

NOTES

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

495 Advanced Statistical Machine Learning and Pattern Recognition

Author:

Thomas Teh (CID: 0124 3008)

Date: January 31, 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

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

1.2.1 Formulation of the Gaussian Mixture Models

Let \mathbf{z} be a K -dimensional binary random variable with 1-of- K representation.

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

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

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

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:

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

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 X , we were also given the values of the latent variables Z , the likelihood and the log likelihood function are given by

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

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

1. Expectation Step:

Taking the expectation on the log likelihood

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

$$= G(\theta)$$

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

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

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

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

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

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

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

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

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 $\boldsymbol{\mu}_k$, 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

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

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

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

where

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

4. Evaluate the log likelihood

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

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

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

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

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

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

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

$$\begin{aligned}
 \boldsymbol{\mu}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \\
 \pi_k &= \frac{N_k}{N} \\
 N_k &= \sum_{n=1}^N \gamma(z_{nk})
 \end{aligned}$$

1.4 Convergence of the EM Algorithm

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \Rightarrow \ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + KL(q \parallel p)$$

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left[\frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right]$$

$$KL(q \parallel p) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left[\frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right]$$

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

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

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.

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

$$p(\mathbf{x}_i | \mathbf{y}_i, \mathbf{W}, \boldsymbol{\mu}, \sigma) p(\mathbf{y}_i) = \left[(2\pi\sigma^2)^{-\frac{F}{2}} e^{-\frac{1}{2\sigma^2}(\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{y}_i)^\top (\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{W}\mathbf{y}_i)} \right] \left[(2\pi)^{-\frac{d}{2}} e^{-\frac{1}{2}\mathbf{y}_i^\top \mathbf{y}_i} \right]$$

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):

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

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

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

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

4. Maximizing the likelihood and solve for the parameters

$$\mathbf{S}_t = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top \quad (\mathbf{S}_t \text{ is the covariance matrix})$$

$$\sigma^2 = \frac{1}{F-d} \sum_{j=d+1}^F \lambda_j$$

$$\mathbf{W}_d = \mathbf{U}_d (\boldsymbol{\Lambda} - \sigma^2 \mathbf{I})^{\frac{1}{2}} \mathbf{V}^\top$$

5. Hence we no longer have a projection but:

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

$$\widehat{\mathbf{x}}_i = \mathbf{W} \mathbb{E}_{p(\mathbf{y}_i | \mathbf{x}_i)}[\mathbf{y}_i] + \boldsymbol{\mu}$$

2.2 EM PPCA

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

$$p(\mathbf{X}, \mathbf{Y} | \boldsymbol{\theta}) = \prod_{i=1}^N p(\mathbf{x}_i | \mathbf{y}_i, \boldsymbol{\theta}) p(\mathbf{y}_i)$$

$$\ln p(\mathbf{X}, \mathbf{Y} | \boldsymbol{\theta}) = \sum_{i=1}^N (\ln p(\mathbf{x}_i | \mathbf{y}_i, \boldsymbol{\theta}) + \ln p(\mathbf{y}_i))$$

$$\ln p(\mathbf{x}_i | \mathbf{y}_i, \boldsymbol{\theta}) = -\frac{F}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i - \boldsymbol{\mu})^\top (\mathbf{x}_i - \mathbf{W}\mathbf{y}_i - \boldsymbol{\mu})$$

$$\ln p(\mathbf{y}_i) = -\frac{1}{2} \mathbf{y}_i^\top \mathbf{y}_i - \frac{D}{2} \ln 2\pi$$

2. Take the expectation on the log likelihood on the joint distribution:

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

$$\mathbb{E}_{p(\mathbf{Y}|\mathbf{X})}[\ln p(\mathbf{X}, \mathbf{Y} | \boldsymbol{\theta})] =$$