



Introduction to Statistical Machine Learning

Cheng Soon Ong & Christian Walder

Machine Learning Research Group

Data61 | CSIRO

and

College of Engineering and Computer Science

The Australian National University

Canberra

February – June 2017

(Many figures from C. M. Bishop, "Pattern Recognition and Machine Learning")



Part XIX

Mixture Models and EM 2



- Starting point is the log of the likelihood function

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

- Critical point of $\ln p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ w.r.t. $\boldsymbol{\mu}_k$

$$0 = \sum_{n=1}^N \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{\gamma(z_{nk})} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

- Therefore

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

where the **effective number of points assigned to Gaussian k** is $N_k = \sum_{n=1}^N \gamma(z_{nk})$.



- Maximum of the log of the likelihood function for

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

- Similarly for the covariance matrix

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

- and for the mixing coefficients π_k (using a Lagrange multiplier as $\sum_k \pi_k = 1$)

$$\pi_k = \frac{N_k}{N}.$$

- This is not a closed form solution because the responsibilities $\gamma(z_{nk})$ depend on $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$.



- Given a Gaussian mixture and data \mathbf{X} , maximise the log likelihood w.r.t. the parameters $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$.

- Initialise the means $\boldsymbol{\mu}_k$, covariances $\boldsymbol{\Sigma}_k$ and mixing coefficients π_k . Evaluate the log likelihood function.
- E step** : Evaluate the $\gamma(z_k)$ using the current parameters

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

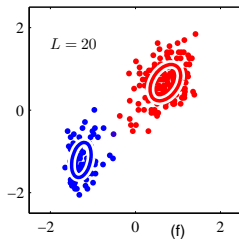
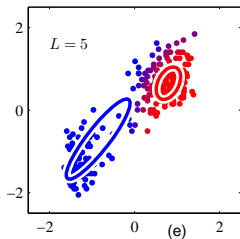
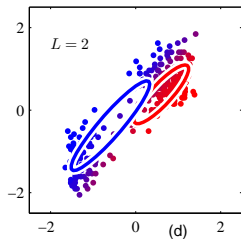
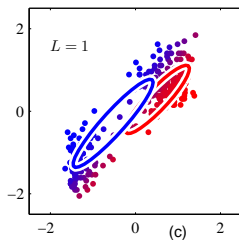
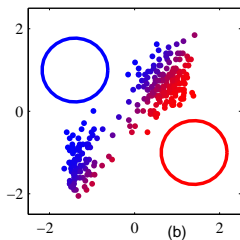
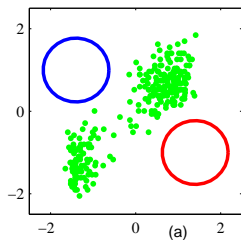
- M step** : Re-estimate the parameters using the current $\gamma(z_k)$

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \pi_k^{\text{new}} = \frac{N_k}{N}$$
$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})(\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T$$

- Evaluate the log likelihood, if not converged then goto 2.

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

EM for Gaussian Mixtures - Example





- Assume a Gaussian mixture model.
- Covariance matrices given by $\epsilon \mathbf{I}$, where ϵ is shared by all components.
- Then

$$p(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{(2\pi\epsilon)^{D/2}} \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\}.$$

- Keep ϵ fixed, do not re-estimate.
- Responsibilities

$$\gamma(z_{nk}) = \frac{\pi_k \exp \left\{ -\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 / 2\epsilon \right\}}{\sum_j \pi_j \exp \left\{ -\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 / 2\epsilon \right\}}$$

- Taking the limit $\epsilon \rightarrow 0$, the term in the denominator for which $\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2$ is the smallest will go to zero most slowly.



- Assume a Gaussian mixture model.

$$\gamma(z_{nk}) = \frac{\pi_k \exp \{ -\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 / 2\epsilon \}}{\sum_j \pi_j \exp \{ -\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 / 2\epsilon \}}$$

- Therefore

$$\gamma(z_{nk}) = \begin{cases} 1 & \text{if } \|\mathbf{x}_n - \boldsymbol{\mu}_k\| < \|\mathbf{x}_n - \boldsymbol{\mu}_j\| \quad \forall j \neq k \\ 0 & \text{otherwise} \end{cases}$$

- Holds independent of π_k as long as none are zero.
- Hard assignment to exactly one cluster : K -means.

$$\lim_{\epsilon \rightarrow 0} \gamma(z_{nk}) = r_{nk}$$



- Set of D binary variables $x_i, i = 1, \dots, D$.
- Each governed by a Bernoulli distribution with parameter μ_i . Therefore

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

- Expectation and covariance

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

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



Mixture of Bernoulli Distributions

- Mixture

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

with

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

- Similar calculation as with mixture of Gaussian

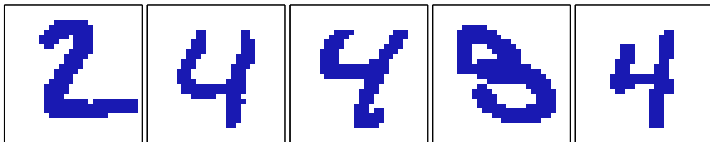
$$\gamma(z_{nk}) = \frac{\pi_k p(\mathbf{x}_n \mid \boldsymbol{\mu}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n \mid \boldsymbol{\mu}_j)}$$

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

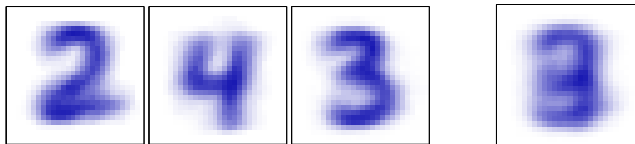
$$\bar{\mathbf{x}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \boldsymbol{\mu}_k = \bar{\mathbf{x}}$$

$$\pi_k = \frac{N_k}{N}$$

EM for Mixture of Bernoulli Distributions - Digits



Examples from a digits data set, each pixel taken only binary values.



Parameters μ_{ki} for each component in the mixture.

Fit to one multivariate Bernoulli distribution.



- EM finds the maximum likelihood solution for models with latent variables.
- Two kinds of variables
 - Observed variables \mathbf{X}
 - Latent variables \mathbf{Z}plus model parameters θ .
- Log likelihood is then

$$\ln p(\mathbf{X} | \theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \theta) \right\}$$

- Optimisation problem due to the log-sum.
- Assume maximisation of the distribution $p(\mathbf{X}, \mathbf{Z} | \theta)$ over the **complete data set** $\{\mathbf{X}, \mathbf{Z}\}$ is straightforward.
- But we only have the **incomplete data set** $\{\mathbf{X}\}$ and the posterior distribution $p(\mathbf{Z} | \mathbf{X}, \theta)$.



- Key idea of EM: As \mathbf{Z} is not observed, work with an 'averaged' version $Q(\theta, \theta^{\text{old}})$ of the complete log-likelihood $\ln p(\mathbf{X}, \mathbf{Z} | \theta)$, averaged over all states of \mathbf{Z} .

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



- 1 Choose an initial setting for the parameters θ^{old} .
- 2 **E step** Evaluate $p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$.
- 3 **M step** Evaluate θ^{new} given by

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

where

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

- 4 Check for convergence of log likelihood or parameter values. If not yet converged, then

$$\theta^{\text{old}} = \theta^{\text{new}}$$

and go to step 2.



- Start with the product rule for the observed variables \mathbf{X} , the unobserved variables \mathbf{Z} , and the parameters θ

$$\ln p(\mathbf{X}, \mathbf{Z} | \theta) = \ln p(\mathbf{Z} | \mathbf{X}, \theta) + \ln p(\mathbf{X} | \theta).$$

- Apply $\sum_{\mathbf{Z}} q(\mathbf{Z})$ with arbitrary $q(\mathbf{Z})$ to the formula

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z} | \theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{Z} | \mathbf{X}, \theta) + \ln p(\mathbf{X} | \theta).$$

- Rewrite as

$$\ln p(\mathbf{X} | \theta) = \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q(\mathbf{Z})}}_{\mathcal{L}(q, \theta)} - \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z} | \mathbf{X}, \theta)}{q(\mathbf{Z})}}_{\text{KL}(q||p)}$$

- $\text{KL}(q||p)$ is the **Kullback-Leibler** divergence.



- ‘Distance’ between two distributions $p(y)$ and $q(y)$

$$\text{KL}(q\|p) = \sum_y q(y) \ln \frac{q(y)}{p(y)} = - \sum_y q(y) \ln \frac{p(y)}{q(y)}$$

$$\text{KL}(q\|p) = \int q(y) \ln \frac{q(y)}{p(y)} dy = - \int q(y) \ln \frac{p(y)}{q(y)} dy$$

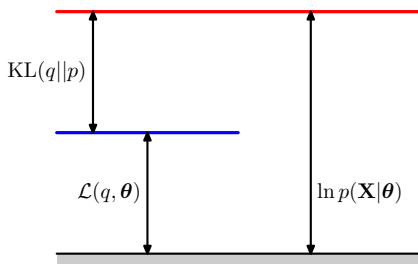
- $\text{KL}(q\|p) \geq 0$
- not symmetric: $\text{KL}(q\|p) \neq \text{KL}(p\|q)$
- $\text{KL}(q\|p) = 0$ iff $q = p$.
- invariant under parameter transformations
- Example: Kullback-Leibler divergence between two normal distributions $q(x) = \mathcal{N}(x | \mu_1, \sigma_1)$ and $p(x) = \mathcal{N}(x | \mu_2, \sigma_2)$

$$\text{KL}(q\|p) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$



- The two parts of $\ln p(\mathbf{X} | \theta)$

$$\ln p(\mathbf{X} | \theta) = \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q(\mathbf{Z})}}_{\mathcal{L}(q, \theta)} - \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z} | \mathbf{X}, \theta)}{q(\mathbf{Z})}}_{\text{KL}(q||p)}$$

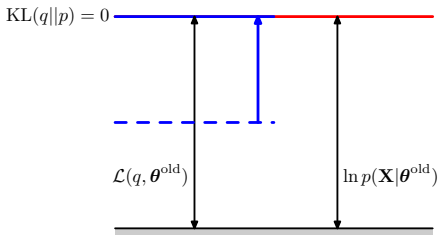




EM Algorithm - E Step

- Hold θ^{old} fixed. Maximise the lower bound $\mathcal{L}(q, \theta^{\text{old}})$ with respect to $q(\cdot)$.
- $\mathcal{L}(q, \theta^{\text{old}})$ is a functional.
- $\ln p(\mathbf{X} | \theta)$ does NOT depend on $q(\cdot)$.
- Maximum for $\mathcal{L}(q, \theta^{\text{old}})$ will occur when the Kullback-Leibler divergence vanishes.
- Therefore, choose $q(\mathbf{Z}) = p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$

$$\ln p(\mathbf{X} | \theta) = \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q(\mathbf{Z})}}_{\mathcal{L}(q, \theta)} - \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z} | \mathbf{X}, \theta)}{q(\mathbf{Z})}}_{\text{KL}(q||p)}$$

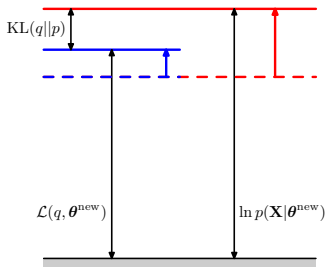




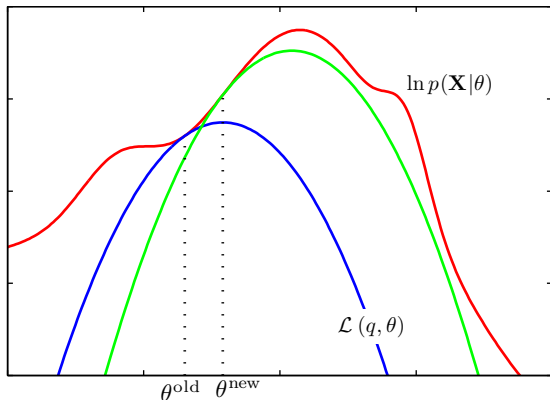
EM Algorithm - M Step

- Hold $q(\cdot) = p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$ fixed. Maximise the lower bound $\mathcal{L}(q, \theta)$ with respect to θ :
 $\theta^{\text{new}} = \arg \max_{\theta} \mathcal{L}(q, \theta^{\text{old}}) = \arg \max_{\theta} \sum_{\mathbf{Z}} q(\cdot) \ln p(\mathbf{X}, \mathbf{Z} | \theta)$
- $\mathcal{L}(q, \theta^{\text{new}}) > \mathcal{L}(q, \theta^{\text{old}})$ unless maximum already reached.
- As $q(\cdot) = p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$ is fixed, $p(\mathbf{Z} | \mathbf{X}, \theta^{\text{new}})$ will not be equal to $q(\cdot)$, and therefore the Kullback-Leiber distance will be greater than zero (unless converged).

$$\ln p(\mathbf{X} | \theta) = \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q(\mathbf{Z})}}_{\mathcal{L}(q, \theta)} - \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z} | \mathbf{X}, \theta)}{q(\mathbf{Z})}}_{\text{KL}(q||p)}$$



EM Algorithm - Parameter View



Red curve : incomplete data likelihood.

Blue curve : After E step. Green curve : After M step.