{y}_i} p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) p(\mathbf{y}_i) d\mathbf{y}_i \end{aligned}$$

2. The joint distribution can be obtained by applying Bayes' Rule.

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{y}_1, \dots, \mathbf{y}_N | \boldsymbol{\theta}) &= \prod_{i=1}^N p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) p(\mathbf{y}_i) \\ p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) p(\mathbf{y}_i) &= \left[(2\pi\sigma^2)^{-\frac{F}{2}} e^{-\frac{1}{2\sigma^2} (\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{y}_i)^\top (\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{y}_i)} \right] \left[(2\pi)^{-\frac{d}{2}} e^{-\frac{1}{2} \mathbf{y}_i^\top \mathbf{y}_i} \right] \end{aligned}$$

3. By collecting the \mathbf{y} in the exponents, completing the squares and then applying the Woodbury identity (more details please refer to notes in 496):

$$p(\mathbf{x}_i | \mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}, \mathbf{D})$$

$$\begin{aligned}
\mathbf{D} &= \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I} \\
p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) &= \mathcal{N}(\mathbf{y}_i | \mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}), \sigma^2 \mathbf{M}^{-1}) \\
\mathbf{M} &= \sigma^2 \mathbf{I} + \mathbf{W}^\top \mathbf{W}
\end{aligned}$$

4. Maximizing the likelihood and solve for the parameters

$$\begin{aligned}
\mathbf{S}_t &= \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top \quad (\mathbf{S}_t \text{ is the covariance matrix}) \\
\sigma^2 &= \frac{1}{F-d} \sum_{j=d+1}^F \lambda_j \\
\mathbf{W}_d &= \mathbf{U}_d (\boldsymbol{\Lambda} - \sigma^2 \mathbf{I})^{\frac{1}{2}} \mathbf{V}^\top
\end{aligned}$$

5. Hence we no longer have a projection but:

$$\begin{aligned}
\mathbb{E}_{p(\mathbf{y}_i | \mathbf{x}_i)}[\mathbf{y}_i] &= \mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \\
\widehat{\mathbf{x}}_i &= \mathbf{W} \mathbb{E}_{p(\mathbf{y}_i | \mathbf{x}_i)}[\mathbf{y}_i] + \boldsymbol{\mu}
\end{aligned}$$

2.2 EM PPCA

1. Write the log likelihood of the joint distribution of the observed and latent variables

$$\begin{aligned}
p(\mathbf{X}, \mathbf{Y} | \boldsymbol{\theta}) &= \prod_{i=1}^N p(\mathbf{x}_i | \mathbf{y}_i, \boldsymbol{\theta}) p(\mathbf{y}_i) \\
\ln p(\mathbf{X}, \mathbf{Y} | \boldsymbol{\theta}) &= \sum_{i=1}^N (\ln p(\mathbf{x}_i | \mathbf{y}_i, \boldsymbol{\theta}) + \ln p(\mathbf{y}_i)) \\
\ln p(\mathbf{x}_i | \mathbf{y}_i, \boldsymbol{\theta}) &= -\frac{F}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i - \boldsymbol{\mu})^\top (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i - \boldsymbol{\mu}) \\
\ln p(\mathbf{y}_i) &= -\frac{1}{2} \mathbf{y}_i^\top \mathbf{y}_i - \frac{D}{2} \ln 2\pi
\end{aligned}$$

2. **E-Step:** Take the expectation on the log likelihood on the joint distribution:

$$\ln p(\mathbf{X}, \mathbf{Y} | \boldsymbol{\theta}) = \sum_{i=1}^N \left[-\frac{F}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i - \boldsymbol{\mu})^\top (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i - \boldsymbol{\mu}) - \frac{1}{2} \mathbf{y}_i^\top \mathbf{y}_i - \frac{D}{2} \ln 2\pi \right]$$

Expanding the above and use the identities $\text{tr}[\mathbf{y}_i^\top \mathbf{W} \mathbf{W}^\top \mathbf{y}_i] = \text{tr}[\mathbf{y}_i \mathbf{y}_i^\top \mathbf{W} \mathbf{W}^\top]$

$$\begin{aligned}
\mathbb{E}_{p(\mathbf{Y} | \mathbf{X})}[\ln p(\mathbf{X}, \mathbf{Y} | \boldsymbol{\theta})] &= -\frac{NF}{2} \ln 2\pi\sigma^2 - \frac{ND}{2} \ln 2\pi \\
&\quad - \sum_{i=1}^N \left\{ \frac{1}{2\sigma^2} [(\mathbf{x}_i - \boldsymbol{\mu})^\top (\mathbf{x}_i - \boldsymbol{\mu}) - 2(\mathbf{x}_i - \boldsymbol{\mu})^\top \mathbf{W} \mathbb{E}[\mathbf{y}_i]] \right.
\end{aligned}$$

$$+ \text{tr}[\mathbb{E}(\mathbf{y}_i \mathbf{y}_i^\top) \mathbf{W}^\top \mathbf{W}] + \frac{1}{2} \text{tr}[\mathbb{E}(\mathbf{y}_i^\top \mathbf{y}_i)] \}$$

We can obtain the moments of \mathbf{y}_i from the earlier derivation of the PPCA

$$p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \mathcal{N}(\mathbf{y}_i | \mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}), \sigma^2 \mathbf{M}^{-1})$$

$$\mathbf{M} = \sigma^2 \mathbf{I} + \mathbf{W}^\top \mathbf{W}$$

$$\mathbb{E}_{p(\mathbf{Y}|\mathbf{X})}[\mathbf{y}_i] = \mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu})$$

$$\mathbb{E}_{p(\mathbf{Y}|\mathbf{X})}[\mathbf{y}_i \mathbf{y}_i^\top] = \sigma^2 \mathbf{M}^{-1} + \mathbb{E}[\mathbf{y}_i] \mathbb{E}[\mathbf{y}_i]^\top$$

3. **M-Step:** Maximize the log likelihood

$$\boldsymbol{\mu} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mathbf{W} \mathbb{E}[\mathbf{y}_i])$$

$$\mathbf{W} = \left[\sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}) \mathbb{E}[\mathbf{y}_i]^\top \right] \left[\sum_{i=1}^N \mathbb{E}[\mathbf{y}_i \mathbf{y}_i^\top] \right]^{-1}$$

$$\sigma^2 = \frac{1}{NF} \sum_{i=1}^N \left\{ \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 - 2 \mathbb{E}[\mathbf{y}_i]^\top \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) + \text{tr}(\mathbb{E}[\mathbf{y}_i \mathbf{y}_i^\top] \mathbf{W}^\top \mathbf{W}) \right\}$$

4. Advantages of EM PPCA:

- A complexity of $O(NFD)$ can be significantly smaller than $O(NF^2)$ (from the computation of the covariance)
- The EM procedure can be extended to factor analysis model
- EM allows us to deal with missing values

3 PPCA Mixture Model



Figure 3: PPCA Mixture Model.

The model is given by:

$$\begin{aligned}
\mathbf{x}_i &= \boldsymbol{\mu}_k + \mathbf{W}_k \mathbf{y}_{ik} + \mathbf{e}_{ik} \\
\mathbf{y}_{ik} &\sim \mathcal{N}(\mathbf{y}_{ik} | \mathbf{0}, \mathbf{I}) \\
\mathbf{e}_{ik} &\sim \mathcal{N}(\mathbf{e}_{ik} | \mathbf{0}, \sigma_k^2 \mathbf{I}), \forall k = 1, \dots, K \\
p(\mathbf{x}_i | \mathbf{z}_{ik} = 1, \mathbf{y}_{ik}, \boldsymbol{\mu}_k, \mathbf{W}_k, \sigma_k^2 \mathbf{I}) &= \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k + \mathbf{W}_k \mathbf{y}_{ik}) \\
p(\mathbf{x}_i | \boldsymbol{\theta}) &= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \mathbf{D}_k) \\
\mathbf{D}_k &= \mathbf{W}_k \mathbf{W}_k^\top + \sigma_k^2 \mathbf{I}
\end{aligned}$$

3.1 Probability Distributions for EM Formulation

1. Priors distributions:

$$\begin{aligned}
p(\mathbf{z}_i | \boldsymbol{\theta}_z) &= \prod_{k=1}^K \pi_k^{z_{ik}} \\
p(\mathbf{Y}_i | \mathbf{z}_i, \boldsymbol{\theta}_z) &= \prod_{k=1}^K p(\mathbf{y}_{ik})^{z_{ik}} = \prod_{k=1}^K \mathcal{N}(\mathbf{y}_{ik} | \mathbf{0}, \mathbf{I})^{z_{ik}}
\end{aligned}$$

2. Conditional distribution:

$$\begin{aligned}
p(\mathbf{x}_i | \mathbf{z}_i, \mathbf{Y}_i, \boldsymbol{\theta}_x) &= \prod_{k=1}^K p(\mathbf{x}_i | \mathbf{z}_{ik} = 1, \mathbf{y}_{ik}, \mathbf{W}_k, \boldsymbol{\mu}_k, \sigma_k^2 \mathbf{I})^{z_{ik}} \\
&= \prod_{k=1}^K \mathcal{N}(\mathbf{x}_{ik} | \mathbf{W}_k \mathbf{y}_{ik} + \boldsymbol{\mu}_k, \sigma_k^2 \mathbf{I})^{z_{ik}}
\end{aligned}$$

3. Marginal distributions per cluster is found by integrating \mathbf{y} out:

$$\begin{aligned}
p(\mathbf{x}_i | \mathbf{z}_{ik} = 1, \mathbf{W}_k, \boldsymbol{\mu}_k, \sigma_k^2 \mathbf{I}) &= \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \mathbf{D}_k) \\
\mathbf{D}_k &= \mathbf{W}_k \mathbf{W}_k^\top + \sigma_k^2 \mathbf{I}
\end{aligned}$$

4. Full marginal distributions:

$$\begin{aligned}
p(\mathbf{x}_i | \boldsymbol{\theta}) &= \sum_{k=1}^K p(\mathbf{z}_{ik} = 1) p(\mathbf{x}_i | \mathbf{z}_{ik} = 1, \boldsymbol{\theta}_x) \\
&= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \mathbf{D}_k)
\end{aligned}$$

5. Posteriors on \mathbf{y}_{ik} (we will take expectation on this in the E step):

$$p(\mathbf{y}_{ik} | \mathbf{x}_{ik}, \mathbf{z}_{ik} = 1, \mathbf{W}_k, \boldsymbol{\mu}_k, \sigma_k^2 \mathbf{I}) = \mathcal{N}(\mathbf{y}_{ik} | \mathbf{M}_k^{-1} \mathbf{W}_k^\top (\mathbf{x}_i - \boldsymbol{\mu}_k), \sigma_k^2 \mathbf{M}_k^{-1})$$

$$\begin{aligned}\mathbb{E}(\mathbf{y}_{ik}) &= \mathbf{M}_k^{-1} \mathbf{W}_k^\top (\mathbf{x}_i - \boldsymbol{\mu}_k) \\ \mathbb{E}(\mathbf{y}_{ik} \mathbf{y}_{ik}^\top) &= \sigma_k^2 \mathbf{M}_k^{-1} + \mathbb{E}(\mathbf{y}_{ik}) \mathbb{E}(\mathbf{y}_{ik})^\top\end{aligned}$$

6. Posteriors on \mathbf{z}_i (we will take expectation on this in the E step):

$$\begin{aligned}p(\mathbf{z}_i | \mathbf{x}_i, \theta) &= \frac{p(\mathbf{x}_i, \mathbf{z}_i | \theta)}{p(\mathbf{x}_i | \theta)} = \frac{p(\mathbf{x}_i | \mathbf{z}_i, \theta) p(\mathbf{z}_i | \theta)}{p(\mathbf{x}_i | \theta)} = \frac{\prod_k^K \left(\mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \mathbf{D}_k) \pi_k \right)^{z_{ik}}}{\sum_{l=1}^K \pi_l \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_l, \mathbf{D}_l)} \\ \mathbb{E}[z_{ik}] &= \sum_{z_{ik}=0,1} z_{ik} p(\mathbf{z}_i | \mathbf{x}_i, \theta) = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \mathbf{D}_k)}{\sum_{l=1}^K \pi_l \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_l, \mathbf{D}_l)}\end{aligned}$$

3.2 Formulation of EM

1. Setting the joint likelihood (decomposition below can be guided by graphical model):

$$\begin{aligned}p(\mathbf{X}, \mathbf{Y}, \mathbf{Z} | \theta) &= p(\mathbf{X} | \mathbf{Z}, \mathbf{Y}, \theta_x) p(\mathbf{Z}, \mathbf{Y} | \theta_Z) = p(\mathbf{X} | \mathbf{Z}, \mathbf{Y}, \theta_x) p(\mathbf{Y} | \mathbf{Z}) p(\mathbf{Z} | \theta_Z) \\ p(\mathbf{Z} | \theta_Z) &= \prod_{i=1}^N \prod_{k=1}^K \pi_k^{z_{ik}} \\ p(\mathbf{Y} | \mathbf{Z}) &= \prod_{i=1}^N p(\mathbf{Y}_i | \mathbf{z}_i, \theta_Z) = \prod_{i=1}^N \prod_{k=1}^3 p(\mathbf{y}_{ik})^{z_{ik}} = \prod_{i=1}^N \prod_{k=1}^3 \mathcal{N}(\mathbf{y}_{ik} | \mathbf{0}, \mathbf{I})^{z_{ik}} \\ p(\mathbf{X} | \mathbf{Z}, \mathbf{Y}, \theta_x) &= \prod_{i=1}^N p(\mathbf{x}_i | \mathbf{z}_i, \mathbf{Y}_i, \theta_x) = \prod_{i=1}^N \prod_{k=1}^K \mathcal{N}(\mathbf{x}_i | \mathbf{W}_k \mathbf{y}_{ik} + \boldsymbol{\mu}_k, \sigma_k^2 \mathbf{I})^{z_{ik}}\end{aligned}$$

Hence the log likelihood of the joint distribution:

$$\begin{aligned}\ln p(\mathbf{X}, \mathbf{Z}, \mathbf{Y} | \theta) &= \sum_{i=1}^N \sum_k^K z_{ik} \left[\ln \mathcal{N}(\mathbf{x}_i | \mathbf{W}_k \mathbf{y}_{ik} + \boldsymbol{\mu}_k, \sigma_k^2 \mathbf{I}) + \ln \mathcal{N}(\mathbf{y}_{ik} | \mathbf{0}, \mathbf{I}) + \ln \pi_k \right] \\ &= \sum_{i=1}^N \sum_k^K z_{ik} \left[-\frac{1}{2\sigma_k^2} (\mathbf{x}_i - \boldsymbol{\mu}_k - \mathbf{W}_k \mathbf{y}_{ik})^\top (\mathbf{x}_i - \boldsymbol{\mu}_k - \mathbf{W}_k \mathbf{y}_{ik}) - \frac{F}{2} \ln 2\pi - F \ln \sigma_k \right] \\ &\quad + \sum_{i=1}^N \sum_k^K z_{ik} \left[-\frac{1}{2} \mathbf{y}_{ik}^\top \mathbf{y}_{ik} - \frac{d}{2} \ln 2\pi \right] + \sum_{i=1}^N \sum_k^K z_{ik} \ln \pi_k\end{aligned}$$

2. Taking the expectation yields

$$\begin{aligned}\mathbb{E}[\ln p(\mathbf{X}, \mathbf{Z}, \mathbf{Y} | \theta)] &= -\frac{1}{2\sigma_k^2} \sum_{i=1}^N \sum_k^K \mathbb{E}[z_{ik}] \left[\|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 - \mathbb{E}[\mathbf{y}_{ik}]^\top \mathbf{W}_k^\top (\mathbf{x}_i - \boldsymbol{\mu}_k) + \text{tr}(\mathbf{W}_k \mathbf{W}_k^\top \mathbb{E}[\mathbf{y}_{ik} \mathbf{y}_{ik}^\top]) \right] \\ &\quad + \sum_{i=1}^N \sum_k^K \left[\frac{F}{2} \ln 2\pi - F \ln \sigma_k \right] + \sum_{i=1}^N \sum_k^K \mathbb{E}[z_{ik}] \left[-\frac{1}{2} \text{tr}(\mathbb{E}[\mathbf{y}_{ik} \mathbf{y}_{ik}^\top]) - \frac{d}{2} \ln 2\pi \right] + \sum_{i=1}^N \sum_k^K \mathbb{E}[z_{ik}] \ln \pi_k\end{aligned}$$

3. The maximization step:

$$\begin{aligned}
 \pi_k &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}[z_{ik}] \\
 \boldsymbol{\mu}_k &= \frac{\sum_{i=1}^N \mathbb{E}[z_{ik}](\mathbf{x}_i - \mathbf{W}_k \mathbb{E}[\mathbf{y}_{ik}])}{\sum_{i=1}^N \mathbb{E}[z_{ik}]} \\
 \mathbf{W}_k &= \left[\sum_{i=1}^N \mathbb{E}[z_{ik}](\mathbf{x}_i - \boldsymbol{\mu}_k) \mathbb{E}[\mathbf{y}_{ik}^\top] \right] \left[\sum_{i=1}^N \mathbb{E}[z_{ik}] \mathbb{E}[\mathbf{y}_{ik} \mathbf{y}_{ik}^\top] \right]^{-1} \\
 \sigma_k^2 &= \frac{1}{F \sum_{i=1}^N \mathbb{E}[z_{ik}]} \sum_{i=1}^N \mathbb{E}[z_{ik}] \left[\|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 - \mathbb{E}[\mathbf{y}_{ik}]^\top \mathbf{W}_k^\top (\mathbf{x}_i - \boldsymbol{\mu}_k) + \text{tr}(\mathbf{W}_k^\top \mathbf{W}_k \mathbb{E}[\mathbf{y}_{ik} \mathbf{y}_{ik}^\top]) \right]
 \end{aligned}$$

4 PPCA with Missing Values

Given the PPCA model:

$$\begin{aligned} \mathbf{x} &= \mathbf{W}\mathbf{y} + \boldsymbol{\mu} + \mathbf{e} \\ \mathbf{e} &\sim \mathcal{N}(\mathbf{e}|\mathbf{0}, \sigma^2 \mathbf{I}) \\ \mathbf{y} &\sim \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{I}) \end{aligned}$$

$$\begin{aligned} p(\mathbf{x}_i|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) &= \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}, \mathbf{D}) \\ p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) &= \mathcal{N}(\mathbf{y}_i|\mathbf{M}^{-1}\mathbf{W}^\top(\mathbf{x}_i - \boldsymbol{\mu}), \sigma^2 \mathbf{M}^{-1}) \\ \mathbf{D} &= \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I} \\ \mathbf{M} &= \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I} \end{aligned}$$

1. We can reformulate the PPCA model to take into account of missing data:

$$\begin{aligned} \mathbf{x} = \begin{bmatrix} \mathbf{x}^o \\ \mathbf{x}^u \end{bmatrix} &\Rightarrow p(\mathbf{x}^o, \mathbf{x}^u) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x}^o \\ \mathbf{x}^u \end{bmatrix} \middle| \begin{bmatrix} \boldsymbol{\mu}^o \\ \boldsymbol{\mu}^u \end{bmatrix}, \begin{bmatrix} \mathbf{D}^{oo} & \mathbf{D}^{ou} \\ \mathbf{D}^{uo} & \mathbf{D}^{uu} \end{bmatrix}\right) \\ p(\mathbf{x}^u|\mathbf{x}^o) &= \mathcal{N}\left(\mathbf{x}^u \middle| \boldsymbol{\mu}^u + \mathbf{D}_{uo}\mathbf{D}_{oo}^{-1}(\mathbf{x}^o - \boldsymbol{\mu}^o), \mathbf{D}_{uu} - \mathbf{D}_{uo}\mathbf{D}_{oo}^{-1}\mathbf{D}_{ou}\right) \end{aligned}$$

2. For convenience, we can rewrite this formulation as derive its first order and second order moments:

$$\begin{aligned} p(\mathbf{x}_i|\mathbf{x}_i^o) &= \mathcal{N}(\mathbf{x}_i|\mathbf{z}_i, \mathbf{Q}) \\ \mathbf{z}_i &= \begin{bmatrix} \mathbf{x}_i^o \\ \boldsymbol{\mu}^u + \mathbf{D}_{uo}\mathbf{D}_{oo}^{-1}(\mathbf{x}_i^o - \boldsymbol{\mu}^o) \end{bmatrix} \\ \mathbf{Q} &= \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{D}_{uu} - \mathbf{D}_{uo}\mathbf{D}_{oo}^{-1}\mathbf{D}_{ou} \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{p(\mathbf{x}_i|\mathbf{x}_i^o)}[\mathbf{x}_i] &= \mathbf{z}_i \\ \mathbb{E}_{p(\mathbf{x}_i|\mathbf{x}_i^o)}[\mathbf{x}_i \mathbf{x}_i^\top] &= \mathbf{Q} + \mathbf{z}_i \mathbf{z}_i^\top \\ \mathbb{E}_{p(\mathbf{x}_i|\mathbf{x}_i^o)}[(\mathbf{x}_i - \mathbf{a})(\mathbf{x}_i - \mathbf{a})^\top] &= \mathbf{Q} + (\mathbf{z}_i - \mathbf{a})(\mathbf{z}_i - \mathbf{a})^\top \end{aligned}$$

3. Writing down the log likelihood for the observed variable \mathbf{X} and the latent variable \mathbf{Y}

$$\begin{aligned} \ln p(\mathbf{X}, \mathbf{Y}|\boldsymbol{\theta}) &= -\frac{1}{2\sigma^2} \sum_{i=1}^N \left\{ \text{tr}[(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top] - 2\text{tr}[\mathbf{y}_i(\mathbf{x}_i - \boldsymbol{\mu})^\top \mathbf{W}] + \text{tr}[\mathbf{W}^\top \mathbf{W} \mathbf{y}_i \mathbf{y}_i^\top] \right\} \\ &\quad - \frac{1}{2} \sum_{i=1}^N \text{tr}[\mathbf{y}_i \mathbf{y}_i^\top] - \frac{NF}{2} \ln 2\pi - NF \ln \sigma - N \ln 2\pi \end{aligned}$$

4. Take the expectation with regards to the probability distribution $p(\mathbf{X}, \mathbf{Y} | \mathbf{X}^o)$

$$\begin{aligned} & \mathbb{E}_{p(\mathbf{x}_i, \mathbf{y}_i | \mathbf{x}_i^o)} [\ln p(\mathbf{X}, \mathbf{Y} | \boldsymbol{\theta})] \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^N \left\{ \text{tr}(\mathbb{E}[(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top]) - 2\text{tr}(\mathbb{E}[\mathbf{y}_i(\mathbf{x}_i - \boldsymbol{\mu})^\top] \mathbf{W}) + \text{tr}(\mathbf{W}^\top \mathbf{W} \mathbb{E}[\mathbf{y}_i \mathbf{y}_i^\top]) \right\} \\ & - \frac{1}{2} \sum_{i=1}^N \text{tr}(\mathbb{E}[\mathbf{y}_i \mathbf{y}_i^\top]) - \frac{NF}{2} \ln 2\pi - NF \ln \sigma - N \ln 2\pi \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{p(\mathbf{x}_i, \mathbf{y}_i | \mathbf{x}_i^o)} [(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top] &= \mathbb{E}_{p(\mathbf{x}_i, \mathbf{y}_i | \mathbf{x}_i^o)} [(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top] \\ &= \mathbf{Q} + (\mathbf{z}_i - \boldsymbol{\mu})(\mathbf{z}_i - \boldsymbol{\mu})^\top \\ \mathbb{E}_{p(\mathbf{x}_i, \mathbf{y}_i | \mathbf{x}_i^o)} [\mathbf{y}_i(\mathbf{x}_i - \boldsymbol{\mu})^\top] &= \mathbb{E}_{p(\mathbf{x}_i | \mathbf{x}_i^o)} \left[\mathbb{E}_{p(\mathbf{y}_i | \mathbf{x}_i^o)} [\mathbf{y}_i(\mathbf{x}_i - \boldsymbol{\mu})^\top] \right] \\ &= \mathbb{E}_{p(\mathbf{x}_i | \mathbf{x}_i^o)} \left[\mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top \right] \\ &= \mathbf{M}^{-1} \mathbf{W}^\top \mathbb{E}_{p(\mathbf{x}_i | \mathbf{x}_i^o)} [(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top] \\ \mathbb{E}_{p(\mathbf{x}_i, \mathbf{y}_i | \mathbf{x}_i^o)} [\mathbf{y}_i \mathbf{y}_i^\top] &= \mathbb{E}_{p(\mathbf{x}_i | \mathbf{x}_i^o)} \left[\mathbb{E}_{p(\mathbf{y}_i | \mathbf{x}_i^o)} [\mathbf{y}_i \mathbf{y}_i^\top] \right] \\ &= \mathbb{E}_{p(\mathbf{x}_i | \mathbf{x}_i^o)} \left[\sigma^2 \mathbf{M}^{-1} + \mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top \mathbf{W} \mathbf{M}^{-1} \right] \\ &= \sigma^2 \mathbf{M}^{-1} + \mathbf{M}^{-1} \mathbf{W}^\top \mathbb{E}_{p(\mathbf{x}_i | \mathbf{x}_i^o)} [(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top] \mathbf{W} \mathbf{M}^{-1} \end{aligned}$$

5. Maximization step:

$$\begin{aligned} \boldsymbol{\mu} &= \frac{1}{N} \sum_{i=1}^N \mathbf{z}_i \\ \mathbf{W} &= \left[\sum_{i=1}^N \mathbb{E}[(\mathbf{x}_i - \boldsymbol{\mu}) \mathbf{y}_i^\top] \right] \left[\sum_{i=1}^N \mathbb{E}[\mathbf{y}_i \mathbf{y}_i^\top] \right]^{-1} \\ \sigma^2 &= \frac{1}{NF} \sum_{i=1}^N \left\{ \text{tr}(\mathbb{E}[(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top]) - 2\text{tr}(\mathbb{E}[\mathbf{y}_i(\mathbf{x}_i - \boldsymbol{\mu})^\top] \mathbf{W}) + \text{tr}(\mathbb{E}[\mathbf{y}_i \mathbf{y}_i^\top] \mathbf{W}^\top \mathbf{W}) \right\} \\ \mathbf{D} &= \mathbf{W} \mathbf{W}^\top + \sigma^2 \mathbf{I} \end{aligned}$$

5 Hidden Markov Model

5.1 Preliminary: Markov Chains

1. Markov Property:

$$p(\mathbf{x}_i | \mathbf{x}_1, \dots, \mathbf{x}_{i-1}) = p(\mathbf{x}_i | \mathbf{x}_{i-1})$$

2. Markov Model: Given the observations $\{\mathbf{x}_t\}_{t=1}^N$

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = p(\mathbf{x}_1) \prod_{t=2}^N p(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad \text{bigram model}$$

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = p(\mathbf{x}_1) p(\mathbf{x}_2 | \mathbf{x}_1) \prod_{t=3}^N p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_{t-2}) \quad \text{trigram model}$$

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = p(\mathbf{x}_1) p(\mathbf{x}_2 | \mathbf{x}_1) p(\mathbf{x}_3 | \mathbf{x}_1, \mathbf{x}_2) \prod_{t=4}^N p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_{t-2}, \mathbf{x}_{t-3}) \quad \text{n-gram model}$$

3. A transition matrix A specifies the probabilities of getting from state i (row) to state j (column) in one step. Note that A is a stochastic matrix, i.e. the row must sum to 1.
4. Stationary distribution is the long term distribution over the states

$$\boldsymbol{\pi}^\top = \boldsymbol{\pi}^\top A$$

Solving for the stationary distribution is the same as eigenanalysis, where $\boldsymbol{\pi}$ is an eigenvector with eigenvalue 1.

$$\begin{aligned} \boldsymbol{\pi}^\top (I - A) &= \mathbf{0} \\ \boldsymbol{\pi}^\top \mathbf{1} &= 1 \end{aligned}$$

5. For a stationary distribution to exist, the markov chain has to be irreducible and aperiodic.
 - Irreducibility: The state transition diagram must be singly connected, i.e. it is possible to move from one state to another ($a_{ij}(t) > 0$)
 - Aperiodicity: $d(i) = \gcd\{t : a_{ii}(t) > 0\} = 1, \forall i$. At some point in time, the probability of a recurrent connection is greater than 0.
6. **Detailed balance condition:** A markov process is called a reversible Markov process if it satisfies the detailed balance equations.

$$\pi_i P_{ij} = \pi_j P_{ji}$$

7. Estimation of the transition matrix from training data (language model)

(a) The probability of a character is given by

$$p(\mathbf{x}_1|\boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{x_{1k}}$$

$$p(\mathbf{x}_t|\mathbf{x}_{t-1}) = \prod_{j=1}^K \prod_{k=1}^K a_{jk}^{x_{(t-1)j}x_{tk}}$$

(b) Probability of a sequence D_l with length T

$$P(D_l|\theta) = p(\mathbf{x}_1^l, \dots, \mathbf{x}_T^l) = p(\mathbf{x}_1^l) \prod_{t=2}^T p(\mathbf{x}_t^l|\mathbf{x}_{t-1}^l) = \prod_{k=1}^K \pi_k^{x_{1k}^l} \prod_{t=2}^T \prod_{j=1}^K \prod_{k=1}^K a_{jk}^{x_{(t-1)j}^l x_{tk}^l}$$

(c) The likelihood function is then

$$p(D_1, \dots, D_N|\theta) = \prod_{l=1}^N p(D_l|\theta) = \prod_{l=1}^N \left\{ \prod_{k=1}^K \pi_k^{x_{1k}^l} \prod_{t=2}^T \prod_{j=1}^K \prod_{k=1}^K a_{jk}^{x_{(t-1)j}^l x_{tk}^l} \right\}$$

$$\begin{aligned} \ln p(D_1, \dots, D_N|\theta) &= \sum_{l=1}^N \sum_{k=1}^K x_{1k}^l \ln \pi_k + \sum_{l=1}^N \sum_{t=2}^T \sum_{j=1}^K \sum_{k=1}^K x_{(t-1)j}^l x_{tk}^l \ln a_{jk} \\ &= \sum_{k=1}^K \left(\sum_{l=1}^N x_{1k}^l \right) \ln \pi_k + \sum_{j=1}^K \sum_{k=1}^K \left(\sum_{l=1}^N \sum_{t=2}^T x_{(t-1)j}^l x_{tk}^l \right) \ln a_{jk} \\ &= \sum_{k=1}^K N_k^1 \ln \pi_k + \sum_{j=1}^K \sum_{k=1}^K N_{jk} \ln a_{jk} \end{aligned}$$

where

$$N_k^1 = \sum_{l=1}^N x_{1k}^l \quad N_{jk} = \sum_{l=1}^N \sum_{t=2}^T x_{(t-1)j}^l x_{tk}^l$$

(d) We just need to solve the optimization below by formulating the Lagrangian:

$$\begin{aligned} \max_{\boldsymbol{\theta}} \quad & \sum_{k=1}^K N_k^1 \ln \pi_k + \sum_{j=1}^K \sum_{k=1}^K N_{jk} \ln a_{jk} \\ \text{s.t.} \quad & \sum_{k=1}^K \pi_k = 1 \end{aligned}$$

$$\sum_{k=1}^K a_{jk} = 1$$

$$\pi_k = \frac{N_k^1}{\sum_{k=1}^K N_k^1} \quad a_{jk} = \frac{N_{jk}}{\sum_{k=1}^K N_{jk}}$$

5.2 Hidden Markov Model

5.2.1 Model Description and Important Properties

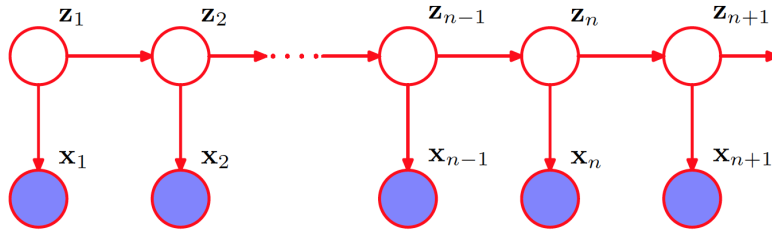


Figure 4: State Space Model. If z are discrete and markovian, then it is a Hidden Markov model.

1. The HMM is a specific instance of the state space model in Figure 4 in which the latent variables are discrete.
2. If we view a single time slice of the model, it corresponds to a mixture distribution densities given by $p(x|z)$, where $p(z_n|z_{n-1})$.
3. Properties of conditional independence

$$\begin{aligned} p(X|z_n) &= p(x_1, \dots, x_n|z_n)p(x_{n+1}, \dots, x_N|z_n) \\ p(x_1, \dots, x_{n-1}|x_n, z_n) &= p(x_1, \dots, x_{n-1}|z_n) \\ p(x_1, \dots, x_{n-1}|z_{n-1}, z_n) &= p(x_1, \dots, x_{n-1}|z_{n-1}) \\ p(x_{n+1}, \dots, x_N|z_n, z_{n+1}) &= p(x_{n+1}, \dots, x_N|z_{n+1}) \\ p(x_{n+2}, \dots, x_N|z_{n+1}, x_{n+1}) &= p(x_{n+2}, \dots, x_N|z_{n+1}) \\ p(X|z_{n-1}, z_n) &= p(x_1, \dots, x_{n-1}|z_{n-1})p(x_n|z_n)p(x_{n+1}, \dots, x_N|z_n) \\ p(x_{N+1}|X, z_{N+1}) &= p(x_{N+1}|z_{N+1}) \\ p(z_{N+1}|z_N, X) &= p(z_{N+1}|z_N) \end{aligned}$$

4. Inferences using the Hidden Markov Model

- **Filtering:** Compute the belief state $p(z_n|x_1, \dots, x_n)$ online as the data streams in. It is known as filtering to reduce the noise in the estimation of the hidden state.

- **Smoothing:** Compute $p(z_t|x_1, \dots, x_N)$ offline, given all the evidence. Reducing the uncertainty by conditioning on the past and future data.
- **Fixed lag smoothing:** Compute $p(z_{t-\ell}|x_1, \dots, x_n)$, where $\ell > 0$ is the lag.
- **MAP estimation:** Compute $\arg \max_{z_1, \dots, z_n} p(z_1, \dots, z_n|x_1, \dots, x_n)$. This is known as Viterbi Decoding.
- **Evaluation:** Find the probability of the evidence $p(x_1, \dots, x_N) = \sum_z p(x, z)$.
- **Prediction:** Predicting the future given the past, we compute $p(z_{n+\delta}|x_1, \dots, x_n)$ and $p(x_{n+\delta}|x_1, \dots, x_n)$.
- **Parameter Estimation:** This is known as the Baum-Welch algorithm. We estimate the parameter A , π and θ .

5.2.2 Filtering and Smoothing: Forwards and Backwards Algorithm

1. We can write the posterior probabilities of the latent variables as below. After that, we will formulate the computation for $\alpha(z_n)$ (forward probabilities) and $\beta(z_n)$ (backward probabilities).

$$\begin{aligned}
 \gamma(z_n) &= p(z_n|X) = \frac{p(X|z_n)p(z_n)}{p(X)} \\
 &= \frac{p(x_1, \dots, x_n|z_n)p(x_{n+1}, \dots, x_N|z_n)p(z_n)}{p(X)} \\
 &= \frac{p(x_1, \dots, x_n, z_n)p(x_{n+1}, \dots, x_N|z_n)}{p(X)} \\
 &= \frac{\alpha(z_n)\beta(z_n)}{p(X)}
 \end{aligned}$$

2. The forward probabilities can be computed recursively based on the below:

$$\begin{aligned}
 \alpha(z_n) &= p(x_1, \dots, x_n, z_n) \\
 &= p(x_1, \dots, x_n|z_n)p(z_n) \\
 &= p(x_n|z_n)p(x_1, \dots, x_{n-1}|z_n)p(z_n) \\
 &= p(x_n|z_n)p(x_1, \dots, x_{n-1}, z_n) \\
 &= p(x_n|z_n) \sum_{z_{n-1}} p(x_1, \dots, x_{n-1}, z_n|z_{n-1})p(z_{n-1}) \\
 &= p(x_n|z_n) \sum_{z_{n-1}} p(x_1, \dots, x_{n-1}|z_{n-1})p(z_n|z_{n-1})p(z_{n-1}) \\
 &= p(x_n|z_n) \sum_{z_{n-1}} p(x_1, \dots, x_{n-1}, z_{n-1})p(z_n|z_{n-1}) \\
 &= p(x_n|z_n) \sum_{z_{n-1}} \alpha(z_{n-1})p(z_n|z_{n-1})
 \end{aligned}$$

The initial condition of the recursive formula is given by:

$$\alpha(z_1) = p(x_1, z_1) = p(z_1)p(x_1|z_1) = \prod_{k=1}^K \{\pi_k p(x_1|z_{1k})\}^{z_{1k}}$$

$$p(x_1, \dots, x_n) = \sum_{z_n} p(x_1, \dots, x_n, z_n) = \sum_{z_n} \alpha(z_n)$$

$$p(z_n|x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n, z_n)}{p(x_1, \dots, x_n)} = \frac{\alpha(z_n)}{\sum_{z_n} \alpha(z_n)} = \tilde{\alpha}(z_n)$$

3. The backward probabilities can be computed as follows:

$$\begin{aligned} \beta(z_n) &= p(x_{n+1}, \dots, x_N | z_n) \\ &= \sum_{z_{n+1}} p(x_{n+1}, \dots, x_N, z_{n+1} | z_n) \\ &= \sum_{z_{n+1}} p(x_{n+1}, \dots, x_N | z_n, z_{n+1}) p(z_{n+1} | z_n) \\ &= \sum_{z_{n+1}} p(x_{n+1}, \dots, x_N | z_{n+1}) p(z_{n+1} | z_n) \\ &= \sum_{z_{n+1}} p(x_{n+2}, \dots, x_N | z_{n+1}) p(x_{n+1} | z_{n+1}) p(z_{n+1} | z_n) \\ &= \sum_{z_{n+1}} \beta(z_{n+1}) p(x_{n+1} | z_{n+1}) p(z_{n+1} | z_n) \end{aligned}$$

5.2.3 Prediction

5.2.4 Baum-Welch Algorithm

5.2.5 Viterbi Algorithm

6 Useful Identities:

- Woodbury Identity:

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)VA^{-1}$$

- Law of Total Expectation:

$$\mathbb{E}_{p(X)}[X] = \mathbb{E}_{p(Y)}[\mathbb{E}_{p(X|Y)}[X|Y]]$$

- Conditional and Margin of Block Distributions Given the distribution:

$$p(x, y) = \mathcal{N}\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}\right)$$

We have

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\boldsymbol{\mu}_{x|y}, \boldsymbol{\Sigma}_{x|y}) \\ \boldsymbol{\mu}_{x|y} &= \boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} (\mathbf{y} - \boldsymbol{\mu}_y) \\ \boldsymbol{\Sigma}_{x|y} &= \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx} \end{aligned}$$

- Product of two Gaussians (note that the underlying random variable must be the same):

$$\mathcal{N}(\mathbf{x}|\mathbf{a}, \mathbf{A}) \mathcal{N}(\mathbf{x}|\mathbf{b}, \mathbf{B}) = c \mathcal{N}(\mathbf{x}|\mathbf{c}, \mathbf{C})$$

$$\begin{aligned} \mathbf{C} &= (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} \\ \mathbf{c} &= \mathbf{C}(\mathbf{A}^{-1} \mathbf{a} + \mathbf{B}^{-1} \mathbf{b}) \\ c &= \mathcal{N}(\mathbf{a}|\mathbf{b}, \mathbf{A} + \mathbf{B}) = \mathcal{N}(\mathbf{b}|\mathbf{a}, \mathbf{A} + \mathbf{B}) \end{aligned}$$

- Derivative of inverse

$$\begin{aligned} \frac{\partial \mathbf{Y}^{-1}}{\partial x} &= -\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial x} \mathbf{Y}^{-1} \\ \frac{\partial \mathbf{a}^\top \mathbf{X}^{-1} \mathbf{B}}{\partial \mathbf{X}} &= -\mathbf{X}^{-\top} \mathbf{a} \mathbf{b}^\top \mathbf{X}^{-\top} \end{aligned}$$

- Derivative of determinant

$$\frac{\partial |\mathbf{X}|}{\partial \mathbf{X}} = |\mathbf{X}| \mathbf{X}^{-\top}$$

$$\mathbf{a}^\top \mathbf{a} = \text{tr}(\mathbf{a} \mathbf{a}^\top)$$