

NOTES

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

495 Advanced Statistical Machine Learning and Pattern Recognition

Author:

Thomas Teh (CID: 0124 3008)

Date: March 9, 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} :

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

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

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

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

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

1.2.2 Maximum Likelihood

The log-likelihood function is given by

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

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

$$\begin{aligned}
 p(\mathbf{x}|\boldsymbol{\mu}) &= \prod_{i=1}^D \mu_i^{x_i} (1 - \mu_i)^{1-x_i} \\
 \mathbb{E}[\mathbf{x}] &= \boldsymbol{\mu} \\
 \mathbb{V}[\mathbf{x}] &= \text{diag}\{\mu_i(1 - \mu_i)\}
 \end{aligned}$$

The mixture of the Bernoulli distributions is given by

$$\begin{aligned}
 p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\pi}) &= \sum_{k=1}^K \pi_k p(\mathbf{x}|\boldsymbol{\mu}_k) \\
 p(\mathbf{x}|\boldsymbol{\mu}_k) &= \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i} \\
 \mathbb{E}[\mathbf{x}] &= \sum_{k=1}^K \pi_k \boldsymbol{\mu}_k \\
 \mathbb{V}[\mathbf{x}] &= \sum_{k=1}^K \pi_k (\boldsymbol{\Sigma}_k + \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top) - \mathbb{E}[\mathbf{X}] \mathbb{E}[\mathbf{X}]^\top \\
 \boldsymbol{\Sigma}_k &= \text{diag}\{\mu_{ki}(1 - \mu_{ki})\}
 \end{aligned}$$

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

$$\begin{aligned}
 p(\mathbf{x}|\mathbf{z}, \boldsymbol{\mu}) &= \prod_{k=1}^K p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k} \\
 p(\mathbf{z}|\boldsymbol{\pi}) &= \prod_{k=1}^K \pi_k^{z_k}
 \end{aligned}$$

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

$$\begin{aligned}
 \mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\pi})] &= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left\{ \ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki})] \right\} \\
 \gamma(z_{nk}) &= \frac{\sum_{z_{nk}} z_{nk} [\pi_k p(\mathbf{x}_n|\boldsymbol{\mu}_k)]^{z_{nk}}}{\sum_{z_{nj}} [\pi_j p(\mathbf{x}_n|\boldsymbol{\mu}_j)]^{z_{nj}}} = \frac{\pi_k p(\mathbf{x}_n|\boldsymbol{\mu}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n|\boldsymbol{\mu}_j)}
 \end{aligned}$$

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:

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

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

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

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

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

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

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

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

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

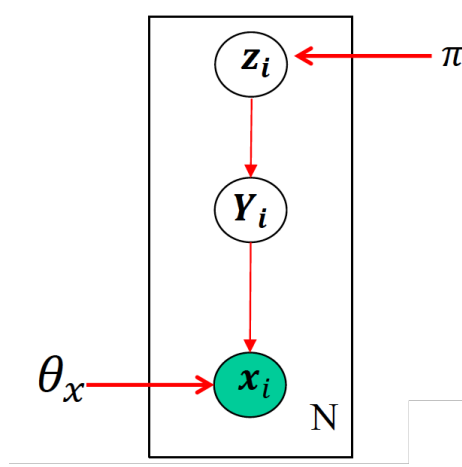


Figure 3: PPCA Mixture Model.

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:

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

$$= \prod_{k=1}^K \mathcal{N}(\mathbf{x}_{ik} | \mathbf{W}_k \mathbf{y}_{ik} + \boldsymbol{\mu}_k, \sigma_k^2 \mathbf{I})^{z_{ik}}$$

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

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

$$\mathbf{D}_k = \mathbf{W}_k \mathbf{W}_k^\top + \sigma_k^2 \mathbf{I}$$

4. Full marginal distributions:

$$p(\mathbf{x}_i | \boldsymbol{\theta}) = \sum_{k=1}^K p(\mathbf{z}_{ik} = 1) p(\mathbf{x}_i | \mathbf{z}_{ik} = 1, \boldsymbol{\theta}_x)$$

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

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

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

$$\mathbb{E}(\mathbf{y}_{ik}) = \mathbf{M}_k^{-1} \mathbf{W}_k^\top (\mathbf{x}_i - \boldsymbol{\mu}_k)$$

$$\mathbb{E}(\mathbf{y}_{ik} \mathbf{y}_{ik}^\top) = \sigma_k^2 \mathbf{M}_k^{-1} + \mathbb{E}(\mathbf{y}_{ik}) \mathbb{E}(\mathbf{y}_{ik})^\top$$

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

$$p(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}) = \frac{p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})}{p(\mathbf{x}_i | \boldsymbol{\theta})} = \frac{p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) p(\mathbf{z}_i | \boldsymbol{\theta})}{p(\mathbf{x}_i | \boldsymbol{\theta})} = \frac{\prod_k^K (\mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \mathbf{D}_k) \pi_k)^{z_{ik}}}{\sum_{l=1}^K \pi_l \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_l, \mathbf{D}_l)}$$

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

3.2 Formulation of EM

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

$$p(\mathbf{X}, \mathbf{Y}, \mathbf{Z} | \boldsymbol{\theta}) = p(\mathbf{X} | \mathbf{Z}, \mathbf{Y}, \boldsymbol{\theta}_x) p(\mathbf{Z}, \mathbf{Y} | \boldsymbol{\theta}_Z) = p(\mathbf{X} | \mathbf{Z}, \mathbf{Y}, \boldsymbol{\theta}_x) p(\mathbf{Y} | \mathbf{Z}) p(\mathbf{Z} | \boldsymbol{\theta}_Z)$$

$$p(\mathbf{Z} | \boldsymbol{\theta}_Z) = \prod_{i=1}^N \prod_{k=1}^K \pi_k^{z_{ik}}$$

$$p(\mathbf{Y} | \mathbf{Z}) = \prod_{i=1}^N p(\mathbf{Y}_i | \mathbf{z}_i, \boldsymbol{\theta}_Z) = \prod_{i=1}^N \prod_{k=1}^3 p(\mathbf{y}_{ik})^{z_{ik}} = \prod_{i=1}^N \prod_{k=1}^3 \mathcal{N}(\mathbf{y}_{ik} | \mathbf{0}, \mathbf{I})^{z_{ik}}$$

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

Hence the log likelihood of the joint distribution:

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

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

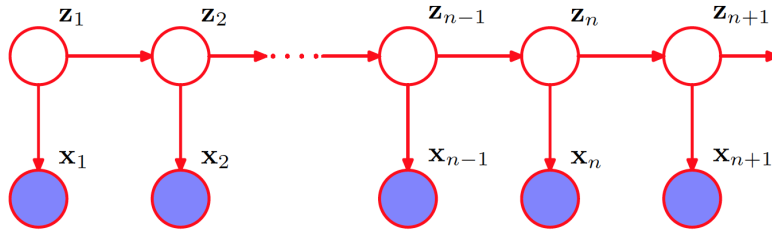


Figure 4: State Space Model. If z are discrete and markovian, then it is a Hidden Markov model.

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 Formulation (Discrete Case)

Let there be K possible states for the latent variable z and S possible states for the observed variable x .

$$\begin{aligned}
 p(z_1 | \pi) &= \prod_{c=1}^K \pi_c^{z_{1,c}} \\
 p(z_n | z_{n-1}, A) &= \prod_{i=1}^K \prod_{j=1}^K a_{ij}^{z_{n-1,i} z_{n,j}} \\
 p(x_n | z_n) &= \prod_{i=1}^S \prod_{j=1}^K b_{ij}^{x_{n,i} z_{n,j}}
 \end{aligned}$$

If x is continuous, we just replace the last equation as

$$p(x_n | z_n) = \prod_{j=1}^K p(x_n | z_{n,j})^{z_{n,j}}$$

Given the M observations of sequence D with size N

$$\begin{aligned}
 p(D, Z | \theta) &= \prod_{l=1}^M p(x_1^l, \dots, x_N^l, z_1^l, \dots, z_N^l | \theta) \\
 &= \prod_{l=1}^M \left(\prod_{n=1}^N p(x_n^l | z_n^l) p(z_1^l) \prod_{n=2}^N p(z_n^l | z_{n-1}^l) \right) \\
 &= \prod_{l=1}^M \left(\prod_{n=1}^N \prod_{i=1}^S \prod_{j=1}^K b_{ij}^{x_{n,i}^l z_{n,j}^l} \prod_{c=1}^K \pi_c^{z_{1,c}^l} \prod_{n=2}^N \prod_{i=1}^K \prod_{j=1}^K a_{ij}^{z_{n-1,i}^l z_{n,j}^l} \right)
 \end{aligned}$$

Taking the log of the likelihood function

$$L = \sum_{l=1}^M \sum_{n=1}^N \sum_{i=1}^S \sum_{j=1}^K x_{n,i}^l z_{n,j}^l \ln b_{ij} + \sum_{l=1}^M \sum_{c=1}^K z_{1,c}^l \ln \pi_c + \sum_{l=1}^M \sum_{n=2}^N \sum_{i=1}^S \sum_{j=1}^K z_{n-1,i}^l z_{n,j}^l \ln a_{ij}$$

Taking expectation of the likelihood function

$$\mathbb{E}[L] = \sum_{l=1}^M \sum_{n=1}^N \sum_{i=1}^S \sum_{j=1}^K x_{n,i}^l \mathbb{E}[z_{n,j}^l] \ln b_{ij} + \sum_{l=1}^M \sum_{c=1}^K \mathbb{E}[z_{1,c}^l] \ln \pi_c + \sum_{l=1}^M \sum_{n=2}^N \sum_{i=1}^S \sum_{j=1}^K \mathbb{E}[z_{n-1,i}^l z_{n,j}^l] \ln a_{ij}$$

$$\begin{aligned} \mathbb{E}[z_{1,c}^l] &= \sum_{z_{1,c}^l} z_{1,c}^l p(z_{1,c}^l | \mathbf{x}_1^l, \dots, \mathbf{x}_T^l) = p(z_{1,c}^l = 1 | \mathbf{x}_1^l, \dots, \mathbf{x}_T^l) \\ &= \frac{p(\mathbf{x}_1^l, z_{1,c}^l = 1) p(\mathbf{x}_2^l, \dots, \mathbf{x}_N^l | z_{1,c}^l = 1)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_N^l)} = \frac{\alpha(z_{1,c}^l) \beta(z_{1,c}^l)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_N^l)} \end{aligned}$$

$$\begin{aligned} \mathbb{E}[z_{n,j}^l] &= \sum_{z_{n,j}^l} z_{n,j}^l p(z_{n,j}^l | \mathbf{x}_1^l, \dots, \mathbf{x}_N^l) = p(z_{n,j}^l = 1 | \mathbf{x}_1^l, \dots, \mathbf{x}_N^l) \\ &= \frac{p(\mathbf{x}_1^l, \dots, \mathbf{x}_n^l, z_{n,j}^l = 1) p(\mathbf{x}_{n+1}^l, \dots, \mathbf{x}_N^l | z_{n,j}^l = 1)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_N^l)} = \frac{\alpha(z_{n,j}^l) \beta(z_{n,j}^l)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_N^l)} \end{aligned}$$

$$\begin{aligned} \mathbb{E}[z_{n-1,i}^l z_{n,j}^l] &= \sum_{z_{n-1,i}^l} \sum_{z_{n,j}^l} z_{n-1,i}^l z_{n,j}^l p(z_{n-1,i}^l z_{n,j}^l | \mathbf{x}_1^l, \dots, \mathbf{x}_N^l) = p(z_{n-1,i}^l = 1, z_{n,j}^l = 1 | \mathbf{x}_1^l, \dots, \mathbf{x}_N^l) \\ &= \frac{\alpha(z_{n-1,i}^l) \prod_{r=1}^S b_{j,r}^{x_{n,r}^l} a_{ij} \beta(z_{n,j}^l)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_N^l)} \end{aligned}$$

Formulating the Lagrangian and solving for the parameters we get

$$\begin{aligned} \mathcal{L} &= \sum_{l=1}^M \sum_{n=1}^N \sum_{i=1}^S \sum_{j=1}^K x_{n,i}^l \mathbb{E}[z_{n,j}^l] \ln b_{ij} + \sum_{l=1}^M \sum_{c=1}^K \mathbb{E}[z_{1,c}^l] \ln \pi_c + \sum_{l=1}^M \sum_{n=2}^N \sum_{i=1}^S \sum_{j=1}^K \mathbb{E}[z_{n-1,i}^l z_{n,j}^l] \ln a_{ij} \\ &\quad + \left(\sum_{j=1}^K \lambda_j \sum_{i=1}^S b_{ij} - 1 \right) + \mu \left(\sum_{c=1}^K \pi_c - 1 \right) + \phi \left(\sum_{k=1}^K a_{jk} - 1 \right) \end{aligned}$$

Taking the derivative w.r.t. to the parameters, we get

$$\begin{aligned} b_{ij} &= \frac{\sum_{l=1}^M \sum_{n=1}^N \mathbb{E}[z_{n,j}^l] x_{n,i}^l}{\sum_{l=1}^M \sum_{n=1}^N \mathbb{E}[z_{n,j}^l]} \\ \pi_c &= \frac{\sum_{l=1}^M \mathbb{E}[z_{1,c}^l]}{\sum_{l=1}^M \sum_{r=1}^K \mathbb{E}[z_{1,r}^l]} \\ a_{ij} &= \frac{\sum_{l=1}^M \sum_{n=2}^N \mathbb{E}[z_{n-1,i}^l z_{n,j}^l]}{\sum_{r=1}^K \sum_{l=1}^M \sum_{n=2}^N \mathbb{E}[z_{n-1,i}^l z_{n,r}^l]} \end{aligned}$$

5.2.3 Formulation (Continuous Case)

Let there be K possible states for the latent variable \mathbf{z} and $p(\mathbf{x}_n|\mathbf{z}_n)$ is a normal distribution.

$$p(\mathbf{z}_1|\boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{z_{1,k}}$$

$$p(\mathbf{z}_n|\mathbf{z}_{n-1}, \mathbf{A}) = \prod_{i=1}^K \prod_{j=1}^K a_{ij}^{z_{n-1,i} z_{n,j}}$$

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

Given the M observations of sequence D with size N

$$\begin{aligned} p(D, \mathbf{Z}|\boldsymbol{\theta}) &= \prod_{l=1}^M p(\mathbf{x}_1^l, \dots, \mathbf{x}_N^l, \mathbf{z}_1^l, \dots, \mathbf{z}_N^l|\boldsymbol{\theta}) \\ &= \prod_{l=1}^M \left(\prod_{n=1}^N p(\mathbf{x}_n^l|\mathbf{z}_n^l) p(\mathbf{z}_1^l) \prod_{n=2}^N p(\mathbf{z}_n^l|\mathbf{z}_{n-1}^l) \right) \\ &= \prod_{l=1}^M \left(\prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_{n,k}^l} \prod_{k=1}^K \pi_k^{z_{1,k}^l} \prod_{n=2}^N \prod_{i=1}^K \prod_{j=1}^K a_{ij}^{z_{n-1,i}^l z_{n,j}^l} \right) \end{aligned}$$

Taking log of the likelihood,

$$\begin{aligned} L &= \sum_{l=1}^M \sum_{n=1}^N \sum_{k=1}^K z_{n,k}^l \ln \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \sum_{l=1}^M \sum_{k=1}^K z_{1,k}^l \ln \pi_k + \sum_{l=1}^M \sum_{n=2}^N \sum_{i=1}^K \sum_{j=1}^K z_{n-1,i}^l z_{n,j}^l \ln a_{ij} \\ \mathbb{E}[L] &= \sum_{l=1}^M \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}[z_{n,k}^l] \ln \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \sum_{l=1}^M \sum_{k=1}^K \mathbb{E}[z_{1,k}^l] \ln \pi_k + \sum_{l=1}^M \sum_{n=2}^N \sum_{i=1}^K \sum_{j=1}^K \mathbb{E}[z_{n-1,i}^l z_{n,j}^l] \ln a_{ij} \end{aligned}$$

$$\begin{aligned} \mathbb{E}[z_{1,c}^l] &= \frac{p(\mathbf{x}_1^l, z_{1,c}^l = 1) p(\mathbf{x}_2^l, \dots, \mathbf{x}_N^l | z_{1,c}^l = 1)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_N^l)} = \frac{\alpha(z_{1,c}^l) \beta(z_{1,c}^l)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_N^l)} \\ \mathbb{E}[z_{n,j}^l] &= \frac{p(\mathbf{x}_1^l, \dots, \mathbf{x}_n^l, z_{n,j}^l = 1) p(\mathbf{x}_{n+1}^l, \dots, \mathbf{x}_N^l | z_{n,j}^l = 1)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_N^l)} = \frac{\alpha(z_{n,j}^l) \beta(z_{n,j}^l)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_N^l)} \\ \mathbb{E}[z_{n-1,i}^l z_{n,j}^l] &= \frac{\alpha(z_{n-1,i}^l) \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) a_{ij} \beta(z_{n,j}^l)}{p(\mathbf{x}_1^l, \dots, \mathbf{x}_N^l)} \end{aligned}$$

Taking the derivative w.r.t. to the parameters, we get

$$\boldsymbol{\mu}_k = \frac{\sum_{l=1}^M \sum_{n=1}^N \mathbb{E}[z_{n,k}^l] \mathbf{x}_n^l}{\sum_{l=1}^M \sum_{n=1}^N \mathbb{E}[z_{n,k}^l]}$$

$$\begin{aligned}\Sigma_k &= \frac{\sum_{l=1}^M \sum_{n=1}^N \mathbb{E}[z_{n,k}^l] (\mathbf{x}_n^l - \boldsymbol{\mu}_k)(\mathbf{x}_n^l - \boldsymbol{\mu}_k)^\top}{\sum_{l=1}^M \sum_{n=1}^N \mathbb{E}[z_{n,k}^l]} \\ \pi_k &= \frac{\sum_{l=1}^M \mathbb{E}[z_{1,k}^l]}{\sum_{l=1}^M \sum_{r=1}^K \mathbb{E}[z_{1,r}^l]} \\ a_{ij} &= \frac{\sum_{l=1}^M \sum_{n=2}^N \mathbb{E}[z_{n-1,i}^l z_{n,j}^l]}{\sum_{r=1}^K \sum_{l=1}^M \sum_{n=2}^N \mathbb{E}[z_{n-1,i}^l z_{n,r}^l]}\end{aligned}$$

5.2.4 Inference

1. Filtering and Smoothing: Forwards and Backwards Algorithm

- (a) We can write the posterior probabilities of the latent variables as below. After that, we will formulate the computation for $\alpha(\mathbf{z}_n)$ (forward probabilities) and $\beta(\mathbf{z}_n)$ (backward probabilities).

$$\begin{aligned}\gamma(\mathbf{z}_n) &= p(\mathbf{z}_n | \mathbf{X}) = \frac{p(\mathbf{X} | \mathbf{z}_n) p(\mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n) p(\mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{\alpha(\mathbf{z}_n) \beta(\mathbf{z}_n)}{p(\mathbf{X})}\end{aligned}$$

- (b) The forward probabilities can be computed recursively based on the below:

$$\begin{aligned}\alpha(\mathbf{z}_n) &= p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n) \\ &= p(\mathbf{x}_1, \dots, \mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n) \\ &= p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_n) p(\mathbf{z}_n) \\ &= p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_n) \\ &= p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{z}_{n-1}) \\ &= p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{z}_{n-1}) \\ &= p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1}) \\ &= p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \alpha(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})\end{aligned}$$

The initial condition of the recursive formula is given by:

$$\begin{aligned}\alpha(z_1) &= p(\mathbf{x}_1, z_1) = p(z_1)p(\mathbf{x}_1|z_1) = \prod_{k=1}^K \{\pi_k p(\mathbf{x}_1|z_{1k})\}^{z_{1k}} \\ p(\mathbf{x}_1, \dots, \mathbf{x}_n) &= \sum_{z_n} p(\mathbf{x}_1, \dots, \mathbf{x}_n, z_n) = \sum_{z_n} \alpha(z_n) \\ p(z_n|\mathbf{x}_1, \dots, \mathbf{x}_n) &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, z_n)}{p(\mathbf{x}_1, \dots, \mathbf{x}_n)} = \frac{\alpha(z_n)}{\sum_{z_n} \alpha(z_n)} = \tilde{\alpha}(z_n)\end{aligned}$$

(c) The backward probabilities can be computed as follows:

$$\begin{aligned}\beta(z_n) &= p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | z_n) \\ &= \sum_{z_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N, z_{n+1} | z_n) \\ &= \sum_{z_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | z_n, z_{n+1}) p(z_{n+1} | z_n) \\ &= \sum_{z_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | z_{n+1}) p(z_{n+1} | z_n) \\ &= \sum_{z_{n+1}} p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | z_{n+1}) p(\mathbf{x}_{n+1} | z_{n+1}) p(z_{n+1} | z_n) \\ &= \sum_{z_{n+1}} \beta(z_{n+1}) p(\mathbf{x}_{n+1} | z_{n+1}) p(z_{n+1} | z_n)\end{aligned}$$

Initial condition for the backward probabilities is given by:

$$p(z_N | \mathbf{x}_1, \dots, \mathbf{x}_N) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_N, z_N)}{p(\mathbf{x}_1, \dots, \mathbf{x}_N)} = \frac{\alpha(z_N)}{p(\mathbf{x}_1, \dots, \mathbf{x}_N)} \Rightarrow \beta(z_N) = 1$$

2. Prediction

$$\begin{aligned}p(\mathbf{x}_{n+2} | \mathbf{x}_1, \dots, \mathbf{x}_n) &= \sum_{z_{n+2}} p(\mathbf{x}_{n+2} | z_{n+2}) p(z_{n+2} | \mathbf{x}_1, \dots, \mathbf{x}_n) \\ p(z_{n+2} | \mathbf{x}_1, \dots, \mathbf{x}_n) &= \sum_{z_{n+1}} \sum_{z_n} p(z_{n+2}, z_{n+1}, z_n | \mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= \sum_{z_{n+1}} \sum_{z_n} p(z_{n+2}, z_{n+1} | \mathbf{x}_1, \dots, \mathbf{x}_n, z_n) p(z_n | \mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= \sum_{z_{n+1}} \sum_{z_n} p(z_{n+2}, z_{n+1} | z_n) \tilde{\alpha}(z_n) \\ &= \sum_{z_{n+1}} \sum_{z_n} p(z_{n+2} | z_{n+1}) p(z_{n+1} | z_n) \tilde{\alpha}(z_n) \\ &= \mathbf{A}^\top \mathbf{A}^\top \tilde{\alpha}\end{aligned}$$

3. Smooth Transition

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

4. Baum-Welch Algorithm / EM

- (a) Initialization: We set all the parameters to random values. For A or π , if any of their elements are initialized to 0, they will remain to be 0.
- (b) Compute both the smoothed and filtered posterior probability recursively
- (c) Compute the expectation $\mathbb{E}[z_{n,k}^l]$ for all M number of sequences.
- (d) Compute the smooth transition probability

5. Scaling Issue

$$\begin{aligned}
p(x_1, \dots, x_n) &= \prod_{k=1}^n c_k \\
c_k &= p(x_k | x_1, \dots, x_{k-1})
\end{aligned}$$

We can then define $\hat{\alpha}(z_n)$ and $\hat{\beta}(z_n)$

$$\begin{aligned}
\alpha(z_n) &= p(z_n | x_1, \dots, x_n) p(x_1, \dots, x_n) = \hat{\alpha}(z_n) \prod_{k=1}^n c_k \\
\beta(z_n) &= \hat{\beta}(z_n) \prod_{k=n+1}^N c_k
\end{aligned}$$

The recursive formulae to use are then:

$$\begin{aligned}
c_n \hat{\alpha}(z_n) &= p(x_n | z_n) \sum_{z_{n-1}} \hat{\alpha}(z_{n-1}) p(z_n | z_{n-1}) \\
c_{n+1} \hat{\beta}(z_n) &= \sum_{z_{n+1}} \hat{\beta}(z_{n+1}) p(x_{n+1} | z_{n+1}) p(z_{n+1} | z_n)
\end{aligned}$$

Hence, the posteriors become:

$$\begin{aligned}
p(z_n | x_1, \dots, x_N) &= \hat{\alpha}(z_n) \hat{\beta}(z_n) \\
p(z_{n-1}, z_n | x_1, \dots, x_N) &= c_n^{-1} \hat{\alpha}(z_{n-1}) p(x_n | z_n) p(z_n | z_{n-1}) \hat{\beta}(z_n)
\end{aligned}$$

6. Viterbi Algorithm

We define the probability of ending in a state j at time n , given that we take the most probable path

$$\begin{aligned}\delta_n(j) &= \max_{z_1, \dots, z_{n-1}} p(z_1, \dots, z_{n-1}, z_n, j = 1 | x_1, \dots, x_n) \\ &= c_n^{-1} p(x_n | z_n, j = 1) \max_i \delta_{n-1}(i) p(z_n, j = 1 | z_{n-1}, i = 1) \\ a_n(j) &= \arg \max_i p(x_n | z_n, j = 1) \delta_{n-1}(i) p(z_n, j = 1 | z_{n-1}, i = 1)\end{aligned}$$

6 Linear Dynamical System (Kalman Filter)

6.1 Definition of a Linear Dynamical System

The model and the transition model are given as follows

$$\begin{aligned} \mathbf{x}_n &= \mathbf{C}\mathbf{z}_n + \mathbf{v}_n \\ \mathbf{z}_n &= \mathbf{A}\mathbf{z}_{n-1} + \mathbf{w}_n \\ \mathbf{z}_1 &= \boldsymbol{\mu}_0 + \mathbf{u} \end{aligned}$$

The distributions of the noise are given by

$$\begin{aligned} \mathbf{v} &\sim \mathcal{N}(\mathbf{v}|\mathbf{0}, \Sigma) \\ \mathbf{w} &\sim \mathcal{N}(\mathbf{w}|\mathbf{0}, \Gamma) \\ \mathbf{u} &\sim \mathcal{N}(\mathbf{u}|\mathbf{0}, \mathbf{V}_0) \end{aligned}$$

The emission probability and the transition probability is given by

$$\begin{aligned} p(\mathbf{z}_n|\mathbf{z}_{n-1}) &= \mathcal{N}(\mathbf{z}_n|\mathbf{A}\mathbf{z}_{n-1}, \Gamma) \\ p(\mathbf{x}_n|\mathbf{z}_n) &= \mathcal{N}(\mathbf{x}_n|\mathbf{C}\mathbf{z}_n, \Sigma) \\ p(\mathbf{z}_1) &= \mathcal{N}(\mathbf{z}_1|\boldsymbol{\mu}_0, \mathbf{V}_0) \end{aligned}$$

6.2 Inference for LDS

6.2.1 Filtering

We can denote the probability $p(\mathbf{z}_n|\mathbf{x}_{1:n})$ as below:

$$\hat{\alpha}(\mathbf{z}_n) = p(\mathbf{z}_n|\mathbf{x}_{1:n}) = \mathcal{N}(\mathbf{z}_n|\boldsymbol{\mu}_n, \mathbf{V}_n)$$

Deriving a recursive relation:

$$\begin{aligned} c_n \hat{\alpha}(\mathbf{z}_n) &= p(\mathbf{x}_n|\mathbf{z}_n) \int_{\mathbf{z}_{n-1}} \hat{\alpha}(\mathbf{z}_{n-1}) p(\mathbf{z}_n|\mathbf{z}_{n-1}) d\mathbf{z}_{n-1} \\ c_n \mathcal{N}(\mathbf{z}_n|\boldsymbol{\mu}_n, \mathbf{V}_n) &= \mathcal{N}(\mathbf{x}_n|\mathbf{C}\mathbf{z}_n, \Sigma) \int_{\mathbf{z}_{n-1}} \mathcal{N}(\mathbf{z}_n|\mathbf{A}\mathbf{z}_{n-1}, \Gamma) \mathcal{N}(\mathbf{z}_{n-1}|\boldsymbol{\mu}_{n-1}, \mathbf{V}_{n-1}) d\mathbf{z}_{n-1} \end{aligned}$$

We can then either complete the square or use the identity of Gaussian multiplication for the integral:

$$\int_{\mathbf{z}_{n-1}} \mathcal{N}(\mathbf{z}_n|\mathbf{A}\mathbf{z}_{n-1}, \Gamma) \mathcal{N}(\mathbf{z}_{n-1}|\boldsymbol{\mu}_{n-1}, \mathbf{V}_{n-1}) d\mathbf{z}_{n-1} = \mathcal{N}(\mathbf{z}_n|\mathbf{A}\boldsymbol{\mu}_{n-1}, \mathbf{P}_{n-1})$$

where,

$$\mathbf{P}_{n-1} = \mathbf{A}\mathbf{V}_{n-1}\mathbf{A}^\top + \Gamma$$

Hence, we have

$$\begin{aligned} c_n \hat{\alpha}(\mathbf{z}_n) &= \mathcal{N}(\mathbf{x}_n | \mathbf{C} \mathbf{z}_n, \Sigma) \mathcal{N}(\mathbf{z}_n | \mathbf{A} \boldsymbol{\mu}_{n-1}, \mathbf{P}_{n-1}) \\ c_n \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}_n, \mathbf{V}_n) &= \mathcal{N}(\mathbf{x}_n | \mathbf{C} \mathbf{z}_n, \Sigma) \mathcal{N}(\mathbf{z}_n | \mathbf{A} \boldsymbol{\mu}_{n-1}, \mathbf{P}_{n-1}) \end{aligned}$$

We can then use the product of Gaussians identity and the Woodbury identity to derive the update rule

$$\begin{aligned} \boldsymbol{\mu}_n &= \mathbf{A} \boldsymbol{\mu}_{n-1} + \mathbf{K}_n (\mathbf{x}_n - \mathbf{C} \mathbf{A} \boldsymbol{\mu}_{n-1}) \\ \mathbf{V}_n &= (\mathbf{I} - \mathbf{K}_n \mathbf{C}) \mathbf{P}_{n-1} \\ c_n &= \mathcal{N}(\mathbf{x}_n | \mathbf{C} \mathbf{A} \boldsymbol{\mu}_{n-1}, \mathbf{C} \mathbf{P}_{n-1} \mathbf{C}^\top + \Sigma) \\ \mathbf{K}_n &= \mathbf{P}_{n-1} \mathbf{C}^\top (\mathbf{C} \mathbf{P}_{n-1} \mathbf{C}^\top + \Sigma)^{-1} \end{aligned}$$

The initial condition for the recursion is given by

$$\begin{aligned} c_1 \hat{\alpha}(\mathbf{z}_1) &= p(\mathbf{z}_1) p(\mathbf{x}_1 | \mathbf{z}_1) \\ \boldsymbol{\mu}_1 &= \boldsymbol{\mu}_0 + \mathbf{K}_1 (\mathbf{x}_1 - \mathbf{C} \boldsymbol{\mu}_0) \\ c_1 &= \mathcal{N}(\mathbf{x}_1 | \mathbf{C} \boldsymbol{\mu}_0, \mathbf{C} \mathbf{V}_0 \mathbf{C}^\top + \Sigma) \\ \mathbf{K}_1 &= \mathbf{V}_0 \mathbf{C}^\top (\mathbf{C} \mathbf{V}_0 \mathbf{C}^\top + \Sigma)^{-1} \end{aligned}$$

6.2.2 Smoothing

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_{1:T}) = \hat{\alpha}(\mathbf{z}_n) \hat{\beta}(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n | \hat{\boldsymbol{\mu}}_n, \hat{\mathbf{V}}_n)$$

$$c_{n+1} \hat{\beta}(\mathbf{z}_n) = \int_{\mathbf{z}_{n+1}} \hat{\beta}(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n) d\mathbf{z}_{n+1}$$

$$\begin{aligned} \hat{\boldsymbol{\mu}} &= \boldsymbol{\mu}_n + \mathbf{J}_n (\hat{\boldsymbol{\mu}}_{n+1} - \mathbf{A} \boldsymbol{\mu}_n) \\ \hat{\mathbf{V}}_n &= \mathbf{V}_n + \mathbf{J}_n (\hat{\mathbf{V}}_{n+1} - \mathbf{P}_n) \mathbf{J}_n^\top \\ \mathbf{J}_n &= \mathbf{V}_n \mathbf{A}^\top \mathbf{P}_n^{-1} \end{aligned}$$

$$\begin{aligned} \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) &= P(\mathbf{z}_n, \mathbf{z}_{n-1} | \mathbf{x}_{1:T}) = c_n^{-1} \hat{\alpha}(\mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) P(\mathbf{z}_n | \mathbf{z}_{n-1}) \hat{\beta}(\mathbf{z}_n) \\ &= \frac{\mathcal{N}(\mathbf{z}_{n-1} | \boldsymbol{\mu}_{n-1}, \mathbf{V}_{n-1}) \mathcal{N}(\mathbf{z}_n | \mathbf{A} \mathbf{z}_{n-1}, \Gamma) \mathcal{N}(\mathbf{x}_n | \mathbf{C} \mathbf{z}_n, \Sigma) \mathcal{N}(\mathbf{z}_n | \hat{\boldsymbol{\mu}}_n, \hat{\mathbf{V}}_n)}{c_n \hat{\alpha}(\mathbf{z}_n)} \\ &= \mathcal{N} \left(\begin{bmatrix} \mathbf{z}_{n-1} \\ \mathbf{z}_n \end{bmatrix} \middle| \begin{bmatrix} \hat{\boldsymbol{\mu}}_{n-1} \\ \hat{\boldsymbol{\mu}}_n \end{bmatrix}, \begin{bmatrix} \hat{\mathbf{V}}_{n-1} & \mathbf{J}_{n-1}^\top \hat{\mathbf{V}}_n \\ (\mathbf{J}_{n-1} \hat{\mathbf{V}}_n)^\top & \hat{\mathbf{V}}_n \end{bmatrix} \right) \end{aligned}$$

6.3 EM Algorithm

First, we can compute the expectations that we will need in the EM algorithm

$$\begin{aligned}\mathbb{E}[\mathbf{z}_n] &= \hat{\boldsymbol{\mu}}_n \\ \mathbb{E}[\mathbf{z}_n \mathbf{z}_{n-1}^\top] &= \mathbf{J}_{n-1} \hat{\mathbf{V}}_n + \hat{\boldsymbol{\mu}}_n \hat{\boldsymbol{\mu}}_{n-1}^\top \\ \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] &= \hat{\mathbf{V}}_n + \hat{\boldsymbol{\mu}}_n \hat{\boldsymbol{\mu}}_n^\top\end{aligned}$$

Deriving the log joint likelihood

$$\begin{aligned}p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) &= p(\mathbf{x}_{1:N}, \mathbf{z}_{1:N} | \boldsymbol{\theta}) \\ &= \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_1) \prod_{n=2}^N p(\mathbf{z}_n | \mathbf{z}_{n-1})\end{aligned}$$

6.3.1 Expectation Step

$$\begin{aligned}\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) &= \ln p(\mathbf{z}_1 | \boldsymbol{\mu}_0, \mathbf{V}_0) + \sum_{n=2}^N \ln p(\mathbf{z}_n | \mathbf{z}_{n-1}, \mathbf{A}, \boldsymbol{\Gamma}) + \sum_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{C}, \boldsymbol{\Sigma}) \\ \mathbb{E}[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})] &= \mathbb{E}[\ln p(\mathbf{z}_1 | \boldsymbol{\mu}_0, \mathbf{V}_0)] + \mathbb{E}\left[\sum_{n=2}^N \ln p(\mathbf{z}_n | \mathbf{z}_{n-1}, \mathbf{A}, \boldsymbol{\Gamma})\right] + \mathbb{E}\left[\sum_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{C}, \boldsymbol{\Sigma})\right]\end{aligned}$$

6.3.2 Maximization Step

To find $\boldsymbol{\mu}_0$ and \mathbf{V}_0 :

$$\begin{aligned}\mathbb{E}[\ln p(\mathbf{z}_1 | \boldsymbol{\mu}_0, \mathbf{V}_0)] &= -\frac{1}{2} \ln |\mathbf{V}_0| - \frac{1}{2} \mathbb{E}_{\mathbf{Z} | \boldsymbol{\theta}^{old}} \left[(\mathbf{z}_1 - \boldsymbol{\mu}_0)^\top \mathbf{V}_0^{-1} (\mathbf{z}_1 - \boldsymbol{\mu}_0) \right] + \text{const.} \\ \boldsymbol{\mu}_0^{new} &= \mathbb{E}[\mathbf{z}_1] \\ \mathbf{V}_0^{new} &= \mathbb{E}[\mathbf{z}_1 \mathbf{z}_1^\top] - \mathbb{E}[\mathbf{z}_1] \mathbb{E}[\mathbf{z}_1]^\top\end{aligned}$$

To find \mathbf{A} and $\boldsymbol{\Gamma}$:

$$\begin{aligned}\mathbb{E}\left[\sum_{n=2}^N \ln p(\mathbf{z}_n | \mathbf{z}_{n-1}, \mathbf{A}, \boldsymbol{\Gamma})\right] &= -\frac{N-1}{2} \ln |\boldsymbol{\Gamma}| - \frac{1}{2} \mathbb{E}_{\mathbf{Z} | \boldsymbol{\theta}^{old}} \left[\sum_{n=2}^N (\mathbf{z}_n - \mathbf{A} \mathbf{z}_{n-1})^\top \boldsymbol{\Gamma}^{-1} (\mathbf{z}_n - \mathbf{A} \mathbf{z}_{n-1}) \right] \\ \mathbf{A}^{new} &= \left(\sum_{n=2}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_{n-1}^\top] \right) \left(\sum_{n=2}^N \mathbb{E}[\mathbf{z}_{n-1} \mathbf{z}_{n-1}^\top] \right)^{-1} \\ \boldsymbol{\Gamma}^{new} &= \frac{1}{N-1} \sum_{n=2}^N \left\{ \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] - \mathbf{A}^{new} \mathbb{E}[\mathbf{z}_{n-1} \mathbf{z}_n^\top] \right\} \\ &\quad - \frac{1}{N-1} \sum_{n=2}^N \left\{ \mathbb{E}[\mathbf{z}_n \mathbf{z}_{n-1}^\top] \mathbf{A}^{new} - \mathbf{A}^{new} \mathbb{E}[\mathbf{z}_{n-1} \mathbf{z}_{n-1}^\top] \mathbf{A}^{new} \right\}\end{aligned}$$

To find \mathbf{C} and Σ

$$\begin{aligned}\mathbb{E}\left[\sum_{n=1}^N p(\mathbf{x}_n|\mathbf{z}_n, \mathbf{C}, \Sigma)\right] &= -\frac{N}{2} \ln |\Sigma| - \frac{1}{2} \mathbb{E}_{\mathbf{Z}|\theta^{old}} \left[\sum_{n=1}^N (\mathbf{x}_n - \mathbf{C}\mathbf{z}_n)^\top \Sigma^{-1} (\mathbf{x}_n - \mathbf{C}\mathbf{z}_n) \right] + \text{const.} \\ \mathbf{C}^{new} &= \left(\sum_{n=1}^N \mathbf{x}_n \mathbb{E}[\mathbf{z}_n^\top] \right) \left(\sum_{n=1}^N \mathbf{x}_n \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] \right)^{-1} \\ \Sigma^{new} &= \frac{1}{N} \sum_{n=1}^N \left\{ \mathbf{x}_n \mathbf{x}_n^\top - \mathbf{C}^{new} \mathbb{E}[\mathbf{z}_n] \mathbf{x}_n^\top - \mathbf{x}_n \mathbb{E}[\mathbf{z}_n^\top] \mathbf{C}^{new} + \mathbf{C}^{new} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] \mathbf{C}^{new} \right\}\end{aligned}$$

7 Markov Random Fields

7.1 Graph Preliminaries

1. A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a finite collection of nodes $\mathcal{V} = \{n_1, \dots, n_N\}$ and set of edges $\mathcal{E} \subset \binom{\mathcal{V}}{2}$
2. Neighbor: Two nodes $n_i, n_j \in \mathcal{V}$ are neighbors $\Leftrightarrow (n_i, n_j) \in \mathcal{E}$. The neighbor of a node is denoted by $\mathcal{N}(n_i) = \{n_j : (n_i, n_j) \in \mathcal{E}\}$
3. Neighbor is a symmetric relation (undirected graph): $n_i \in \mathcal{N}(n_j) \Leftrightarrow n_j \in \mathcal{N}(n_i)$
4. Complete graph: $\forall n_i \in \mathcal{V}, \mathcal{N}(n_i) = \{(n_i, n_j), j = \{1, \dots, N\} \setminus \{i\}\}$
5. A clique is a complete subgraph of \mathcal{G} . Maximal clique is a clique with maximal number of nodes, i.e. cannot add any other node while still retaining complete connectedness.

7.2 Definition and Properties

1. Markov Random Field: Undirected graphical model in which each node corresponds to a random variable or a collection of random variables, and the edges identify conditional dependencies
2. Let A, B, C be three disjoint subsets of \mathcal{V} . Let x_A denote the collection of random variables in A . Conditional independence is defined as

$$p(x_A, x_B | x_C) = p(x_A | x_C) p(x_B | x_C)$$

3. Pairwise Markovianity: $(n_i, n_j) \notin \mathcal{E} \Rightarrow x_i$ and x_j are independent when conditioned on all other variables

$$p(x_i, x_j | \mathbf{x}_{\setminus \{i, j\}}) = p(x_i | \mathbf{x}_{\setminus \{i, j\}}) p(x_j | \mathbf{x}_{\setminus \{i, j\}})$$

4. Local Markovianity: Given its neighborhood, a variable is independent on the rest of the variables

$$p(x_i | \mathbf{x}_{\mathcal{V} \setminus \{i\}}) = p(x_i | \mathbf{x}_{\mathcal{N}(i)})$$

5. Global Markovianity: Let A, B, C be three disjoint subsets of \mathcal{V} . If C separates A from B , $\Rightarrow p(x_A, x_B | x_C) = p(x_A | x_C) p(x_B | x_C)$, then $p(x_A, x_B | x_C)$ is global markov w.r.t to graph \mathcal{G} .
6. Consider a random field \mathbf{x} on a graph \mathcal{G} , s.t. $p(\mathbf{x} > 0)$. Let \mathcal{C} denote the set of all maximal cliques of the graph.

- If the field has the local Markov property, then $p(\mathbf{x})$ can be written as a Gibbs distribution

$$p(\mathbf{x}) = \frac{1}{Z} \prod_C \Psi_C(\mathbf{x}_C) = \frac{1}{Z} \exp \left(- \sum_C V_C(\mathbf{x}_C) \right)$$

$$\Psi_C(\mathbf{x}_C) = \exp(-V_C(\mathbf{x}_C)) \text{ (potential functions)}$$

$$Z = \sum_{\mathbf{x}} \prod_C \Psi_C(\mathbf{x}_C) \text{ (partition functions)}$$

If $p(\mathbf{x})$ can be written as Gibbs distribution for the cliques of some graph, then it has the global Markov property.

7.3 The Ising Model

1. The Ising model considers an idealized system of interacting particles, arranged onto a regular planar grid.
2. Each particle can have only two possible states and it interacts only with its nearest neighbors.
3. The contribution of each particle to the total energy of the system depends upon the orientation of its spin compared to its neighbors.
4. If \mathbf{x}_i denotes the possible states of the particle, assuming no external force, the total energy of the system is given by

$$E(\mathbf{x}_1, \dots, \mathbf{x}_N) = -\frac{\beta}{2} \sum_{k=1}^N \sum_{j \in N(k)} \mathbf{x}_k^\top \mathbf{x}_j$$

The probability of the configuration is given by

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = \frac{1}{Z(\beta)} \exp(-E(\mathbf{x}_1, \dots, \mathbf{x}_N))$$

5. If there is external force

$$E(\mathbf{x}_1, \dots, \mathbf{x}_N) = -\frac{\beta}{2} \sum_{k=1}^N \sum_{j \in N(k)} \mathbf{x}_k^\top \mathbf{x}_j - \lambda \sum_{k=1}^N \mathbf{x}_k^\top \mathbf{v}_k$$

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = \frac{1}{Z(\beta)} \exp(-E(\mathbf{x}_1, \dots, \mathbf{x}_N))$$

where \mathbf{v}_k is a different constant per site.

7.3.1 Application: Change Detection

Let \mathbf{z} be the latent variable and \mathbf{x} be the pixels,

$$\mathbf{z} = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\}$$

The joint likelihood is given by

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{z}_1, \dots, \mathbf{z}_N) &= p(\mathbf{x}_1, \dots, \mathbf{x}_N | \mathbf{z}_1, \dots, \mathbf{z}_N) p(\mathbf{z}_1, \dots, \mathbf{z}_N) \\ &= \left(\prod_{i=1}^N p(\mathbf{x}_i | \mathbf{z}_i) \right) \left(\frac{1}{Z(\beta)} \exp \left(\frac{1}{2} \sum_{k=1}^N \sum_{j \in N(k)} \mathbf{z}_k^\top \mathbf{z}_j \right) \right) \end{aligned}$$

Assuming the $p(\mathbf{x}|\mathbf{z})$ is Gaussian,

$$\begin{aligned} p(\mathbf{x}_i | \mathbf{z}_{i,k}) &= \prod_{k=1}^K \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k)^{z_{i,k}} \\ p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{z}_1, \dots, \mathbf{z}_N) &= p(\mathbf{x}_1, \dots, \mathbf{x}_N | \mathbf{z}_1, \dots, \mathbf{z}_N) p(\mathbf{z}_1, \dots, \mathbf{z}_N) \\ &= \left(\prod_{i=1}^N \prod_{k=1}^K \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k)^{z_{i,k}} \right) \left(\frac{1}{Z(\beta)} \exp \left(\frac{1}{2} \sum_{k=1}^N \sum_{j \in N(k)} \mathbf{z}_k^\top \mathbf{z}_j \right) \right) \end{aligned}$$

However, $Z(\beta)$ cannot be computed. Typical solution to resolve this is to use mean field approximation

$$\begin{aligned} p(\mathbf{Z}) &= \frac{1}{Z(\beta)} \exp \left(\frac{1}{2} \beta \sum_{k=1}^N \sum_{j \in N(k)} \mathbf{z}_k^\top \mathbf{z}_j \right) \\ &\approx \prod_{k=1}^N p(\mathbf{z}_k | \mathbb{E}[\mathbf{z}_j, \forall j \in N(k)]) \\ p(\mathbf{z}_k | \mathbb{E}[\mathbf{z}_j, \forall j \in N(k)]) &= \frac{\exp \left(\frac{1}{2} \beta \sum_{k=1}^N \sum_{j \in N(k)} \mathbf{z}_k^\top \mathbb{E}[\mathbf{z}_j] \right)}{\sum_{\mathbf{z}_k} \exp \left(\frac{1}{2} \beta \sum_{k=1}^N \sum_{j \in N(k)} \mathbf{z}_k^\top \mathbb{E}[\mathbf{z}_j] \right)} \end{aligned}$$

7.4 Energy Minimization for MRFs

Generally we want to minimize the energy function

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{z}_1, \dots, \mathbf{z}_N) &= p(\mathbf{x}_1, \dots, \mathbf{x}_N | \mathbf{z}_1, \dots, \mathbf{z}_N) p(\mathbf{z}_1, \dots, \mathbf{z}_N) \\ &= \left(\prod_{i=1}^N p(\mathbf{x}_i | \mathbf{z}_i) \right) \left(\frac{1}{Z(\beta)} \exp \left(\frac{1}{2} \sum_{k=1}^N \sum_{j \in N(k)} \mathbf{z}_k^\top \mathbf{z}_j \right) \right) \\ &= \frac{1}{Z(\beta)} \exp \left(\frac{1}{2} \sum_{k=1}^N \sum_{j \in N(k)} \mathbf{z}_k^\top \mathbf{z}_j - \sum_{k=1}^N \ln p(\mathbf{x}_k | \mathbf{z}_k) \right) \end{aligned}$$

There are two constraints

1. Data constraint: Labeling should reflect the observation
2. Smoothness constraint: Labeling should reflect spatial consistency

7.5 Gaussian Markov Random Fields

1. A Gaussian random field $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ that satisfies

$$p(\mathbf{x}_i | \{\mathbf{x}_j : j \neq i\}) = p(\mathbf{x}_i | \{\mathbf{x}_j : j \in N(i)\})$$

is a Gaussian Markov random field.

2. Another definition: A random vector $\mathbf{x} = (x_1, \dots, x_N)^\top$ is called a Gaussian Markov random field with respect to graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ with mean $\boldsymbol{\mu}$ and a positive semi-definite precision matrix \mathbf{Q} if its density has the form

$$p(\mathbf{x}) = (2\pi)^{-\frac{N}{2}} |\mathbf{Q}| \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top (\mathbf{x} - \boldsymbol{\mu})\right)$$

and

$$\begin{aligned} q_{ij} &\neq 0 \Leftrightarrow n_i \in \mathcal{N}(n_j) \\ q_{ij} &= 0 \Leftrightarrow \mathbf{x}_i \perp \mathbf{x}_j | \mathbf{x}_{\setminus\{i,j\}} \end{aligned}$$

3. Construction of \mathbf{Q} : Let A be the adjacency matrix of the graph (we assume that there's a recurrent edge to each node) and D is the degree matrix of the graph,

$$\begin{aligned} A = [a_{ij}] &\Rightarrow a_{ij} = \begin{cases} 1 & n_i \in \mathcal{N}(j) \\ 1 & i = j \\ 0 & \text{otherwise} \end{cases} \\ D = [D_{ij}] &\Rightarrow a_{ij} = \begin{cases} \sum_j a_{ij} & i = j \\ 0 & \text{otherwise} \end{cases} \\ \mathbf{Q} &= \mathbf{D} - \mathbf{A} \end{aligned}$$

Note that \mathbf{Q} is a sparse matrix.

4. We can also specify a Gaussian Markov random field via clique potentials:

$$V_C(\mathbf{x}_C) = \frac{1}{2} \left(\sum_{i \in C} \alpha_i^C x_i \right)^2 = \frac{1}{2} \left(\sum_{i \in \mathcal{V}} \alpha_i^C x_i \right)^2$$

where

$$i \notin C \Rightarrow \alpha_i^C = 0$$

Then the exponent of the GMRF density becomes

$$\begin{aligned}
 -\sum_{C \in \mathcal{C}} V_C(\mathbf{x}_C) &= -\frac{1}{2} \sum_{C \in \mathcal{C}} \left(\sum_{i \in \mathcal{V}} \alpha_i^C x_i \right)^2 \\
 &= -\frac{1}{2} \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{V}} \left(\sum_{C \in \mathcal{C}} \alpha_i^C \alpha_j^C \right) x_i x_j \\
 &= -\frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x}
 \end{aligned}$$

5. Example of a GMRF: AR(1)

$$\begin{aligned}
 \mathbf{x}_1 &\sim \mathcal{N}(0, (1 - \phi^2)^{-1}) \\
 \mathbf{x}_t | \mathbf{x}_{1:t-1} &\sim \mathcal{N}(\phi \mathbf{x}_{t-1}, 1) \\
 \mathbf{Q} &= \begin{bmatrix} 1 & \phi & 0 & \dots & 0 \\ -\phi & 1 + \phi^2 & -\phi & \ddots & 0 \\ 0 & \ddots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & -\phi & 1 \end{bmatrix}
 \end{aligned}$$

7.5.1 Application: Image Denoising

1. Consider an image to be a rectangular lattice with first-order pixel neighbourhoods:

- GMRF is used as the smoothing prior
- Cliques: pairs of vertically or horizontally adjacent pixels
- Clique potentials: squares of first-order differences

$$V_{\{(i,j), (i,j-1)\}}(x_{i,j}, x_{i,j-1}) = \frac{1}{2} (x_{i,j} - x_{i,j-1})^2$$

- \mathbf{Q} will be a matrix with tridiagonal blocks.

2. The model is given by

$$\begin{aligned}
 \mathbf{x} &= \mathbf{H}\mathbf{y} + \mathbf{u}, \mathbf{u} \sim \mathcal{N}(\mathbf{u} | \mathbf{0}, \sigma^2 \mathbf{I}) \text{ (observation model)} \\
 p(\mathbf{x}) &= (2\pi)^{-\frac{N}{2}} |\mathbf{Q}| \exp\left(-\frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x}\right) \text{ (GMRF prior)} \\
 \mathbb{E}[\mathbf{y}] &= [\sigma^2 \mathbf{Q} + \mathbf{H}^\top \mathbf{H}]^{-1} \mathbf{H}^\top \mathbf{y} \text{ (posterior mean)}
 \end{aligned}$$

3. The model is good for deblurring. However, when it comes to denoising, it oversmooth the image as the edge discontinuities are smoothed out. Preservation of discontinuities can be solved by

- Other prior models

- Hidden / latent binary random variables v to turn off clique potentials

$$V(x_{i,j}, x_{i,j-1}, v_{i,j}) = \frac{1}{2}(1 - v_{i,j})(x_{i,j} - x_{i,j-1})^2$$

- Robust potential functions (L_2 vs L_1 norm)

4. Some other possible potentials

- Convex potentials

$V(x) = x^p , p \in [1, 2]$	Generalized Gaussians
$V(x) = \begin{cases} x^2 & x < a \\ 2a x - a^2 & x \geq a \end{cases}$	Stevenson
$V(x) = 2a^2 \log \cosh(x/a)$	Green

- Non-convex potentials

$V(x) = [\min(x , a)]^2$	Blake, Zisserman
$V(x) = \frac{x^2}{x^2 + a^2}$	Geman, McClure

8 Useful Identities:

- Woodbury Identity:

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

$$(A + BC)^{-1} = A^{-1} - A^{-1}B(I + CA^{-1}B)^{-1}CA^{-1}$$

- Matrix Inversion Identities:

$$(A^{-1} + B^{-1})^{-1} = A(A + B)^{-1}B = B(A + B)^{-1}A$$

- 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:

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}\right)$$

We have

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_{x|y}, \Sigma_{x|y})$$

$$\boldsymbol{\mu}_{x|y} = \boldsymbol{\mu}_x + \Sigma_{xy}\Sigma_{yy}^{-1}(\mathbf{y} - \boldsymbol{\mu}_y)$$

$$\Sigma_{x|y} = \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx}$$

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

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

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b})$$

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

- Derivative of inverse

$$\frac{\partial \mathbf{Y}^{-1}}{\partial \mathbf{x}} = -\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial \mathbf{x}} \mathbf{Y}^{-1}$$

$$\frac{\partial \mathbf{a}^\top \mathbf{X}^{-1} \mathbf{B}}{\partial \mathbf{X}} = -\mathbf{X}^{-\top} \mathbf{a} \mathbf{b}^\top \mathbf{X}^{-\top}$$

- Derivative of trace

$$\frac{\partial \text{tr}(\mathbf{X}\mathbf{A})}{\partial \mathbf{X}} = \mathbf{A}^\top$$

$$\frac{\partial \text{tr}(\mathbf{X}^\top \mathbf{A})}{\partial \mathbf{X}} = \mathbf{A}$$

$$\frac{\partial \text{tr}(\mathbf{A}\mathbf{X}\mathbf{B})}{\partial \mathbf{X}} = \mathbf{A}^\top \mathbf{B}^\top$$

- Derivative of determinant

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

- Matrix identities for trace

$$\begin{aligned} \mathbf{a}^\top \mathbf{a} &= \text{tr}(\mathbf{a} \mathbf{a}^\top) \\ \text{tr}(\mathbf{A} \mathbf{B} \mathbf{C}) &= \text{tr}(\mathbf{B} \mathbf{C} \mathbf{A}) = \text{tr}(\mathbf{C} \mathbf{A} \mathbf{B}) \end{aligned}$$