

Faculty of Science



Mixture Modeling and EM Algorithm

Christian Igel
igel@diku.dk

Department of Computer Science
University of Copenhagen

pictures from C. M. Bishop: *Pattern Recognition and Machine Learning*, Springer, 2006



Outline

- 1 Density Estimation
- 2 Mixture Modeling
- 3 Learning Mixtures with Expectation Maximization
- 4 General Expectation Maximization



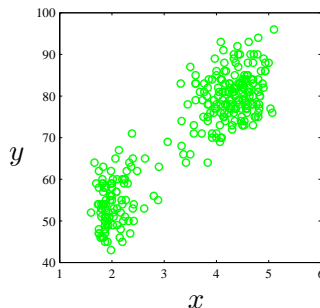
Outline

- 1 Density Estimation
- 2 Mixture Modeling
- 3 Learning Mixtures with Expectation Maximization
- 4 General Expectation Maximization



Example: Old Faithful

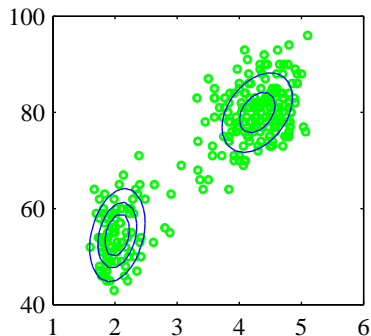
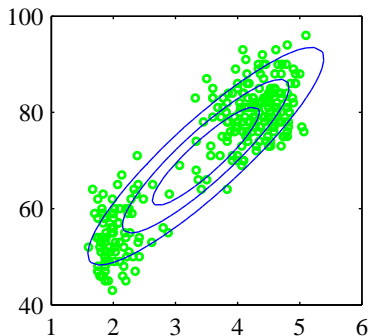
- Hydrothermal geyser in Yellowstone National Park, Wyoming, USA



- x -axis duration of eruption in minutes
- y -axis time to next eruption in minutes



Density estimation



Outline

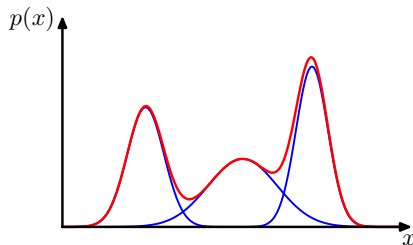
- 1 Density Estimation
- 2 Mixture Modeling**
- 3 Learning Mixtures with Expectation Maximization
- 4 General Expectation Maximization



Mixture modeling

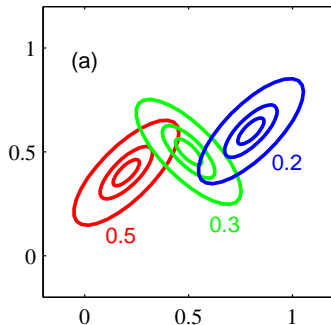
Mixture modeling: Describing complex distributions by convex combinations of simpler distributions:

$$p(\mathbf{x}) = \sum_{k=1}^K p(k)p(\mathbf{x}|k) \quad , \quad \sum_{k=1}^K p(k) = 1$$



Red line: $p(x)$

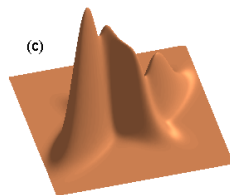
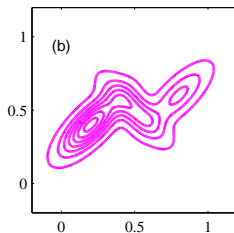
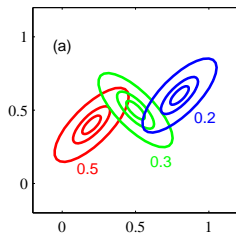
Blue lines: $p(k)p(x|k)$, $k = 1, 2, 3$



Gaussians Mixture Models (GMMs)

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

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



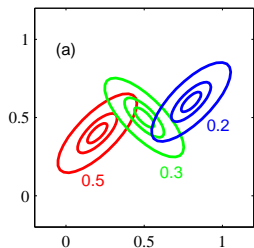
Generative process

To sample GMM distribution

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

we take a generative view:

- 1 First draw a component number k with relative probabilities π_k ,
- 2 then draw a random vector \mathbf{x} from the given component with density $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$.



Outline

- 1 Density Estimation
- 2 Mixture Modeling
- 3 Learning Mixtures with Expectation Maximization**
- 4 General Expectation Maximization



Maximum likelihood learning

- Training set is $S = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_N\}$
- Likelihood function (assuming i.i.d. data) is given by

$$p(S | \boldsymbol{\theta}) = \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\theta})$$

- Parameters $\boldsymbol{\theta} = \{\pi_1, \dots, \pi_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K\}$
- Cost function is the negative logarithmic likelihood (notice sum inside log):

$$\begin{aligned} E(\boldsymbol{\theta}) &= - \sum_{n=1}^N \log p(\mathbf{x}_n | \boldsymbol{\theta}) = - \sum_{n=1}^N \log \sum_{k=1}^K p(\mathbf{x}_n | \boldsymbol{\theta}_k) \pi_k \\ &= - \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned}$$



GMM maximum likelihood

- Necessary condition for an extremum is that the partial derivatives w.r.t. the parameters

$$\boldsymbol{\theta} = \{\pi_k, \dots, \boldsymbol{\mu}_k, \dots, \boldsymbol{\Sigma}_k, \dots\} \text{ of}$$

$$E(\boldsymbol{\theta}) = - \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

are zero.

- We define the *responsibility* ($\gamma(z_{nk})$ in Bishop's textbook)

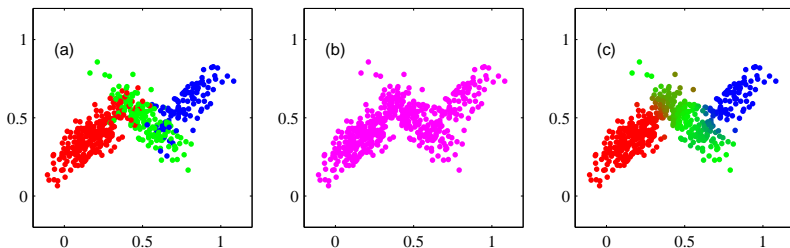
$$\gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})} \in [0, 1] .$$

- Responsibilities depend on GMM parameters; let us assume them to be fixed for the moment.



Responsibility

$$\begin{aligned}\gamma_{nk} &= \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})} \\ &= \frac{p(k)p(\mathbf{x}_n | k)}{\sum_{k'} p(k')p(\mathbf{x}_n | k')} = \frac{p(k)p(\mathbf{x}_n | k)}{p(\mathbf{x}_n)} \\ &= p(\text{pattern } n \text{ generated by component } k \mid \mathbf{x}_n) \in [0, 1]\end{aligned}$$



Maximizing likelihood I

Using (recall $\frac{\partial}{\partial \mathbf{z}} \mathbf{z}^\top \mathbf{B} \mathbf{z} = (\mathbf{B} + \mathbf{B}^\top) \mathbf{z}$)

$$\frac{\partial}{\partial \boldsymbol{\mu}_k} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

setting

$$\frac{\partial}{\partial \boldsymbol{\mu}_k} E(\boldsymbol{\theta}) = - \frac{\partial}{\partial \boldsymbol{\mu}_k} \left\{ \sum_{n=1}^N \log \sum_{k'=1}^K \pi_{k'} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'}) \right\}$$

to zero yields

$$0 = \sum_{n=1}^N \gamma_{nk} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad .$$



Maximizing likelihood II

- From $0 = \sum_{n=1}^N \gamma_{nk} \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)$ we conclude:

$$\boldsymbol{\mu}_k \leftarrow \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n \text{ with } N_k = \sum_{n=1}^N \gamma_{nk}$$

- The same can be done for the other parameters, but is more difficult (and an excellent exercise). We get:

$$\Sigma_k \leftarrow \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \text{ and } \pi_k \leftarrow \frac{N_k}{N}$$

- This is not the gradient of $E(\boldsymbol{\theta})$, because responsibilities depend on $\boldsymbol{\theta}$.
- We do not do std. gradient descent on negative logarithmic likelihood $E(\boldsymbol{\theta})$, but apply an iterative, two-step optimization scheme.



Expectation Maximization (EM) for GMM

EM Algorithm for GMMs

init parameters

while *termination criterion not met* **do**

/* E step */

for $k = 1, \dots, K$ **do**

$$\gamma_{nk} \leftarrow \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k'=1}^K \pi_{k'} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})}$$

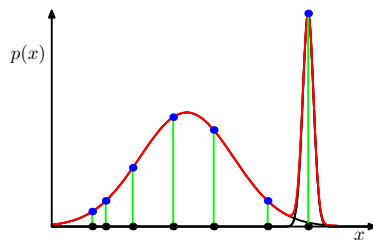
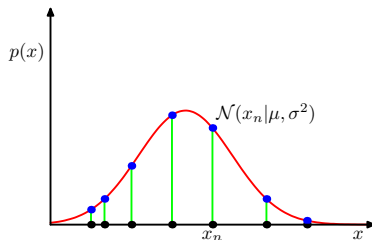
/* M step */

$$N_k \leftarrow \sum_{n=1}^N \gamma_{nk} \quad , \quad \pi_k \leftarrow \frac{N_k}{N}$$

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



Nature of the maximum likelihood solution

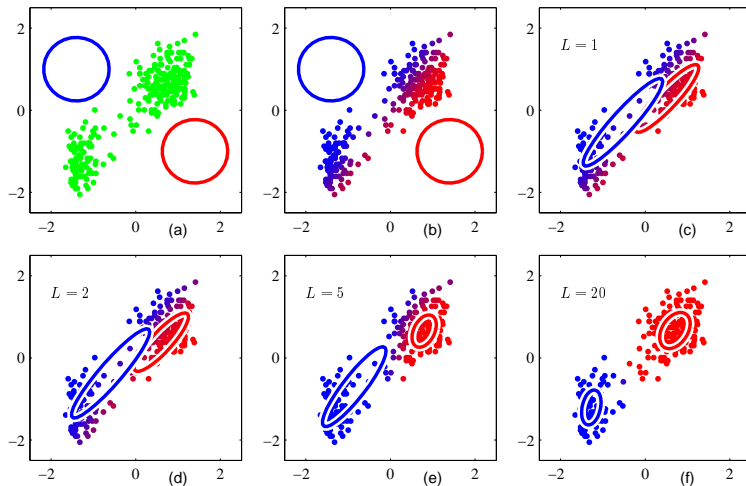


$$E(\boldsymbol{\theta}) = - \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Consider cost when $\boldsymbol{\mu}_k = \mathbf{x}_n$, $\pi_k > 0$ and $\boldsymbol{\Sigma}_k \rightarrow 0$



GMM for Old Faithful



Outline

- 1 Density Estimation
- 2 Mixture Modeling
- 3 Learning Mixtures with Expectation Maximization
- 4 General Expectation Maximization**



Problem statement

- Given a set of observed \mathbf{X} and hidden (latent) \mathbf{Z} random variables as well as a probabilistic model p with parameters θ .
- We want to compute and/or maximize the likelihood

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

or its logarithm

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

- This is difficult, especially because of the sum which prevents the logarithm to act directly on the joint distribution.
- In contrast, $\ln p(\mathbf{X}, \mathbf{Z}|\theta)$ may be easy to compute if p belongs to the exponential family of distributions.



Kullback-Leibler divergence

Kullback-Leibler (KL) divergence between two distribution p and q over \mathcal{Z}

$$\text{KL}(q||p) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z})}{q(\mathbf{Z})} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{q(\mathbf{Z})}{p(\mathbf{Z})}$$

(sum turns to integral for continuous random variables) is

- a (non-symmetric) measure of difference between distributions,
- always positive, zero iff the distributions are the same.



Variational lower bound

Proposition

Given some distribution $q(\mathbf{Z})$ over the hidden variables, it holds

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{F}(q, \boldsymbol{\theta}) + \text{KL}(q\|p_{\mathbf{Z}|\mathbf{X}})$$

with Kullback-Leibler divergence

$$\text{KL}(q\|p_{\mathbf{Z}|\mathbf{X}}) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})}$$

and variational lower bound (“free energy”)

$$\mathcal{F}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}, \mathbf{X}|\boldsymbol{\theta})}{q(\mathbf{Z})} .$$



Proof

We substitute

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

into the definition of the variational lower bound

$$\mathcal{F}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}, \mathbf{X}|\boldsymbol{\theta})}{q(\mathbf{Z})} .$$

and get

$$\begin{aligned} \mathcal{F}(q, \boldsymbol{\theta}) &= \sum_{\mathbf{Z}} q(\mathbf{Z}) (\ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) + \ln p(\mathbf{X}|\boldsymbol{\theta}) - \ln q(\mathbf{Z})) \\ &= -\text{KL}(q||p_{\mathbf{Z}|\mathbf{X}}) + \ln p(\mathbf{X}|\boldsymbol{\theta}) . \end{aligned}$$

Substituting this into the decomposition shows its correctness.



Kullback-Leibler divergence and free energy

- Free energy \mathcal{F} differs from KL divergence as it
 - has the reverse sign and
 - contains the joint instead of the conditional distribution.
- The properties of the KL divergence imply the lower bound property:

$$\mathcal{F}(q(\mathbf{Z}), \boldsymbol{\theta}) \leq \ln p(\mathbf{X}|\boldsymbol{\theta})$$

- $\ln p(\mathbf{X}|\boldsymbol{\theta})$ is also called logarithmic *evidence*, and the variational lower bound is referred to as *evidence lower bound (ELBO)*



General EM algorithm

- The EM algorithm is an iterative method for increasing $p(\mathbf{X}|\boldsymbol{\theta})$ by adapting $\boldsymbol{\theta}$.
- In each iteration n , an E step (expectation step) and an M step (maximization step) is performed.

EM Algorithm

```
1 init  $\boldsymbol{\theta}^{(0)}$ ,  $n \leftarrow 0$ 
2 while termination criterion not met do
3    $q^{(n+1)} \leftarrow \operatorname{argmax}_q \mathcal{F}(q, \boldsymbol{\theta}^{(n)})$  (E step)
4    $\boldsymbol{\theta}^{(n+1)} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta}} \mathcal{F}(q^{(n+1)}, \boldsymbol{\theta})$  (M step)
5    $n \leftarrow n + 1$ 
```



E step

- We have

$$\begin{aligned}\mathcal{F}(q, \boldsymbol{\theta}) &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}, \mathbf{X} | \boldsymbol{\theta})}{q(\mathbf{Z})} \\ &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}) p(\mathbf{X} | \boldsymbol{\theta})}{q(\mathbf{Z})} \\ &= -\text{KL}(q \| p_{\mathbf{Z} | \mathbf{X}}) + \ln p(\mathbf{X} | \boldsymbol{\theta}) .\end{aligned}$$

- This expression is maximized w.r.t. q if the Kullback-Leibler divergence vanishes, that is, if $q(\mathbf{Z}) = p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})$.
- That is, after an optimal E step, we have:

$$\mathcal{F}(q^{(n+1)}, \boldsymbol{\theta}^{(n)}) = \ln p(\mathbf{X} | \boldsymbol{\theta}^{(n)})$$



M step

- Plugging $q^{(n+1)}(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{(n)})$ into

$$\mathcal{F}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{Z}, \mathbf{X}|\boldsymbol{\theta}) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln q(\mathbf{Z})$$

gives

$$\begin{aligned} \mathcal{F}(q^{(n+1)}, \boldsymbol{\theta}) = & \underbrace{\sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{(n)}) \ln p(\mathbf{Z}, \mathbf{X}|\boldsymbol{\theta})}_{Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(n)})} - \\ & \underbrace{\sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{(n)}) \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{(n)})}_{-H(q^{(n+1)})} \end{aligned}$$

- As $-H(q^{(n+1)})$ is independent of $\boldsymbol{\theta}$, optimizing $\mathcal{F}(q^{(n+1)}, \boldsymbol{\theta})$ w.r.t. to $\boldsymbol{\theta}$ in the M step corresponds to

$$\boldsymbol{\theta}^{(n+1)} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(n)}) .$$



EM increases likelihood

- Assume the algorithm has not converged and $\boldsymbol{\theta}^{(n+1)} \neq \boldsymbol{\theta}^{(n)}$.
- The M step gave
$$\mathcal{F}(q^{(n+1)}, \boldsymbol{\theta}^{(n+1)}) > \mathcal{F}(q^{(n+1)}, \boldsymbol{\theta}^{(n)}) = \ln p(\mathbf{X}|\boldsymbol{\theta}^{(n)}).$$
- Further, $q^{(n+1)}(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{(n)})$ and $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{(n+1)})$ have a positive Kullback-Leibler divergence.
- Thus, we have

$$\begin{aligned} \ln p(\mathbf{X}|\boldsymbol{\theta}^{(n+1)}) = & \underbrace{\mathcal{F}(q^{(n+1)}, \boldsymbol{\theta}^{(n+1)})}_{> \ln p(\mathbf{X}|\boldsymbol{\theta}^{(n)})} + \underbrace{\text{KL}(q^{(n+1)} \| p_{\mathbf{Z}|\mathbf{X}}^{(n+1)})}_{> 0} > \ln p(\mathbf{X}|\boldsymbol{\theta}^{(n)}) \end{aligned}$$

showing the increase of both the log-likelihood $\ln p(\mathbf{X}|\boldsymbol{\theta})$ as well as its lower bound $\mathcal{F}(q, \boldsymbol{\theta})$.



GMM learning and general EM

- Hidden/latent discrete random variables $\mathbf{Z} = (z_1, \dots, z_N)$
- 1-hot encoding: $z_i \in \mathbb{R}^K$, $[z_i]_k = 1$ and $[z_i]_j = 0$ for $j \neq k$ if x_i was generated by component k
- Distribution $q(\mathbf{Z})$ can be described by $N \times K$ values q_{nk} giving the probability that pattern n was generated by component k .
- $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{(i)})$ is given by $q_{nk} = \gamma_{nk}^{(i)}$ maximizing \mathcal{F} in the E step
- M step directly corresponds to maximizing the likelihood for fixed responsibilities.



Concluding remarks

- EM algorithm is a general technique for maximum likelihood estimation.
- It can be applied to adapt Gaussian mixture models; k -means clustering can be interpreted in the EM framework.
- The general EM algorithm is the basis for *variational methods*.



Note

$$\begin{aligned}\mathcal{F}(q, \boldsymbol{\theta}) &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}, \mathbf{X} | \boldsymbol{\theta})}{q(\mathbf{Z})} \\&= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\theta}) p(\mathbf{Z} | \boldsymbol{\theta})}{q(\mathbf{Z})} \\&= \mathbb{E}_{q(\mathbf{Z})} (\ln p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\theta}) + \ln p(\mathbf{Z} | \boldsymbol{\theta}) - \ln q(\mathbf{Z})) \\&= \mathbb{E}_{q(\mathbf{Z})} (\ln p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\theta})) + \mathbb{E}_{q(\mathbf{Z})} (\ln p(\mathbf{Z} | \boldsymbol{\theta}) - \ln q(\mathbf{Z})) \\&= \mathbb{E}_{q(\mathbf{Z})} (\ln p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\theta})) - \text{KL}(q(\mathbf{Z}) \| p(\mathbf{Z} | \boldsymbol{\theta}))\end{aligned}$$

All considerations hold if q depends on \mathbf{X} . Thus, we can replace $q(\mathbf{Z})$ by $q(\mathbf{Z} | \mathbf{X})$. If we additionally have a prior $p(\mathbf{Z} | \boldsymbol{\theta})$ that does not depend on $\boldsymbol{\theta}$, we get the ELBO:

$$\mathbb{E}_{q(\mathbf{Z} | \mathbf{X})} (\ln p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\theta})) - \text{KL}(q(\mathbf{Z} | \mathbf{X}) \| p(\mathbf{Z}))$$

This resembles the common form of the ELBO for variational autoencoders (VAEs).

