

DSC291: Machine Learning with Few Labels

Unsupervised Learning

Zhiting Hu

Lecture 14, May 3, 2024

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Recap: Expectation Maximization (EM)

- Supervised MLE is easy: $\max_{\theta} \ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z}|\theta)$
 - Observe both \mathbf{x} and \mathbf{z}
- Unsupervised MLE is hard: $\max_{\theta} \ell(\theta; \mathbf{x}) = \log p(\mathbf{x}|\theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\theta)$
 - Observe only \mathbf{x}
- EM, intuitively:

E-step: $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x}, \theta)$

*We don't actually observe q ,
let's estimate it*

M-step: $\max_{\theta} \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{z}|\theta)]$

*Let's "pretend" we also observe
 \mathbf{z} (its distribution)*

Recap: Expectation Maximization (EM)

- Supervised MLE is easy: $\max_{\theta} \ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z}|\theta)$
 - Observe both \mathbf{x} and \mathbf{z}
- Unsupervised MLE is hard: $\max_{\theta} \ell(\theta; \mathbf{x}) = \log p(\mathbf{x}|\theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\theta)$
 - Observe only \mathbf{x}
- EM, intuitively:

→ E-step: $q^{t+1}(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x}, \theta^t)$

*We don't actually observe q ,
let's estimate it*

↳ M-step: $\max_{\theta} \mathbb{E}_{q^{t+1}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{z}|\theta)]$

*Let's "pretend" we also observe
 \mathbf{z} (its distribution)*

*This is an iterative
process*

Recap: Expectation Maximization (EM)

- The EM algorithm is coordinate-descent on $F(q, \theta)$

- E-step: $q^{t+1} = \arg \min_q F(q, \theta^t) = \underline{p(\mathbf{z}|\mathbf{x}, \theta^t)}$

- the posterior distribution over the latent variables given the data and the current parameters

- M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta) = \underline{\operatorname{argmax}_{\theta} \mathbb{E}_{q^{t+1}}[\log p(\mathbf{x}, \mathbf{z}|\theta)]}$

ELBO

$$\begin{aligned} \ell(\theta; \mathbf{x}) &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta)) \\ &= \underline{-F(q, \theta)} + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta)) \end{aligned}$$

Example: Gaussian Mixture Models (GMMs)

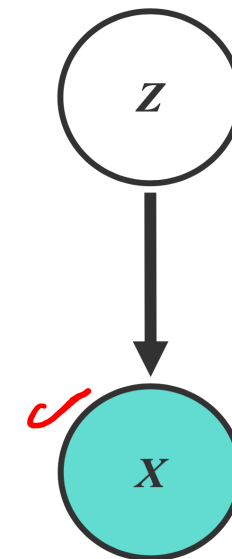
- Consider a mixture of K Gaussian components:

- Z is a latent class indicator vector:

$$Z \in \{0, 1, \dots, K-1\}$$

$$p(z_n) = \text{multi}(z_n : \pi) = \prod_k (\pi_k)^{z_n^k}$$

- X is a conditional Gaussian variable with a class-specific mean/covariance



$$p(x_n | z_n^k = 1, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

- The likelihood of a sample:

$$= \sum_k p(z_n = k | \pi, \mu, \Sigma)$$

$$p(x_n | \mu, \Sigma) = \sum_k p(z_n^k = 1 | \pi) p(x_n | z_n^k = 1, \mu, \Sigma)$$

$$= \sum_{z_n} \prod_k \left((\pi_k)^{z_n^k} N(x_n : \mu_k, \Sigma_k)^{z_n^k} \right) = \sum_k \pi_k N(x_n | \mu_k, \Sigma_k)$$

$$\log p(x_n | \mu, \Sigma)$$

mixture proportion mixture component

$$p(z|x, \pi, \mu, \Sigma) = p(z|x, \pi, \mu, \Sigma)$$

Example: Gaussian Mixture Models (GMMs)

- E-step: computing the posterior of z_n given the current estimate of the parameters (i.e., π, μ, Σ)

\downarrow
 $p(z^k|x, \theta)$
 $k \in \{0, 1, \dots, K-1\}$

Bayes Theorem \star

$$p(z|x) = \frac{p(z)p(x|z)}{p(x)}$$

$z \in \{0, 1, \dots, K-1\}$

$$p(z^k = 1 | \mathbf{x}) = \frac{p(z^k = 1)p(\mathbf{x} | z^k = 1)}{p(\mathbf{x})} = \frac{p(z^k = 1)p(\mathbf{x} | z^k = 1)}{\sum_j p(z^j = 1)p(\mathbf{x} | z^j = 1)}$$

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

Example: Gaussian Mixture Models (GMMs)

- M-step: computing the parameters given the current estimate of z_n
 - Once we have $q^{t+1}(z^k|x) = p(z^k|x, \theta^t) = \gamma^k$, we can compute the expected likelihood:

$$\begin{aligned} \Rightarrow \theta^{t+1} &= \operatorname{argmax}_{\theta} \sum_k q^{t+1}(z^k = 1|x) \log p(x, z^k = 1|\theta) = \mathbb{E}_{q^{t+1}}[\log P(x, z^k|\theta)] \\ &= \mathbb{E}_{q^{t+1}} [\log (p(\mathbf{x}, z | \theta))] \\ &= \sum_k \gamma_k (\log p(z^k = 1|\theta) + \log P(\mathbf{x} | z^k = 1, \theta)) \\ &= \sum_k \gamma_k \log \pi_k + \sum_k \gamma_k \log \mathcal{N}(\mathbf{x}; \mu_k, \Sigma_k) \end{aligned}$$

$\theta = (\pi, \mu, \Sigma)$

- We need to fit K Gaussians, just need to weight examples by γ_k

$\{\mu^k, \Sigma^k\}_{k=1}^K$

Example: Gaussian Mixture Models (GMMs)

- M-step: computing the parameters given the current estimate of z_n

$$\pi_k^* = \arg \max \langle l_c(\boldsymbol{\theta}) \rangle,$$

$$\Rightarrow \frac{\partial}{\partial \pi_k} \langle l_c(\boldsymbol{\theta}) \rangle = 0, \forall k, \quad \text{s.t.} \sum_k \pi_k = 1$$

$$\Rightarrow \pi_k^* = \frac{\sum_n \langle z_n^k \rangle_{q^{(t)}}}{N} = \frac{\sum_n \gamma_n^{k(t)}}{N} = \frac{\langle n_k \rangle}{N}$$

$$\mu_k^* = \arg \max \langle l(\boldsymbol{\theta}) \rangle,$$

$$\Rightarrow \mu_k^{(t+1)} = \frac{\sum_n \gamma_n^{k(t)} x_n}{\sum_n \gamma_n^{k(t)}}$$

$$\Sigma_k^* = \arg \max \langle l(\boldsymbol{\theta}) \rangle,$$

$$\Rightarrow \Sigma_k^{(t+1)} = \frac{\sum_n \gamma_n^{k(t)} (x_n - \mu_k^{(t+1)})(x_n - \mu_k^{(t+1)})^T}{\sum_n \gamma_n^{k(t)}}$$

$x_n \rightarrow \gamma_n^k$
 $\pi_k = \frac{\sum_n \gamma_n^k}{N}$

$\gamma_n = (\gamma_n^0, \gamma_n^1, \dots, \gamma_n^{K-1})$

Fact:
 $\frac{\partial \log |A^{-1}|}{\partial A^{-1}} = A^T$
 $\frac{\partial \mathbf{x}^T A \mathbf{x}}{\partial A} = \mathbf{x} \mathbf{x}^T$

$x_n \rightarrow \gamma_n^k$
 $k = 0, \dots, K-1$

EM Algorithm for GMM: Quick Summary

- Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k
- Iterate until convergence:
 - E-step: Evaluate the posterior given current parameters

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

Bayes' rule

γ_k

n

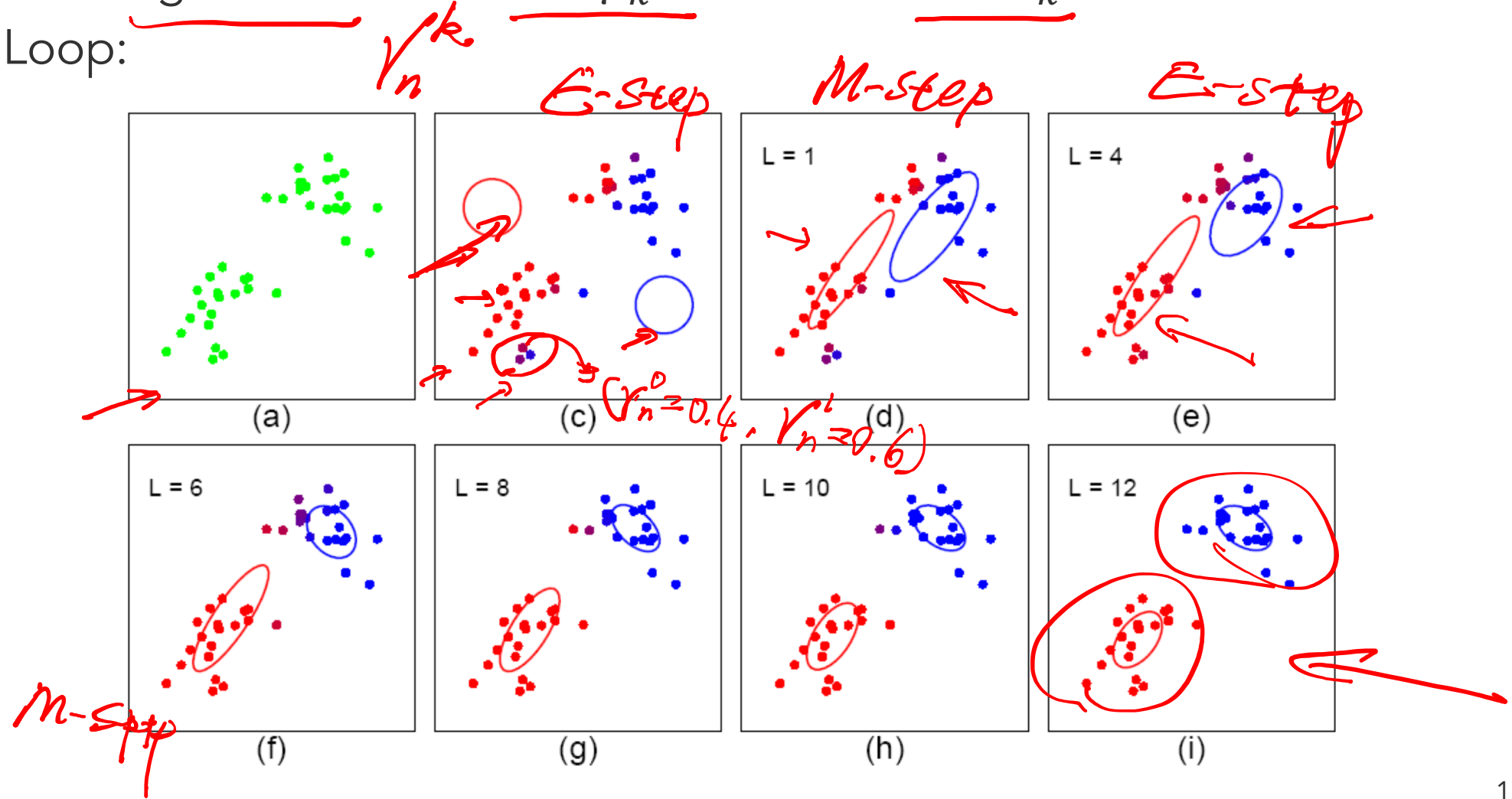
- M-step: Re-estimate the parameters given current posterior

k -means ✓

KNN \rightarrow k -nearest neighbors ✓

Example: Gaussian Mixture Models (GMMs)

- Start: "guess" the centroid μ_k and covariance Σ_k of each of the K clusters
- Loop:



Summary: EM Algorithm

- A way of maximizing likelihood function for latent variable models. Finds MLE of parameters when the original (hard) problem can be broken up into two (easy) pieces
 - Estimate some "missing" or "unobserved" data from observed data and current parameters.
 - Using this "complete" data, find the maximum likelihood parameter estimates.

Summary: EM Algorithm

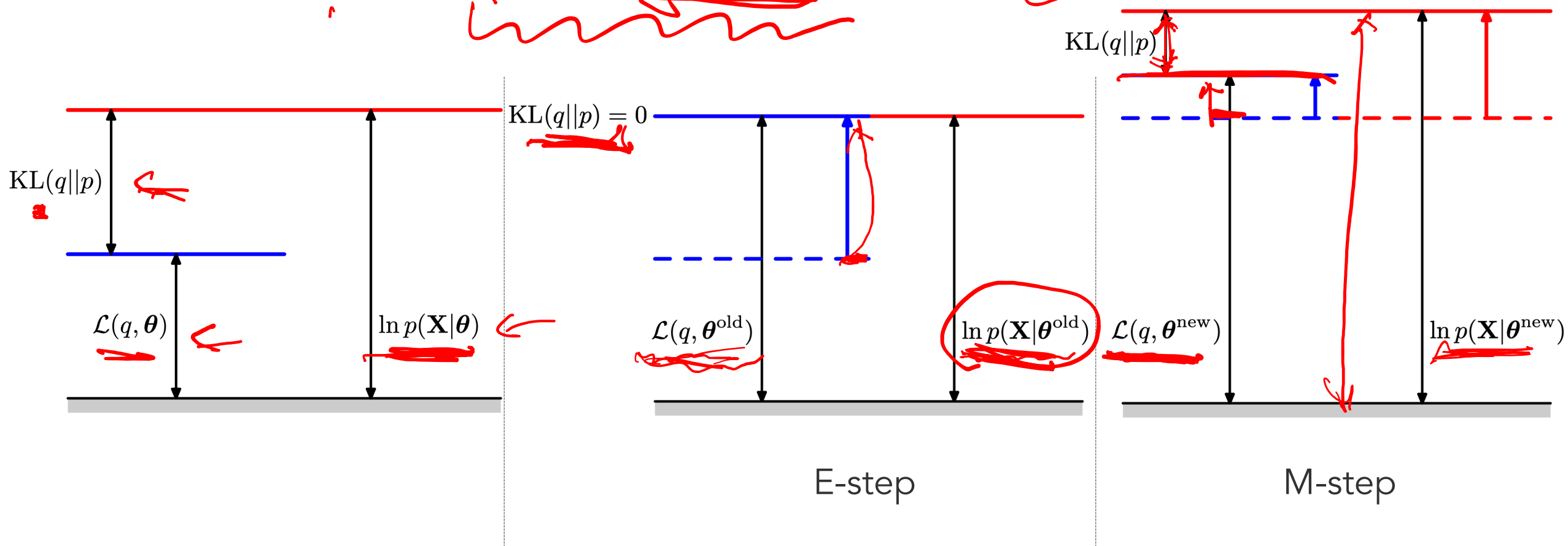
- The EM algorithm is coordinate-decent on $F(q, \theta)$
 - E-step: $q^{t+1} = \arg \min_q F(q, \theta^t) = p(\mathbf{z}|\mathbf{x}, \theta^t)$
 - M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta) = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q^{t+1}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$

$$\begin{aligned} \ell(\theta; \mathbf{x}) &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta)) \\ &= -F(q, \theta) + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta)) \end{aligned}$$

$ELBO \leq \text{Likelihood}$

Each EM iteration guarantees to improve the likelihood

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(z|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, z|\theta)}{q(z|\mathbf{x})} \right] + \text{KL}(q(z|\mathbf{x}) || p(z|\mathbf{x}, \theta))$$



Summary: EM Algorithm

- The EM algorithm is coordinate-decent on $F(q, \theta)$

- E-step: $q^{t+1} = \arg \min_q F(q, \theta^t) = p(\mathbf{z}|\mathbf{x}, \theta^t)$

- M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta) = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q^{t+1}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

$$= -F(q, \theta) + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

- Limitation: need to be able to compute $p(\mathbf{z}|\mathbf{x}, \theta)$, not possible for more complicated models --- solution: Variational inference

Variational Inference

Inference

- Given a model, the goals of inference can include:
 - Computing the likelihood of observed data $p(\mathbf{x}^*)$
 - Computing the marginal distribution over a given subset of variables in the model $p(\mathbf{x}_A)$
 - Computing the conditional distribution over a subsets of nodes given a disjoint subset of nodes $p(\mathbf{x}_A|\mathbf{x}_B)$
 - Computing a mode of the density (for the above distributions) $\operatorname{argmax}_{\mathbf{x}} p(\mathbf{x})$
 -

Variational Inference

- Observed variables \mathbf{x} , latent variables \mathbf{z}
- Variational (Bayesian) inference, a.k.a. **variational Bayes**, is most often used to **approximately** infer the conditional distribution over the latent variables given the observations (and parameters)
 - i.e., the **posterior distribution** over the latent variables

$$p(\mathbf{z}|\mathbf{x}, \theta) = \frac{p(\mathbf{z}, \mathbf{x}|\theta)}{\sum_{\mathbf{z}} p(\mathbf{z}, \mathbf{x}|\theta)}$$

Motivating Example

- Why do we often need to use an approximate inference methods (such as variational Bayes) to compute the posterior distribution?
- It's because we cannot directly compute the posterior distribution for many interesting models
 - I.e. the posterior density is in an intractable form (often involving integrals) which cannot be easily analytically solved.
- As a motivating example, we will try to compute the posterior for a (Bayesian) mixture of Gaussians.

Bayesian mixture of Gaussians

- The mean μ_k is treated as a (latent) random variable

$$\mu_k \sim \mathcal{N}(0, \tau^2) \text{ for } k = 1, \dots, K$$

- For each data $i = 1, \dots, n$

$$z_i \sim \text{Cat}(\pi).$$

$$x_i \sim \mathcal{N}(\mu_{z_i}, \sigma^2).$$

- We have
 - observed variables $x_{1:n}$
 - latent variables $\mu_{1:k}$ and $z_{1:n}$
 - parameters $\{\tau^2, \pi, \sigma^2\}$

- $p(x_{1:n}, z_{1:n}, \mu_{1:k} | \tau^2, \pi, \sigma^2) = \prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i | z_i, \mu_{1:K})$

Bayesian mixture of Gaussians

- We can write the posterior distribution as

$$p(\mu_{1:K}, z_{1:n} | x_{1:n}) = \frac{\prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i | z_i, \mu_{1:K})}{\int_{\mu_{1:K}} \sum_{z_{1:n}} \prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i | z_i, \mu_{1:K})}$$

- The numerator can be computed for any choice of the latent variables
- The problem is the denominator (the marginal probability of the observations)
 - This integral cannot easily be computed analytically
- We need some approximation..

Variational Inference

The main idea behind variational inference:

- Choose a family of distributions over the latent variables $z_{1:m}$ with its own set of variational parameters ν , i.e.

$$q(z_{1:m}|\nu)$$

- Then, we find the setting of the parameters that makes our approximation q closest to the posterior distribution.
 - This is where optimization algorithms come in.
- Then we can use q with the fitted parameters in place of the posterior.
 - E.g. to form predictions about future data, or to investigate the posterior distribution over the hidden variables, find modes, etc.

Variational Inference

- We want to minimize the KL divergence between our approximation $q(\mathbf{z}|\mathbf{x})$ and our posterior $p(\mathbf{z}|\mathbf{x})$

$$\text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}))$$

- But we can't actually minimize this quantity w.r.t q because $p(\mathbf{z}|\mathbf{x})$ is unknown

Evidence Lower Bound (ELBO)

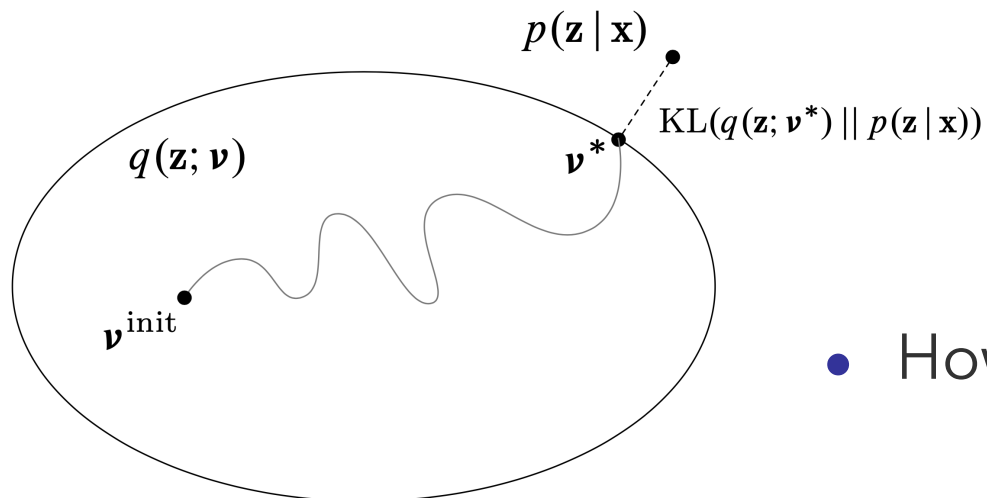
$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

- The ELBO is equal to the negative KL divergence up to a constant $\ell(\theta; \mathbf{x})$
- We maximize the ELBO over q to find an "optimal approximation" to $p(\mathbf{z}|\mathbf{x})$

Variational Inference

- Choose a family of distributions over the latent variables \mathbf{z} with its own set of variational parameters \mathbf{v} , i.e. $q(\mathbf{z}|\mathbf{x}, \mathbf{v})$
- We maximize the ELBO over q to find an “optimal approximation” to $p(\mathbf{z}|\mathbf{x})$

$$\begin{aligned} & \operatorname{argmax}_{\mathbf{v}} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} \right] \\ & = \operatorname{argmax}_{\mathbf{v}} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} [\log p(\mathbf{x}, \mathbf{z}|\theta)] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} [\log q(\mathbf{z}|\mathbf{x}, \mathbf{v})] \end{aligned}$$



- How do we choose the variational family $q(\mathbf{z}|\mathbf{x}, \mathbf{v})$?

Questions?