

Latent Variables

- A system with observed data \mathbf{X}
 - may be far easier to understand in terms of additional variables \mathbf{Z} corresponding to \mathbf{X} ,
 - but they are not observed (**latent**).
- For example, in a mixture of Gaussians,
 - For a single sample \mathbf{x} , the latent variable \mathbf{z} specifies which Gaussian generated the sample \mathbf{x} .
 - The *responsibility* is the **posterior** $p(\mathbf{z}|\mathbf{x})$.

Latent Variables

- A system with observed variables \mathbf{X}
 - may be easier to understand with latent variables \mathbf{Z} , but they are not observed (**latent**).
- Notations:
 - We denote the set of all observed data by \mathbf{X} , in which the n^{th} row represents \mathbf{x}_n^T .
 - Similarly we denote the set of all latent variables by \mathbf{Z} , with a corresponding row \mathbf{z}_n^T .
 - Note: we use lowercase symbol for single sample (\mathbf{x}), matrix symbol for all data (\mathbf{X}).

Learning a Latent Variable Model

- We find model parameters by maximizing the log-likelihood of observed data $\log p(\mathbf{X} \mid \theta)$.
- If we had complete data $\{\mathbf{X}, \mathbf{Z}\}$, we could easily maximize the *complete* data likelihood $p(\mathbf{X}, \mathbf{Z} \mid \theta)$.
- Unfortunately, with *incomplete* data (\mathbf{X} only), we must marginalize over \mathbf{Z} , so

$$\log p(\mathbf{X} \mid \theta) = \log \left[\sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} \mid \theta) \right]$$

(the sum inside the log makes it hard.)

The EM Algorithm in General

- Expectation-Maximization (EM) is a general recipe for finding the parameters that maximize the (log-) likelihood of *latent* variable models
- To find a parameter θ that maximizes the likelihood $p(\mathbf{X} \mid \theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} \mid \theta)$, the EM algorithm first introduces a new (variable) distribution $q(\mathbf{Z})$ over the latent variables.
- A lower bound $\mathcal{L}(q, \theta)$ for the log-likelihood $\log p(\mathbf{X} \mid \theta)$ is established based on q and θ .
- Then, $q(\mathbf{Z})$ and θ are alternately updated (keeping the other fixed) so that $\mathcal{L}(q, \theta)$ is maximized (similar to coordinate ascent) until convergence.

The EM Algorithm in General

- Our goal is to maximize $p(\mathbf{X} \mid \theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} \mid \theta)$
- For ***any distribution*** $q(\mathbf{Z})$ over latent variables:

$$\begin{aligned}\log p(\mathbf{X} \mid \theta) &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X} \mid \theta) \\&= \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z} \mid \theta)}{p(\mathbf{Z} \mid \mathbf{X}, \theta)} \\&= \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z} \mid \theta)}{q(\mathbf{Z})} \frac{q(\mathbf{Z})}{p(\mathbf{Z} \mid \mathbf{X}, \theta)} \\&= \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z} \mid \theta)}{q(\mathbf{Z})} + \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{q(\mathbf{Z})}{p(\mathbf{Z} \mid \mathbf{X}, \theta)} \\&= \mathcal{L}(q, \theta) + KL(q(\mathbf{Z}) \parallel p(\mathbf{Z} \mid \mathbf{X}, \theta)) \\&\geq \mathcal{L}(q, \theta)\end{aligned}$$

Note: KL Divergence

Let p and q be probability distributions of a random variable Z .

$$\begin{aligned} KL(q \parallel p) &= \mathbb{E}_{z \sim q(z)} \left[\log \frac{q(z)}{p(z)} \right] = \sum_z q(z) \log \frac{q(z)}{p(z)} \\ &= - \sum_z q(z) \log p(z) + \sum_z q(z) \log q(z) \end{aligned}$$

This is one way to measure the **dissimilarity** of two probability distributions.

Remarks: (note: the first can be proved using Jensen's inequality)

- $KL(q \parallel p) \geq 0$, with equality iff $p = q$.
- $KL(q \parallel p) \neq KL(p \parallel q)$ in general

Background note: Jensen's Inequality

- If f is convex, then for any θ_i s.t. $0 \leq \theta_i \leq 1$ ($\forall i$),
$$\theta_1 + \theta_2 + \cdots + \theta_k = 1$$

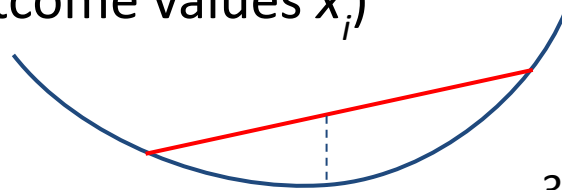
$$f(\theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_k x_k) \leq \theta_1 f(x_1) + \cdots + \theta_k f(x_k)$$

- It can be seen as a generalization of the definition of convex function:

$$f \text{ is convex} \iff f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y) \text{ for all } 0 \leq \theta \leq 1$$

- Jensen's inequality can be written in expectation form
(think of θ_i as probability mass for different outcome values x_i)

$$f(\mathbb{E}[x]) \leq \mathbb{E}[f(x)]$$



Background note: Jensen's Inequality

- If f is convex, then for any θ_i s.t. $0 \leq \theta_i \leq 1$ ($\forall i$),
$$\theta_1 + \theta_2 + \cdots + \theta_k = 1$$

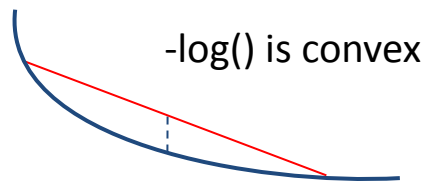
$$f(\theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_k x_k) \leq \theta_1 f(x_1) + \cdots + \theta_k f(x_k)$$

- Jensen's inequality can be written in expectation form

$$f(\mathbb{E}[x]) \leq \mathbb{E}[f(x)]$$

- To show $KL(q \parallel p)$ is non-negative for any p, q ,
plug in $f(\dots) = -\log(\dots)$ and the following:

$$\theta_i = q(z), x_i = \frac{p(z)}{q(z)}$$



$$-\log(\mathbb{E}[x]) \leq \mathbb{E}[-\log(x)]$$

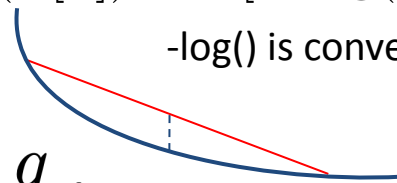
Non-negativity of KL divergence

- Jensen's inequality can be written in expectation form for a convex function f

$$f(\mathbb{E}[x]) \leq \mathbb{E}[f(x)]$$

$$-\log(\mathbb{E}[x]) \leq \mathbb{E}[-\log(x)]$$

$-\log()$ is convex



- To show $KL(q \parallel p)$ is non-negative for any p, q , plug in $f(\dots) = -\log(\dots)$ and the following: $\theta_i = q(z)$, $x_i = \frac{p(z)}{q(z)}$

$$\begin{aligned} KL(q \parallel p) &= \sum_z q(z) \log\left(\frac{q(z)}{p(z)}\right) \\ &= \sum_z q(z) \left(-\log\left(\frac{p(z)}{q(z)}\right)\right) \\ &\geq -\log\left(\underbrace{\sum_z q(z) \frac{p(z)}{q(z)}}_{=\sum_z p(z)=1}\right) \\ &= 0 \end{aligned}$$



Jensen's inequality for $-\log()$:

$$-\log(\mathbb{E}[x]) \leq \mathbb{E}[-\log(x)]$$

i.e., plugin

$$\begin{aligned} -\log\left(\sum_i \theta_i x_i\right) &\leq \sum_i \theta_i (-\log(x_i)) \\ \text{with } \theta_i &= q(z), x_i = \frac{p(z)}{q(z)} \end{aligned}$$

The EM Algorithm in a nutshell

- We have shown that: [variational lower bound]

$$\begin{aligned}\log p(\mathbf{X} \mid \theta) &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z} \mid \theta)}{q(\mathbf{Z})} + \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{q(\mathbf{Z})}{p(\mathbf{Z} \mid \mathbf{X}, \theta)} \\ &= \mathcal{L}(q, \theta) + KL(q(\mathbf{Z}) \parallel p(\mathbf{Z} \mid \mathbf{X}, \theta)) \\ &\geq \mathcal{L}(q, \theta) \quad \text{Evidence Lower bound (ELBO) or variational lower bound}\end{aligned}$$

with equality holding if and only if $q(\mathbf{Z}) = p(\mathbf{Z} \mid \mathbf{X}, \theta)$

- **EM algorithm:**

* E: expectation

* M: maximization

Repeat alternating optimization until convergence:

- E-step: for fixed θ , find q that maximizes $\mathcal{L}(q, \theta)$
- M-step: for fixed q , find θ that maximizes $\mathcal{L}(q, \theta)$

The EM Algorithm: E-step

- We have shown that: [variational lower bound]

$$\begin{aligned}\log p(\mathbf{X} \mid \theta) &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z} \mid \theta)}{q(\mathbf{Z})} + \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{q(\mathbf{Z})}{p(\mathbf{Z} \mid \mathbf{X}, \theta)} \\ &= \mathcal{L}(q, \theta) + KL(q(\mathbf{Z}) \parallel p(\mathbf{Z} \mid \mathbf{X}, \theta)) \\ &\geq \mathcal{L}(q, \theta) \quad \text{Evidence Lower bound (ELBO) or variational lower bound}\end{aligned}$$

with equality holding if and only if $q(\mathbf{Z}) = p(\mathbf{Z} \mid \mathbf{X}, \theta)$

- **(E-step)** For a fixed θ , which q maximizes $\mathcal{L}(q, \theta)$?
 $\Rightarrow p(\mathbf{Z} \mid \mathbf{X}, \theta)$, because all other q would make $\mathcal{L}(q, \theta)$ strictly less than $\log p(\mathbf{X} \mid \theta)$

The EM Algorithm: M-step

- We also note that for a fixed q , the $\mathcal{L}(q, \theta)$ term can be decomposed into two terms:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z} \mid \theta)}{q(\mathbf{Z})} \\ &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z} \mid \theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})\end{aligned}$$

- (1) A weighted sum of $\log p(\mathbf{X}, \mathbf{Z} \mid \theta)$.

This is tractable and can be optimized w.r.t θ

- (2) Entropy of $q(\mathbf{Z})$ which is independent of θ since q is fixed.

- **(M-step)** Thus, when q is fixed, we can find θ that maximizes $\mathcal{L}(q, \theta)$.

The EM Algorithm: summary

- Initialize parameters θ randomly
- Repeat until convergence:
(optimize $\mathcal{L}(q, \theta)$ w.r.t. q and θ alternately.)
 - “E-step”: Set $q(\mathbf{Z}) = p(\mathbf{Z} \mid \mathbf{X}, \theta)$ compute posterior \rightarrow optimal $q(\mathbf{Z})!$
 - “M-step”: Update θ via the following maximization

$$\operatorname{argmax}_{\theta} \mathcal{L}(q, \theta) = \operatorname{argmax}_{\theta} \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z} \mid \theta)$$

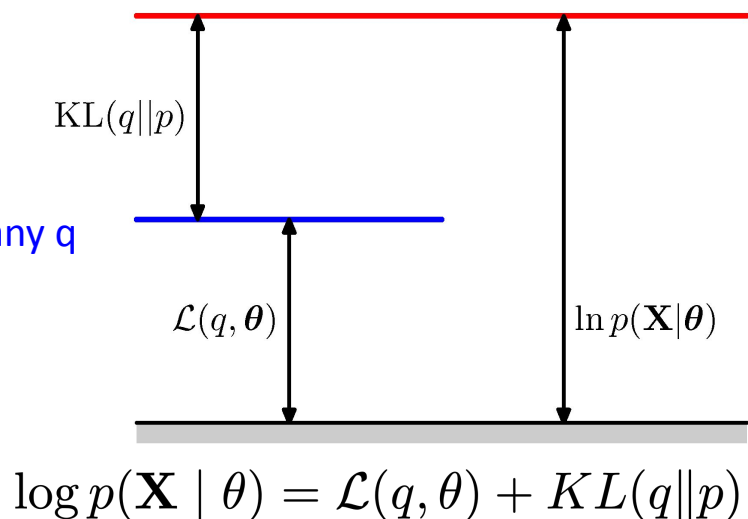
use $q(\mathbf{Z})$ as (fractional) pseudo-counts and maximize the “data completion” log-likelihood

- Note we have assumed that $p(\mathbf{Z} \mid \mathbf{X}, \theta)$ is tractable (i.e., find exact posterior $p(\mathbf{Z} \mid \mathbf{X}, \theta)$). Q. What if it is not?

Visualize the Decomposition

- Note: $KL(q||p) \geq 0$
 - with equality only when $q=p$.

for any q



- Thus, $\mathcal{L}(q, \theta)$ is a lower bound on $\log p(\mathbf{X} | \theta)$

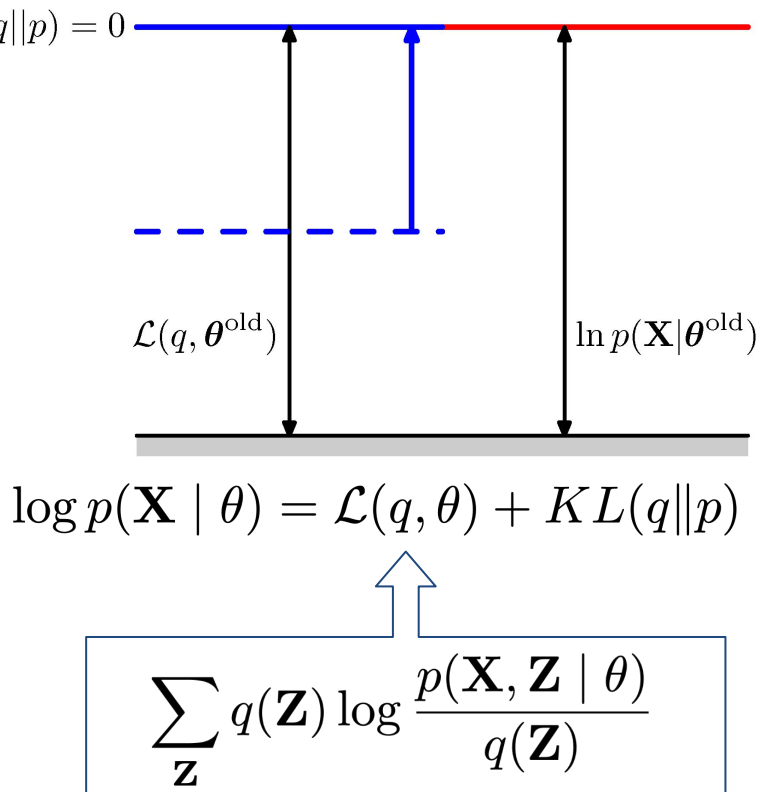
which EM tries to maximize.

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q(\mathbf{Z})}$$

Visualize the E-Step

- E-step: for fixed θ , find q that maximizes $\mathcal{L}(q, \theta)$

for $q(\mathbf{Z}) = P(\mathbf{Z}|\mathbf{X}, \theta)$



- E-Step changes $q(\mathbf{Z})$ to maximize $\mathcal{L}(q, \theta)$

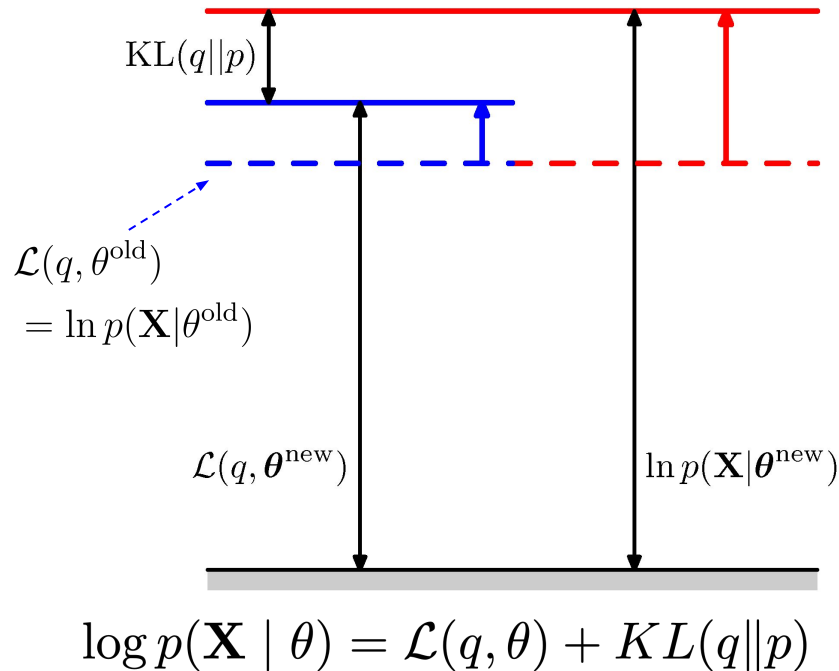
- So maximized when

$$KL(q||p) = 0$$

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

Visualize the M-Step

- M-step: for fixed q , find θ that maximizes $\mathcal{L}(q, \theta)$



- Holding $q(\mathbf{Z})$ constant; increase $\mathcal{L}(q, \theta)$
- Updating θ will make $\log p(\mathbf{X} | \theta)$ increase!

$$\ln p(\mathbf{X} | \theta^{\text{new}}) \geq \ln p(\mathbf{X} | \theta^{\text{old}})$$
- But now $p \neq q$
- so $KL(q || p) > 0$

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q(\mathbf{Z})}$$

The EM Algorithm: Multiple data-points

- Variational lower bound for a single example \mathbf{x} :

$$\begin{aligned}\log p(\mathbf{x}|\theta) &= \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{z}, \mathbf{x}|\theta)}{q(\mathbf{z})} + KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}, \theta)) \\ &\geq \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{z}, \mathbf{x}|\theta)}{q(\mathbf{z})}\end{aligned}$$

- Lower bound on the log-likelihood of the *entire* training data $\mathcal{D} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$:

$$\begin{aligned}\log p(\mathcal{D}|\theta) &= \sum_n \log p(\mathbf{x}^{(n)}|\theta) = \sum_n \sum_{\mathbf{z}} q^{(n)}(\mathbf{z}) \log \frac{p(\mathbf{z}, \mathbf{x}^{(n)}|\theta)}{q^{(n)}(\mathbf{z})} + \sum_n KL(q^{(n)}(\mathbf{z})||p(\mathbf{z}|\mathbf{x}^{(n)}, \theta)) \\ &\geq \sum_n \sum_{\mathbf{z}} q^{(n)}(\mathbf{z}) \log \frac{p(\mathbf{z}, \mathbf{x}^{(n)}|\theta)}{q^{(n)}(\mathbf{z})}\end{aligned}$$

Note that different $q^{(n)}$ is used for each n

The EM Algorithm: Multiple data-points

$$\begin{aligned}\log p(\mathcal{D}|\theta) &= \sum_n \log p(\mathbf{x}^{(n)}|\theta) = \sum_n \sum_{\mathbf{z}} q^{(n)}(\mathbf{z}) \log \frac{p(\mathbf{z}, \mathbf{x}^{(n)}|\theta)}{q^{(n)}(\mathbf{z})} + \sum_n KL(q^{(n)}(\mathbf{z})||p(\mathbf{z}|\mathbf{x}^{(n)}, \theta)) \\ &\geq \sum_n \sum_{\mathbf{z}} q^{(n)}(\mathbf{z}) \log \frac{p(\mathbf{z}, \mathbf{x}^{(n)}|\theta)}{q^{(n)}(\mathbf{z})}\end{aligned}$$

- Initialize random parameters θ
- Repeat until convergence:
 - **“E-step”**: Set $q^{(n)}(\mathbf{z}) = p(\mathbf{z} | \mathbf{x}^{(n)}, \theta)$,
for each training sample n .
 - **“M-step”**: Update θ via the following maximization:

$$\arg \max_{\theta} \sum_n \sum_{\mathbf{z}} q^{(n)}(\mathbf{z}) \log p(\mathbf{z}, \mathbf{x}^{(n)} | \theta)$$