

General EM (cont'd)

STATS 303 Statistical Machine Learning

Spring 2022

Lecture 10

general EM

complete dataset with latent variables

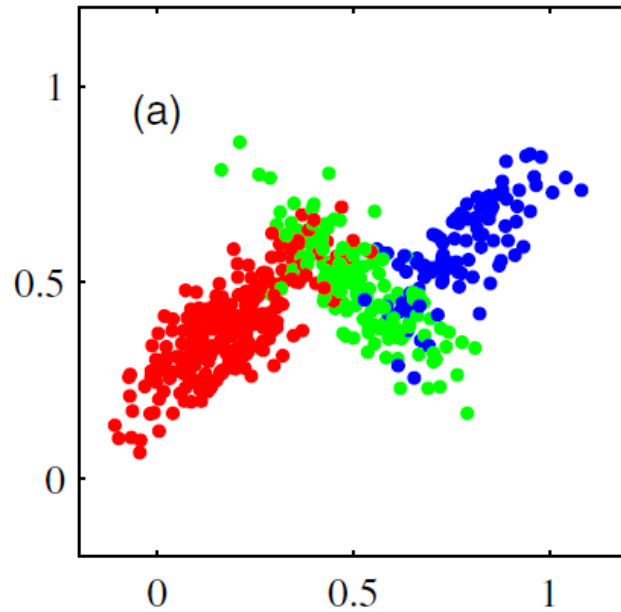
- Suppose \mathbf{X} is the data matrix, and \mathbf{Z} the corresponding latent variables (assumed to be discrete). Then

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

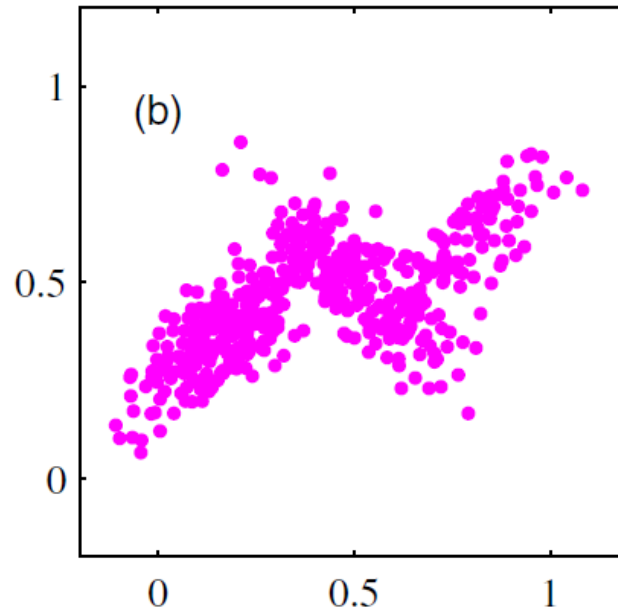
- $\{\mathbf{X}, \mathbf{Z}\}$ is called the **complete** data set; \mathbf{X} is **incomplete**
- In practice, we are not given the complete data set; the only way we estimate \mathbf{Z} is by the **posterior**

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

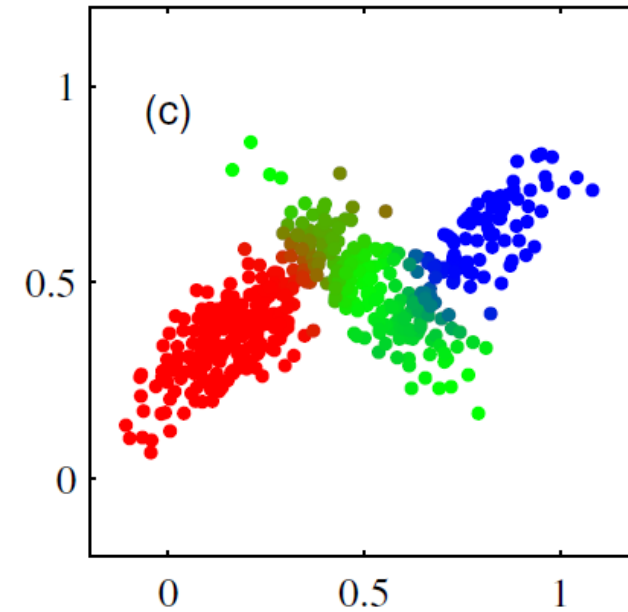
complete dataset with latent variables



complete data set
with both $\{\mathbf{X}, \mathbf{Z}\}$



incomplete data
set with only \mathbf{X}



\mathbf{X} with the
posterior $p(\mathbf{Z}|\mathbf{X}, \theta)$

general EM

(E-step): expectation

1. for fixed parameters, find $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$

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

$$\mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})} \ln p(\mathbf{X}, \mathbf{Z}|\theta)$$

(M-step): maximization

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

GMM revisited

- The likelihood of the complete data set is

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

$p(\mathbf{z} | \boldsymbol{\pi}) p(\mathbf{x} | \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$

- Taking logarithm yields

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

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

GMM revisited

$$p(z_n = c_k | x_n, \mu, \Sigma, \pi) \propto \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)$$

- By Bayes' theorem, $p(\mathbf{Z} | \mathbf{X}, \mu, \Sigma, \pi) \propto \prod_{n=1}^N \prod_{k=1}^K [\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)]^{z_{nk}}$

- What does this imply for each \mathbf{z}_n ?

$$p(z_n | x_n, \mu, \Sigma, \pi) \propto \prod_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)^{z_{nk}}$$

- Under the posterior distribution, the conditional expectation of the indicator z_{nk} is given by

$$\mathbb{E}[z_{nk}] = \frac{p(z_{nk} = 1 | x_n, \mu, \Sigma, \pi)}{\sum_{j=1}^K p(z_{nj} = 1 | x_n, \mu, \Sigma, \pi)} = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)} = \gamma(z_{nk})$$

Eq (9.39) in [Bi] is not very straight-forward. Take care when reading it.

GMM revisited

- The **expectation** of the complete log likelihood is thus

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

$\nwarrow \sim p(\mathbf{z} | \mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})$

- Maximizing the above yields the same parameters $\left\{ \begin{array}{l} \boldsymbol{\mu}_k \\ \boldsymbol{\Sigma}_k \\ \pi_k \end{array} \right.$ as before.

general EM (the general view)

- The **expectation-maximization** algorithm, or **EM** algorithm, is a general technique for finding maximum likelihood solutions for probabilistic models having latent variables.

- Consider

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

- Assuming \mathbf{Z} is discrete and dealing with $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ is easier than $p(\mathbf{X}|\boldsymbol{\theta})$.
- Introduce a distribution $q(\mathbf{Z})$ over the latent variables.

general EM

- Claim

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

KL divergence from q to p



$p(z|x, \theta)$

where

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

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

$$= \mathbb{E}_{q(z)} \ln \left(\frac{q(z)}{p(z|x, \theta)} \right)$$

general EM

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

- To prove the above claim, first note that

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) = \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) + \ln p(\mathbf{X}|\boldsymbol{\theta}) \quad \leftarrow \text{since } p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) \cdot p(\mathbf{X}|\boldsymbol{\theta})$$

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

$$= \sum_{\mathbf{z}} q(\mathbf{z}) \left(\ln p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) + \ln p(\mathbf{x}|\boldsymbol{\theta}) - \ln q(\mathbf{z}) \right)$$

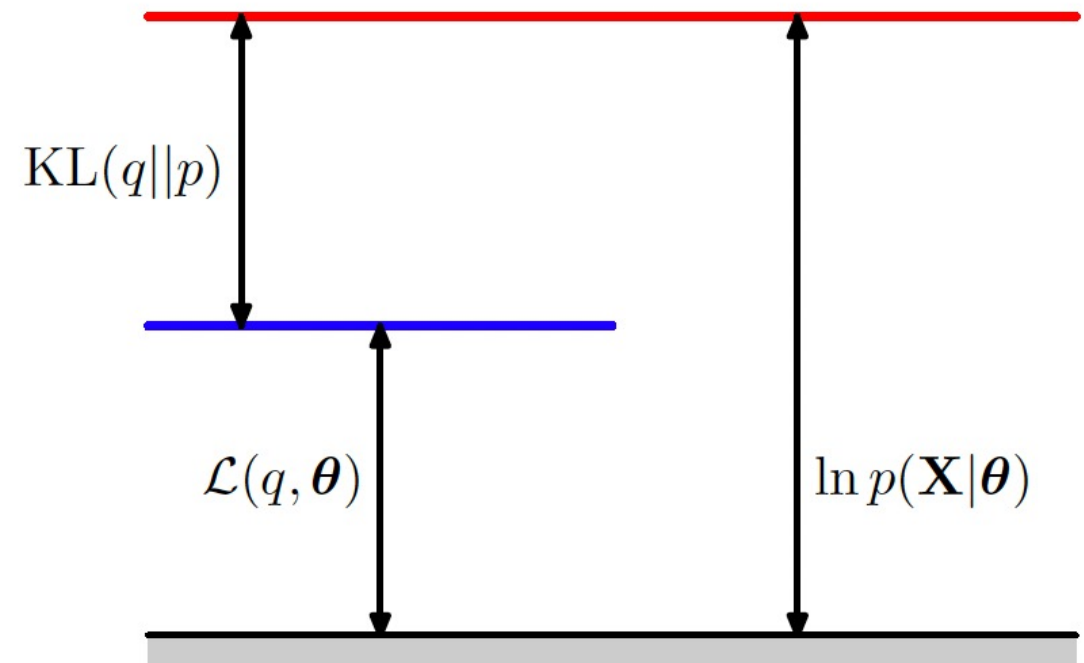
$$= \sum_{\mathbf{z}} q(\mathbf{z}) \ln \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{z})} + \left(\sum_{\mathbf{z}} q(\mathbf{z}) \right) \ln p(\mathbf{x}|\boldsymbol{\theta})$$

$$= -\text{KL}(q||p) + \ln p(\mathbf{x}|\boldsymbol{\theta})$$

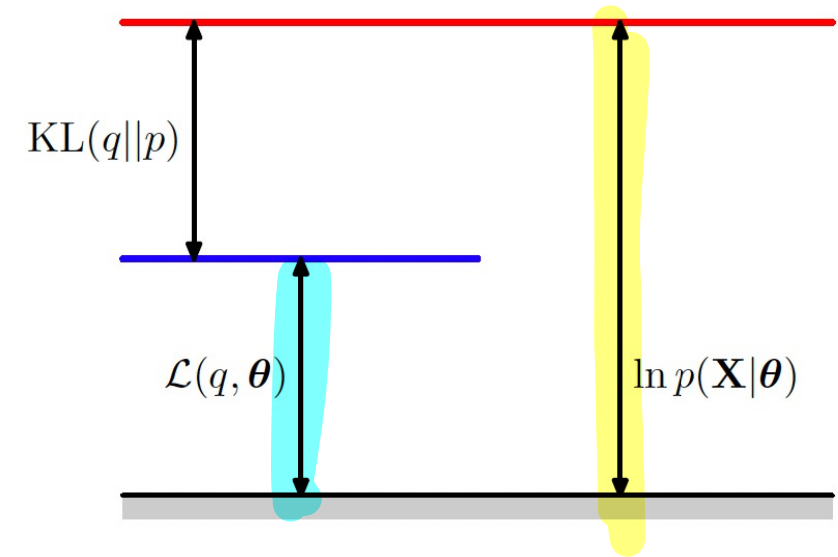


general EM

- Now we have proved that $\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q||p)$
- Fact: $\text{KL}(q || p) \geq 0$
- Therefore, $\mathcal{L}(q, \boldsymbol{\theta})$ is a lower bound of $\ln p(\mathbf{X}|\boldsymbol{\theta})$



general EM



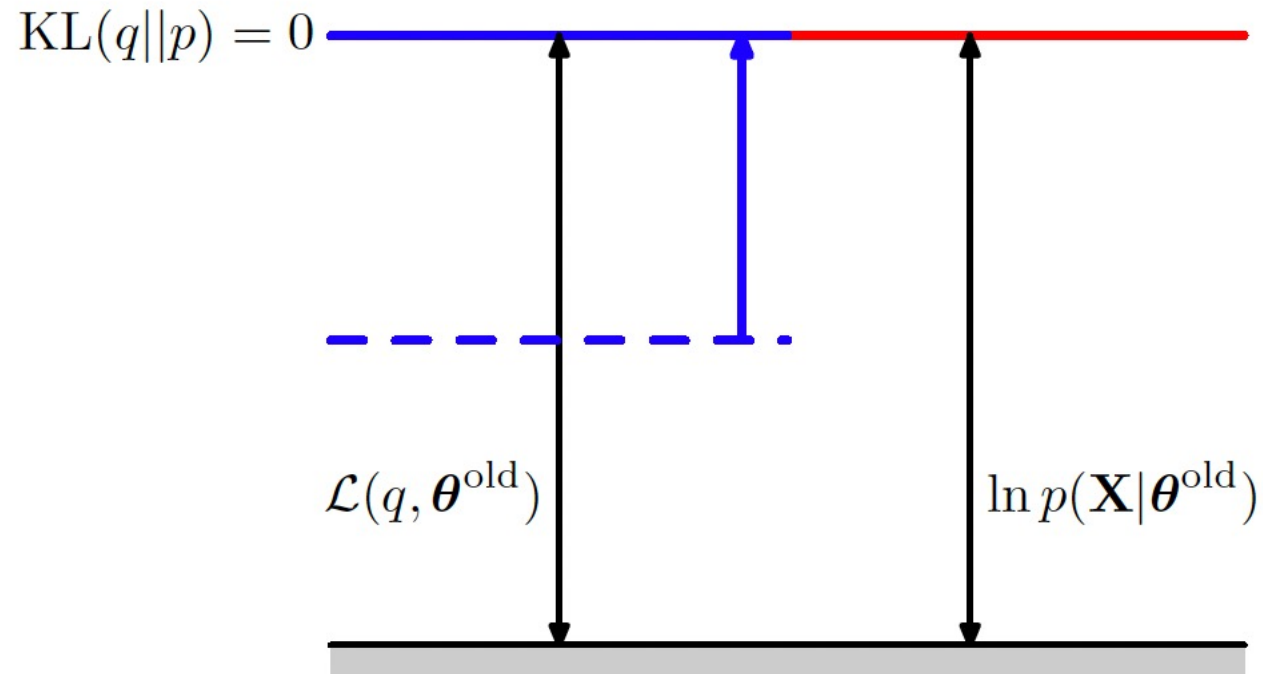
(E-step): expectation

maximize $\mathcal{L}(q, \theta^{\text{old}})$ with respect to q

This q will be pushed to $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$!!!

general EM

Illustration of the E step of the EM algorithm. The q distribution is set equal to the posterior distribution for the current parameter values θ^{old} , causing the lower bound to move up to the same value as the log likelihood function, with the KL divergence vanishing.



general EM

(E-step): expectation

maximize $\mathcal{L}(q, \theta^{\text{old}})$ with respect to q

(M-step): maximization

fix q and maximize $\mathcal{L}(q, \theta)$ with respect to θ

general EM

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

- After the E Step, $q(\mathbf{Z}) = p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$, and therefore the lower bound takes the form

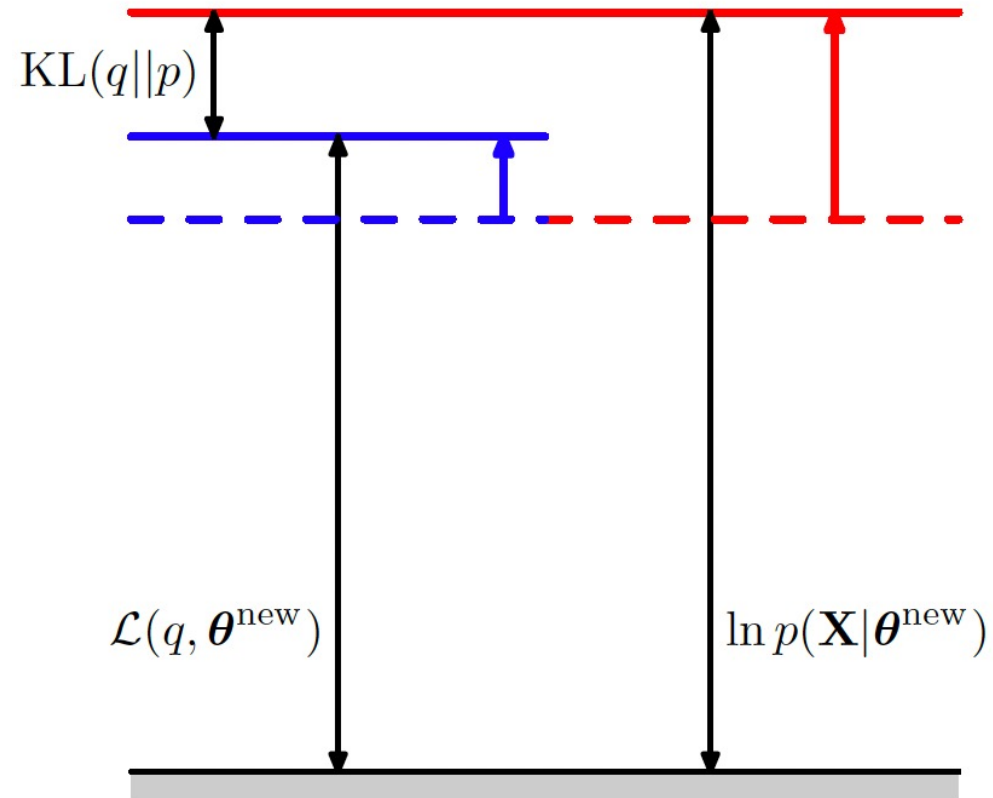
$$\mathcal{L}(q, \theta) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \theta) - \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$$

$$= Q(\theta, \theta^{\text{old}}) + \text{const}$$

$$\mathbb{E}_{p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})} \ln p(\mathbf{X}, \mathbf{Z} | \theta)$$

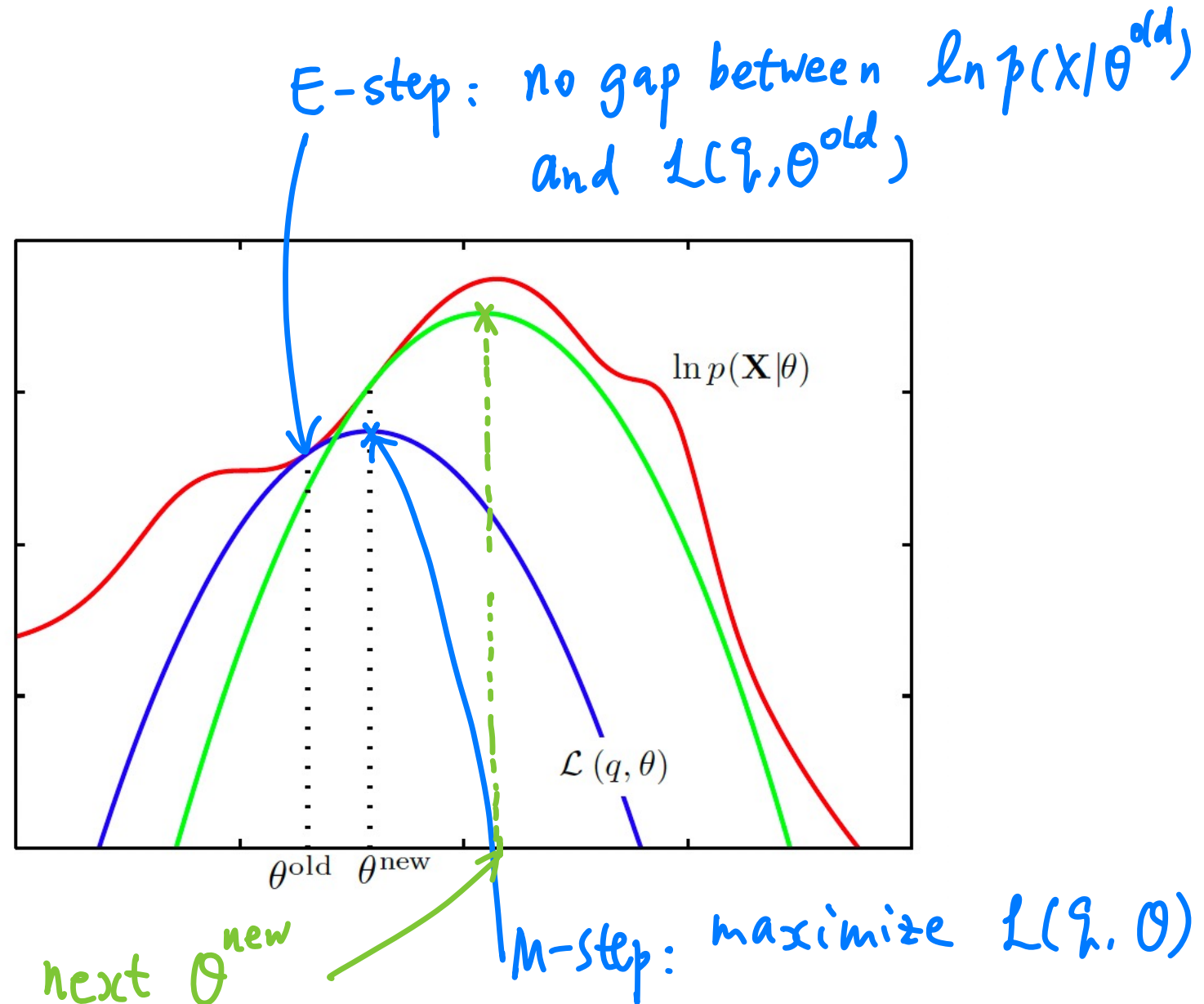
general EM

Illustration of the M step of the EM algorithm. The distribution $q(\mathbf{Z})$ is held fixed and the lower bound $\mathcal{L}(q, \theta)$ is maximized with respect to the parameter vector θ to give a revised value θ^{new} . Because the KL divergence is nonnegative, this causes the log likelihood $\ln p(\mathbf{X}|\theta)$ to increase by at least as much as the lower bound does.



general EM

The EM algorithm involves alternately computing a lower bound on the log likelihood for the current parameter values and then maximizing this bound to obtain the new parameter values. See the text for a full discussion.



spectral clustering

no need to have exact locations

- For both K-means and EM, we use the location of the data instances $\{\mathbf{x}_n\}_{n=1}^N$
- However, in a distance-based method, we make decision only based on the relative locations
- Goal of this section: we want to get low dimensional features disregarding the original locations

guide for projection

- Given $x_n \in \mathbb{R}^d$, $n = 1, \dots, N$.
- Suppose we want to look at low-dimension projection $z_n \in \mathbb{R}^k$, $n = 1, \dots, N$.
- Suppose further that we are given a **similarity score** for each pair of instances
 - $w_{nm} = w_{mn}$ is the similarity between x_n and x_m (assume $w_{nn} = 0$)
- The goal is to **minimize** $\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \|z_n - z_m\|^2 w_{nm}$

case $k = 1: z_n \in \mathbb{R}$

$$\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N |z_n - z_m|^2 w_{nm} = ? \quad \left(\text{let } d_n := \sum_{m=1}^N w_{nm} = \sum_{m=1}^N w_{mn} \right)$$

$$= \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N z_n^2 w_{nm} + z_m^2 w_{nm} - 2 z_n z_m w_{nm}$$

$$= \frac{1}{2} \left(\sum_{n=1}^N z_n^2 \underbrace{\sum_{m=1}^N w_{nm}}_{d_n} + \sum_{m=1}^N z_m^2 \underbrace{\sum_{n=1}^N w_{nm}}_{d_m} - 2 \sum_{n=1}^N \sum_{m=1}^N z_n z_m w_{nm} \right)$$

$$= \sum_{n=1}^N z_n d_n z_n - \sum_{n=1}^N \sum_{m=1}^N z_n w_{nm} z_m$$

$$= z^T D z - z^T W z \quad \text{where } D = \text{diag}(d_1, \dots, d_N), \quad W = (w_{nm})_{n,m=1}^N$$

case $k = 1: z_n \in \mathbb{R}$

- Now that $\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N |z_n - z_m|^2 w_{nm} = z^T (\mathbf{D} - \mathbf{W}) z$ where
 - \mathbf{W} is the symmetric matrix whose (n, m) -th entry is w_{nm} (let's define $w_{nn} = 0$ for each n)
 - \mathbf{D} is the diagonal matrix with $D_{nn} = \sum_{m \neq n} w_{nm}$
 - Let $\mathbf{L} = \mathbf{D} - \mathbf{W}$, we call \mathbf{L} the graph Laplacian of the similarity graph of our data. *(just a remark)*

eigenvector of Laplacian

- Now we need to minimize $\mathbf{z}^T \mathbf{L} \mathbf{z}$.
- $\mathbf{z} = (z_1, \dots, z_N)^T$ should be an eigenvector of \mathbf{L} .
 - Note that there is a trivial eigenvector $\mathbf{c} = \frac{1}{\sqrt{N}} (1, 1, \dots, 1)^T$
 - $\mathbf{L} \mathbf{c} = \mathbf{D} \mathbf{c} - \mathbf{W} \mathbf{c}$, whose n -th entry is $\frac{1}{\sqrt{N}} \left[D_{nn} - \sum_{m \neq n} w_{nm} \right] = 0$
 - but this is not interesting because this means z_n is the same for different data point \mathbf{x}_n .
 - We should have some \mathbf{z} that is orthogonal to the above trivial choice.
 - Therefore, we take \mathbf{z} to be the eigenvector corresponding to the second smallest eigenvalue.

Laplacian eigenmap

- The map that maps $\{\mathbf{x}_n\}_{n=1}^N$ to $\{\mathbf{z}_n\}_{n=1}^N$ is called the **Laplacian eigenmap** where $\mathbf{z} = (z_1, \dots, z_N)^T$ is the eigenvector corresponding to the second smallest eigenvalue of L .
- In case $k > 1$ (more than one feature in \mathbf{z}), we need to take the next features in \mathbf{z} to be the eigenvectors corresponding to the next smallest eigenvalues.
- Specifically, we will get a feature matrix \mathbf{Z} , whose size is N -by- k , and whose k columns are given by the eigenvectors corresponding to **the 2nd to the $(k + 1)$ -th smallest eigenvalues** of L .

spectral clustering

- Once we have the features $\{\mathbf{z}_n\}_{n=1}^N$, we can do clustering on the features instead of on the original data points $\{\mathbf{x}_n\}_{n=1}^N$.
- For instance, we can apply K-means on $\{\mathbf{z}_n\}_{n=1}^N$.

normalized Laplacian

- In practice, instead of working with the Laplacian $L = D - W$, we may want to work with a normalized version of Laplacian, e.g.
 - $L_{\text{rw}} = I - D^{-1}W$
 - $L_{\text{sym}} = I - D^{-1/2}WD^{-1/2}$

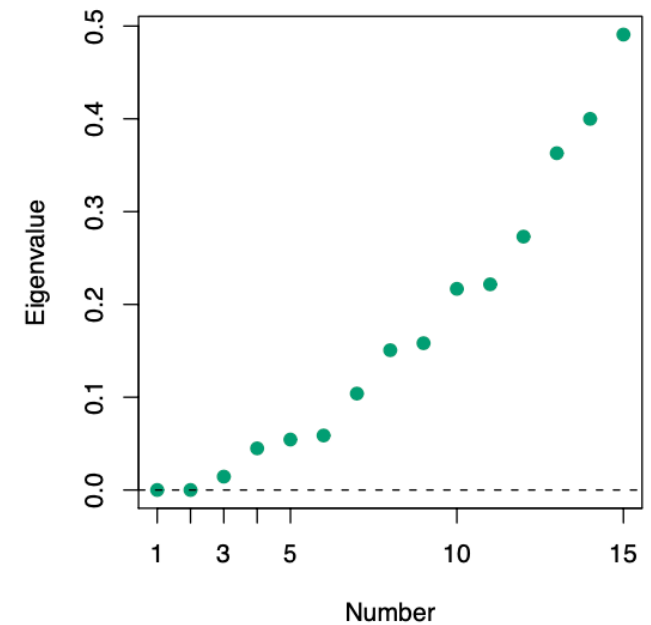
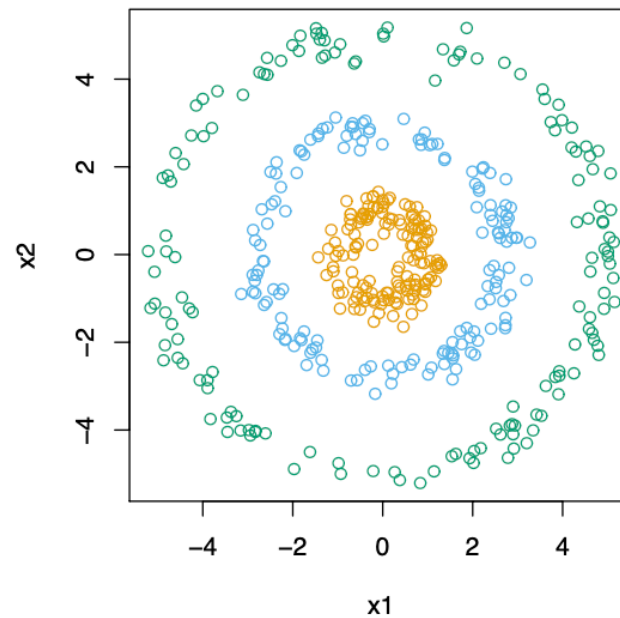
choice of similarity score

- We only consider similarities locally (in a neighborhood):
 - either set $w_{nm} = 0$ for $\| \mathbf{x}_n - \mathbf{x}_m \| > \epsilon$ for some preset threshold ϵ ;
 - or using k-NN: $w_{nm} \neq 0$ if only if \mathbf{x}_m is among the k nearest neighbors of \mathbf{x}_n or vice versa.

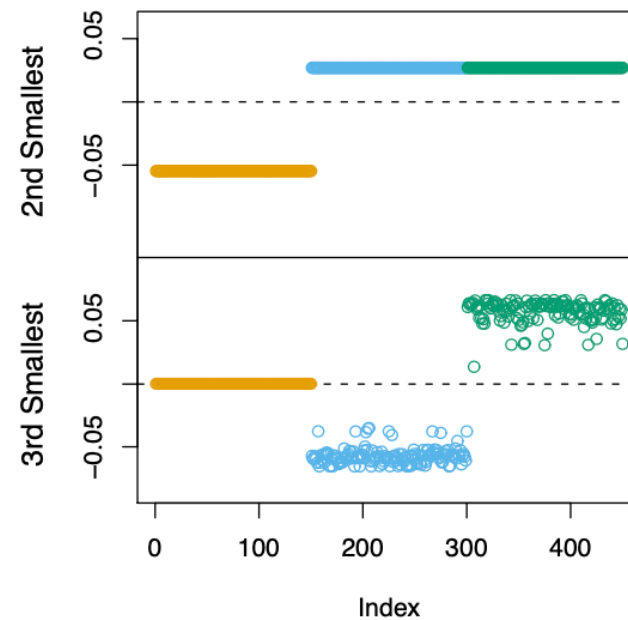
- For $w_{nm} \neq 0$, one possible similarity score is to set

$$w_{nm} = \exp \left(- \frac{\| \mathbf{x}_n - \mathbf{x}_m \|^2}{2\sigma^2} \right)$$

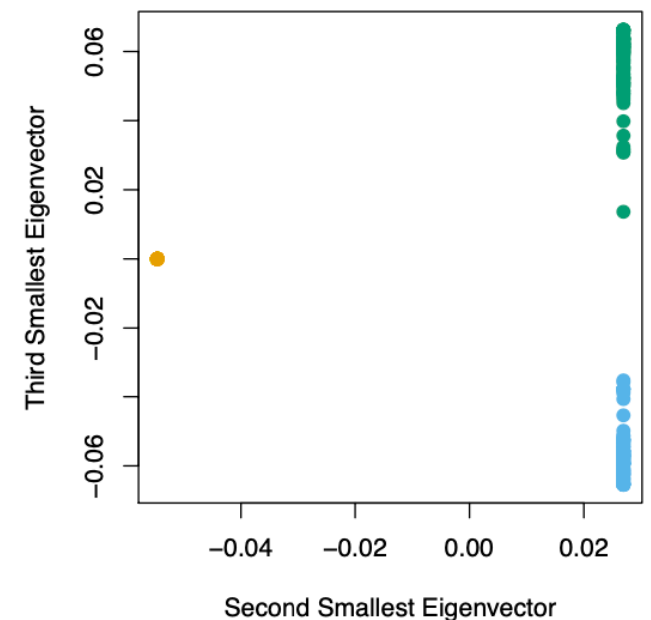
example: three
concentric
clusters



Eigenvectors



Spectral Clustering





Questions?

Reference

- *K-means:*
 - [Al] Ch.7.3
 - [HaTF] Ch.13.2.1
 - [Bi] Ch.9.1
- *EM:*
 - [Al] Ch.7.2, 7.4
 - [HaTF] Ch.13.2.3
 - [Bi] Ch.9.2-9.4
- *Spectral clustering:*
 - [Al] Ch.6.12 7.7
 - [HaTF] Ch.14.5.3