

DSC190: Machine Learning with Few Labels

Supervised Learning, Unsupervised Learning

Zhiting Hu

Lecture 2, September 28, 2021

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Logistics

- Office hours
 - Zhiting Hu: Tuesday 3-4pm, SDSC ~~247E~~249E
 - Meng Song: Wednesday 2:30-3:30pm, CSE 4109
- Project
- Presentation

Outline

- Supervised Learning
 - Maximum likelihood estimation (MLE)
 - Duality between MLE and Maximum Entropy Principle
- Unsupervised learning
 - Maximum likelihood estimation (MLE) with latent variables
 - EM algorithm for MLE

Supervised Learning

- Model to be learned $p_{\theta}(\mathbf{x})$
- Observe full data $\mathcal{D} = \{ \mathbf{x}^* \}$
 - i.i.d: independent, identically distributed
- Maximum Likelihood Estimation (MLE)
 - The most classical learning algorithm

$$\min_{\theta} - \mathbb{E}_{\mathbf{x}^* \sim \mathcal{D}} \left[\log p_{\theta}(\mathbf{x}^*) \right]$$

- MLE is closely connected to the Maximum Entropy (MaxEnt) principle

Exponential Family

- A distribution

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = h(\mathbf{x}) \exp\{\boldsymbol{\theta} \cdot T(\mathbf{x})\} / Z(\boldsymbol{\theta})$$

is an exponential family distribution

- $\boldsymbol{\theta} \in R^d$: natural (canonical) parameter
 - $T(\mathbf{x}) \in R^d$: sufficient statistics, features of data \mathbf{x}
 - $Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}, y} h(\mathbf{x}) \exp\{\boldsymbol{\theta} \cdot T(\mathbf{x})\}$: normalization factor
- Examples: Bernoulli, multinomial, Gaussian, Poisson, gamma,...

Example: Multivariate Gaussian Distribution

- For a continuous vector random variable $\mathbf{x} \in R^k$

$$p(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right\}$$

Moment parameter

$$= \frac{1}{(2\pi)^{k/2}} \exp\left\{-\frac{1}{2} \text{tr}(\Sigma^{-1} \mathbf{x} \mathbf{x}^T) + \mu^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu^T \Sigma^{-1} \mu - \log|\Sigma|\right\}$$

- Exponential family representation

$$\boldsymbol{\theta} = \left[\Sigma^{-1} \mu; -\frac{1}{2} \text{vec}(\Sigma^{-1}) \right] = [\boldsymbol{\theta}_1, \text{vec}(\boldsymbol{\theta}_2)], \quad \boldsymbol{\theta}_1 = \Sigma^{-1} \mu \text{ and } \boldsymbol{\theta}_2^- = -\frac{1}{2} \Sigma^{-1}$$

$$T(\mathbf{x}) = [\mathbf{x}; \text{vec}(\mathbf{x} \mathbf{x}^T)]$$

$$A(\boldsymbol{\theta}) = \frac{1}{2} \mu^T \Sigma^{-1} \mu + \log |\Sigma| = -\frac{1}{2} \text{tr}(\boldsymbol{\theta}_2 \boldsymbol{\theta}_1 \boldsymbol{\theta}_1^T) - \frac{1}{2} \log(-2\boldsymbol{\theta}_2)$$

$$h(\mathbf{x}) = (2\pi)^{-k/2}$$

Example: Multinomial Distribution

- For a binary vector random variable $\mathbf{x} \in \text{multi}(\mathbf{x}|\pi)$

$$\begin{aligned} p(\mathbf{x}|\pi) &= \pi_1^{x_1} \pi_2^{x_2} \cdots \pi_K^{x_K} = \exp\left\{\sum_k x_k \ln \pi_k\right\} \\ &= \exp\left\{\sum_{k=1}^{K-1} x_k \ln \pi_k + \left(1 - \sum_{k=1}^{K-1} x_k\right) \ln\left(1 - \sum_{k=1}^{K-1} \pi_k\right)\right\} \\ &= \exp\left\{\sum_{k=1}^{K-1} x_k \ln\left(\frac{\pi_k}{1 - \sum_{k=1}^{K-1} \pi_k}\right) + \ln\left(1 - \sum_{k=1}^{K-1} \pi_k\right)\right\} \end{aligned}$$

- Exponential family representation $\boldsymbol{\theta} = [\ln(\pi_k/\pi_K); 0]$

$$T(\mathbf{x}) = [\mathbf{x}]$$

$$A(\boldsymbol{\theta}) = -\ln\left(1 - \sum_{k=1}^{K-1} \pi_k\right) = \ln\left(\sum_{k=1}^K e^{\theta_k}\right)$$

$$h(\mathbf{x}) = 1$$

Maximum Likelihood for Exponential Family

$m(\mathbf{x})$: the number of times \mathbf{x} is observed in D

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) &= \sum_{\mathbf{x}} m(\mathbf{x}) \log p(\mathbf{x} | \boldsymbol{\theta}) \\ &= \sum_{\mathbf{x}} m(\mathbf{x}) \left(\sum_i \theta_i T_i(\mathbf{x}) - \log Z(\boldsymbol{\theta}) \right) \\ &= \sum_{\mathbf{x}} m(\mathbf{x}) \sum_i \theta_i T_i(\mathbf{x}) - N \log Z(\boldsymbol{\theta})\end{aligned}$$

- Take gradient and set to 0

$$\begin{aligned}\frac{\partial}{\partial \theta_i} \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) &= \sum_{\mathbf{x}} m(\mathbf{x}) T_i(\mathbf{x}) - N \frac{\partial}{\partial \theta_i} \log Z(\boldsymbol{\theta}) \\ &= \sum_{\mathbf{x}} m(\mathbf{x}) T_i(\mathbf{x}) - N \sum_{\mathbf{x}} p(\mathbf{x} | \boldsymbol{\theta}) T_i(\mathbf{x})\end{aligned}$$

At MLE, the expectations of the sufficient statistics under the model must match empirical feature average

$$\Rightarrow \boxed{\sum_{\mathbf{x}} p(\mathbf{x} | \boldsymbol{\theta}) T_i(\mathbf{x})} = \sum_{\mathbf{x}} \frac{m(\mathbf{x})}{N} T_i(\mathbf{x}) = \boxed{\sum_{\mathbf{x}} \tilde{p}(\mathbf{x} | \boldsymbol{\theta}) T_i(\mathbf{x})}$$

Maximum Entropy (MaxEnt)

- Given \mathcal{D} , to estimate $p(\mathbf{x})$
- We can approach the problem from an entirely different point of view. Begin with some fixed feature expectations:

$$\sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) = \sum_{\mathbf{x}} \frac{m(\mathbf{x})}{N} T_i(\mathbf{x}) := \alpha_i$$

- There may exist many distributions which satisfy them. Which one should we select?
 - MaxEnt principle: the most uncertain or flexible one, i.e., the one with maximum entropy
- This yields a new optimization problem:
 - This is a variational definition of a distribution!

$$\begin{aligned} \max_p \quad & H(p(\mathbf{x})) = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) \\ \text{s.t.} \quad & \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) = \alpha_i \\ & \sum_{\mathbf{x}} p(\mathbf{x}) = 1 \end{aligned}$$

Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(\mathbf{x})} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left(\sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left(\sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(\mathbf{x})} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left(\sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left(\sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

$$\frac{\partial L}{\partial p(\mathbf{x})} = 1 + \log p(\mathbf{x}) - \sum_i \theta_i T_i(\mathbf{x}) - \mu$$

$$p^*(\mathbf{x}) = e^{\mu-1} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\}$$

$$Z(\boldsymbol{\theta}) = e^{\mu-1} = \sum_{\mathbf{x}} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\} \quad \left(\text{since } \sum_{\mathbf{x}} p^*(\mathbf{x}) = 1 \right)$$

$$p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(\mathbf{x})} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left(\sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left(\sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

$$\frac{\partial L}{\partial p(\mathbf{x})} = 1 + \log p(\mathbf{x}) - \sum_i \theta_i T_i(\mathbf{x}) - \mu$$

$$p^*(\mathbf{x}) = e^{\mu-1} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\}$$

$$Z(\boldsymbol{\theta}) = e^{\mu-1} = \sum_{\mathbf{x}} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\} \quad \left(\text{since } \sum_{\mathbf{x}} p^*(\mathbf{x}) = 1 \right)$$

$$p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

- So feature constraints + MaxEnt \Rightarrow **exponential family**.
- Problem is strictly convex w.r.t. $p(\mathbf{x})$, so solution is unique.

Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(\mathbf{x})} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left(\sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left(\sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

$$p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

plug $p(\mathbf{x} | \boldsymbol{\theta})$ back into L , and since $\sum_{\mathbf{x}} \frac{m(\mathbf{x})}{N} T_i(\mathbf{x}) := \alpha_i$:

$$\max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \sum_{\mathbf{x}} m(\mathbf{x}) \sum_i \theta_i T_i(\mathbf{x}) - N \log Z(\boldsymbol{\theta})$$

- Recovers precisely the MLE problem of exponential family

• So feature constraints + MaxEnt \Rightarrow **exponential family**.

• Problem is strictly convex w.r.t. $p(\mathbf{x})$, so solution is unique.

Constraints from Data

- We have seen a case of **convex duality**:
 - In one case, we assume exponential family and show that Maximum Likelihood implies model expectations must match empirical expectations.
 - In the other case, we assume model expectations must match empirical feature counts and show that MaxEnt implies exponential family distribution.

A more general MaxEnt problem

$$\min_p \text{KL}(p(\mathbf{x}) \| h(\mathbf{x}))$$

$$\stackrel{\text{def}}{=} \sum_{\mathbf{x}} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{h(\mathbf{x})} = -\text{H}(p) - \sum_{\mathbf{x}} p(\mathbf{x}) \log h(\mathbf{x})$$

$$\text{s.t.} \quad \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) = \alpha_i$$

$$\sum_{\mathbf{x}} p(\mathbf{x}) = 1$$

$$\Rightarrow p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(\mathbf{x}) \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

Summary

- Maximum entropy is dual to maximum likelihood of exponential family distributions
- This provides an alternative view of the problem of fitting a model into data:
 - The data instances in the training set are treated as constraints, and the learning problem is treated as a constrained optimization problem.
 - We'll revisit this optimization-theoretic view of learning repeatedly in the future!

$$\begin{aligned} \max_p \quad & H(p(\mathbf{x})) = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) \\ \text{s.t.} \quad & \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) = \alpha_i \\ & \sum_{\mathbf{x}} p(\mathbf{x}) = 1 \end{aligned}$$

Unsupervised Learning

- Each data instance is partitioned into two parts:
 - observed variables \mathbf{x}
 - latent (unobserved) variables \mathbf{z}
- Want to learn a model $p_{\theta}(\mathbf{x}, \mathbf{z})$

Latent (unobserved) variables

- A variable can be unobserved (latent) because:
 - imaginary quantity: meant to provide some simplified and abstractive view of the data generation process
 - e.g., speech recognition models, mixture models, ...

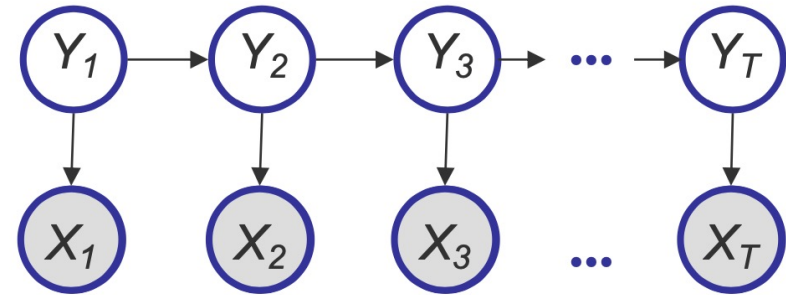
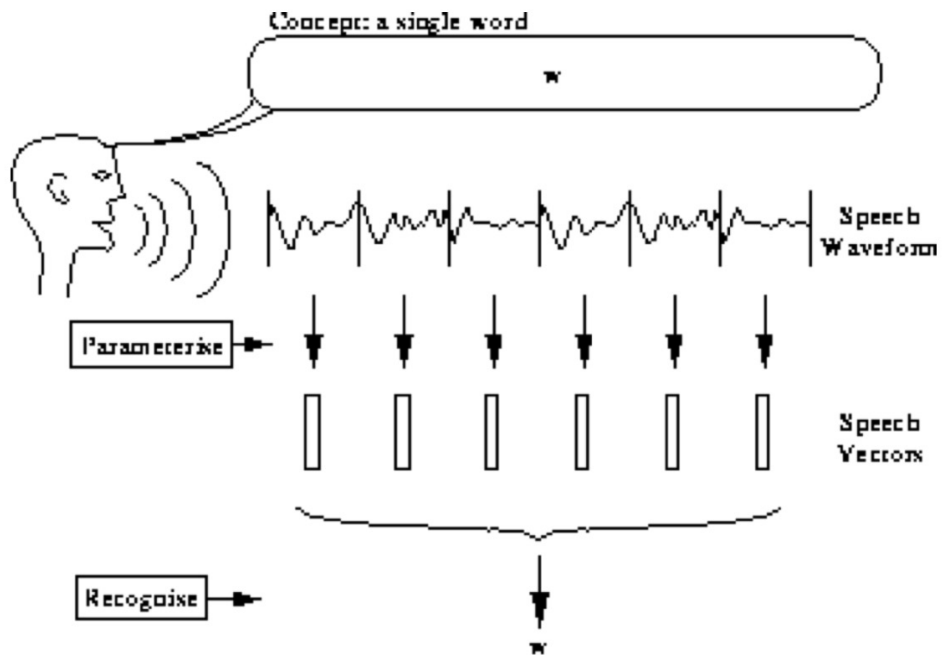
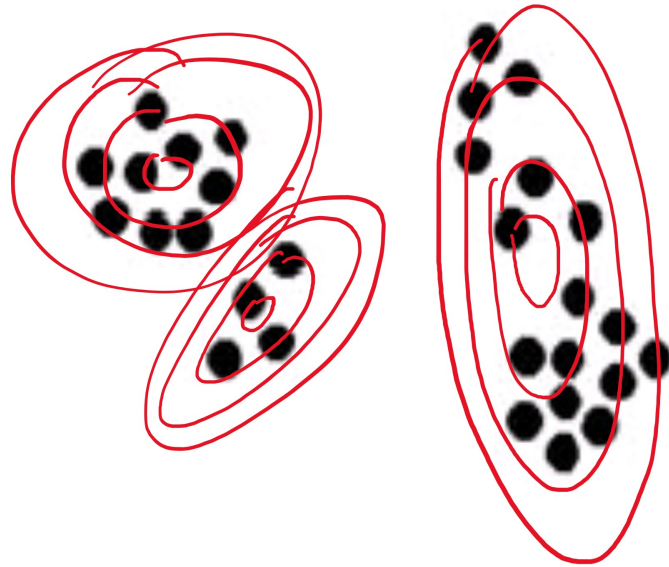


Fig. 1.2 Isolated Word Problem

Latent (unobserved) variables

- A variable can be unobserved (latent) because:
 - imaginary quantity: meant to provide some simplified and abstractive view of the data generation process
 - e.g., speech recognition models, mixture models, ...



Latent (unobserved) variables

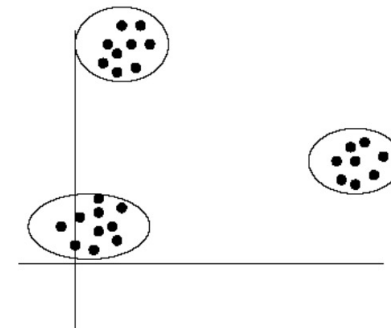
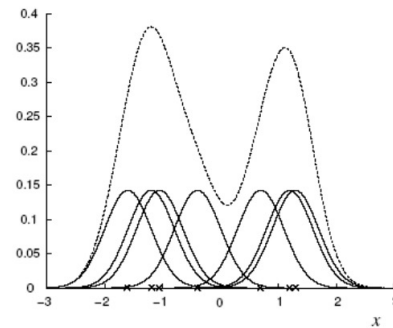
- A variable can be unobserved (latent) because:
 - imaginary quantity: meant to provide some simplified and abstractive view of the data generation process
 - e.g., speech recognition models, mixture models, ...
 - a real-world object (and/or phenomena), but difficult or impossible to measure
 - e.g., the temperature of a star, causes of a disease, evolutionary ancestors ...
 - a real-world object (and/or phenomena), but sometimes wasn't measured, because of faulty sensors, etc.
- Discrete latent variables can be used to partition/cluster data into sub-groups
- Continuous latent variables (factors) can be used for dimensionality reduction (e.g., factor analysis, etc.)

Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:

$$p(x_n | \mu, \Sigma) = \sum_k \pi_k N(x, | \mu_k, \Sigma_k)$$

mixture proportion mixture component



- This model can be used for unsupervised clustering.
 - This model (fit by AutoClass) has been used to discover new kinds of stars in astronomical data, etc.

Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:
 - Z is a latent class indicator vector:

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

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

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

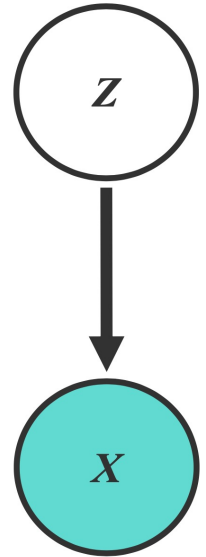
Parameters to be learned:

- The likelihood of a sample:

$$\begin{aligned} p(x_n | \mu, \Sigma) &= \sum_k p(z^k = \mathbf{1} | \pi) p(x_n | z^k = \mathbf{1}, \mu, \Sigma) \\ &= \sum_{z_n} \prod_k \left((\pi_k)^{z_n^k} N(x_n : \mu_k, \Sigma_k)^{z_n^k} \right) = \sum_k \pi_k N(x_n | \mu_k, \Sigma_k) \end{aligned}$$

mixture proportion

mixture component



Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components: $p(x_n | \mu, \Sigma) = \sum_k \pi_k N(x_n | \mu_k, \Sigma_k)$
- Recall MLE for completely observed data

- Data log-likelihood:
$$\ell(\theta; D) = \log \prod_n p(z_n, x_n) = \log \prod_n p(z_n | \pi) p(x_n | z_n, \mu, \sigma)$$

$$= \sum_n \log \prod_k \pi_k^{z_n^k} + \sum_n \log \prod_k N(x_n; \mu_k, \sigma)^{z_n^k}$$

$$= \sum_n \sum_k z_n^k \log \pi_k - \sum_n \sum_k z_n^k \frac{1}{2\sigma^2} (x_n - \mu_k)^2 + C$$

- MLE:

$$\hat{\pi}_{k,MLE} = \arg \max_{\pi} \ell(\theta; D),$$

$$\hat{\mu}_{k,MLE} = \arg \max_{\mu} \ell(\theta; D)$$

$$\hat{\sigma}_{k,MLE} = \arg \max_{\sigma} \ell(\theta; D)$$

$$\Rightarrow \hat{\mu}_{k,MLE} = \frac{\sum_n z_n^k x_n}{\sum_n z_n^k}$$

- What if we do not know z_n ?

Why is Learning Harder?

- **Complete log likelihood:** if both \mathbf{x} and \mathbf{z} can be observed, then

$$\ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z}|\theta) = \log p(\mathbf{z}|\theta_z) + \log p(\mathbf{x}|\mathbf{z}, \theta_x)$$

- Decomposes into a sum of factors, the parameter for each factor can be estimated separately
- But given that \mathbf{z} is not observed, $\ell_c(\theta; \mathbf{x}, \mathbf{z})$ is a random quantity, cannot be maximized directly
- **Incomplete (or marginal) log likelihood:** with \mathbf{z} unobserved, our objective becomes the log of a marginal probability:

$$\ell(\theta; \mathbf{x}) = \log p(\mathbf{x}|\theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\theta)$$

- All parameters become coupled together
- In other models when \mathbf{z} is complex (continuous) variables (as we'll see later), marginalization over \mathbf{z} is intractable.

Expectation Maximization (EM)

- For any distribution $q(\mathbf{z}|\mathbf{x})$, define **expected complete log likelihood**:

$$\mathbb{E}_q[\ell_c(\theta; \mathbf{x}, \mathbf{z})] = \sum_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$$

- A deterministic function of θ
- Inherit the factorizability of $\ell_c(\theta; \mathbf{x}, \mathbf{z})$
- Use this as the surrogate objective
- Does maximizing this surrogate yield a maximizer of the likelihood?

Expectation Maximization (EM)

- For any distribution $q(\mathbf{z}|\mathbf{x})$, define **expected complete log likelihood**:

$$\mathbb{E}_q[\ell_c(\theta; \mathbf{x}, \mathbf{z})] = \sum_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$$

- Jensen's inequality

$$\ell(\theta; \mathbf{x}) = \log p(\mathbf{x} | \theta)$$

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

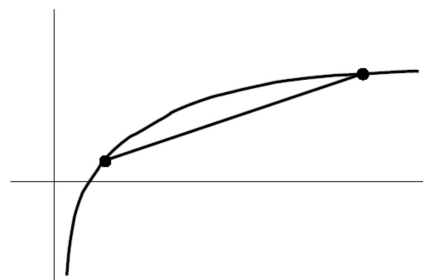
$$= \log \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z} | \mathbf{x})}$$

$$\geq \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z} | \mathbf{x})}$$

Evidence Lower Bound (ELBO)

$$= \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log p(\mathbf{x}, \mathbf{z} | \theta) - \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log q(\mathbf{z} | \mathbf{x})$$

$$= \mathbb{E}_q[\ell_c(\theta; \mathbf{x}, \mathbf{z})] + H(q)$$



Expectation Maximization (EM)

- For any distribution $q(\mathbf{z}|\mathbf{x})$, define **expected complete log likelihood**:

$$\mathbb{E}_q[\ell_c(\theta; \mathbf{x}, \mathbf{z})] = \sum_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$$

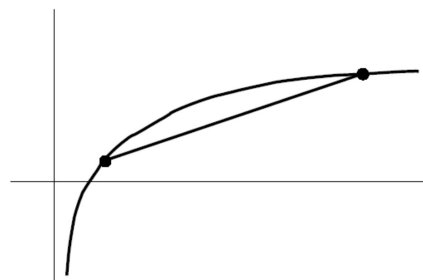
- Jensen's inequality

$$\ell(\theta; \mathbf{x}) = \log p(\mathbf{x} | \theta)$$

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

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

$$\geq \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z} | \mathbf{x})}$$



- Indeed we have

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

Lower Bound and Free Energy

- For fixed data \mathbf{x} , define a functional called the (variational) free energy:

$$F(q, \theta) = -\mathbb{E}_q[\ell_c(\theta; \mathbf{x}, \mathbf{z})] - H(q) \geq \ell(\theta; \mathbf{x})$$

- The EM algorithm is coordinate-descent on F
 - At each step t :

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

- M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta)$

E-step: minimization of $F(q, \theta)$ w.r.t q

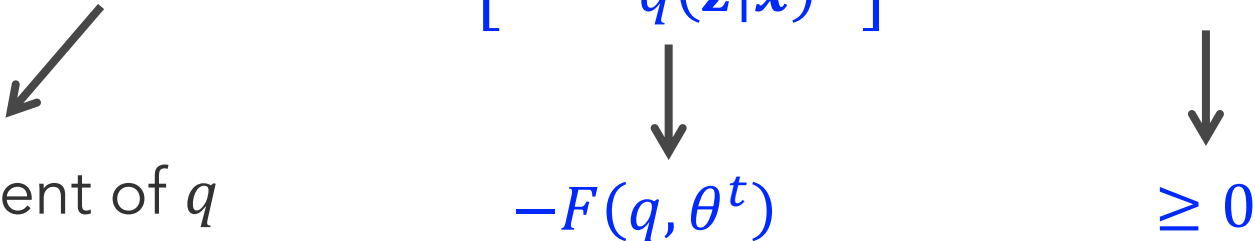
- Claim:

$$q^{t+1} = \operatorname{argmin}_q F(q, \theta^t) = p(\mathbf{z}|\mathbf{x}, \theta^t)$$

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

- Proof (easy): recall

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



Independent of q $-F(q, \theta^t)$ ≥ 0

- $F(q, \theta^t)$ is minimized when $\text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta^t)) = 0$, which is achieved only when $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x}, \theta^t)$

M-step: minimization of $F(q, \theta)$ w.r.t θ

- Note that the free energy breaks into two terms:

$$F(q, \theta) = -\mathbb{E}_q[\ell_c(\theta; \mathbf{x}, \mathbf{z})] - H(q) \geq \ell(\theta; \mathbf{x})$$

- The first term is the expected complete log likelihood and the second term, which does not depend on q , is the entropy.
- Thus, in the M-step, maximizing with respect to θ for fixed q we only need to consider the first term:

$$\theta^{t+1} = \operatorname{argmax}_{\theta} \mathbb{E}_q[\ell_c(\theta; \mathbf{x}, \mathbf{z})] = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q^{t+1}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$$

- Under optimal q^{t+1} , this is equivalent to solving a standard MLE of fully observed model $p(\mathbf{x}, \mathbf{z}|\theta)$, with \mathbf{z} replaced by its expectation w.r.t $p(\mathbf{z}|\mathbf{x}, \theta^t)$

Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:
 - Z is a latent class indicator vector:

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

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

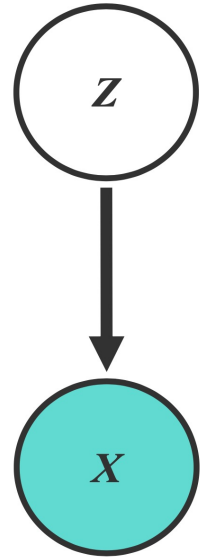
$$p(x_n | z_n^k = \mathbf{1}, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

- The likelihood of a sample:

$$\begin{aligned} p(x_n | \mu, \Sigma) &= \sum_k p(z^k = \mathbf{1} | \pi) p(x_n | z^k = \mathbf{1}, \mu, \Sigma) \\ &= \sum_{z_n} \prod_k \left((\pi_k)^{z_n^k} N(x_n : \mu_k, \Sigma_k)^{z_n^k} \right) = \sum_k \pi_k N(x_n | \mu_k, \Sigma_k) \end{aligned}$$

mixture component

mixture proportion



Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components
- The expected complete log likelihood

$$\begin{aligned}\mathbb{E}_q [\ell_c(\boldsymbol{\theta}; \mathbf{x}, \mathbf{z})] &= \sum_n \mathbb{E}_q [\log p(z_n | \boldsymbol{\pi})] + \sum_n \mathbb{E}_q [\log p(x_n | z_n, \boldsymbol{\mu}, \boldsymbol{\Sigma})] \\ &= \sum_n \sum_k \mathbb{E}_q [z_n^k] \log \pi_k - \frac{1}{2} \sum_n \sum_k \mathbb{E}_q [z_n^k] \left((x_n - \mu_k)^T \boldsymbol{\Sigma}_k^{-1} (x_n - \mu_k) + \log |\boldsymbol{\Sigma}_k| + C \right)\end{aligned}$$

- E-step: computing the posterior of z_n given the current estimate of the parameters (i.e., $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$)

$$p(z_n^k = 1 | \mathbf{x}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}) = \frac{\pi_k^{(t)} N(x_n, | \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)})}{\sum_i \pi_i^{(t)} N(x_n, | \boldsymbol{\mu}_i^{(t)}, \boldsymbol{\Sigma}_i^{(t)})}$$

$\nearrow p(z_n^k = 1, \mathbf{x}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)})$
 $\searrow p(\mathbf{x}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)})$

Example: Gaussian Mixture Models (GMMs)

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

$$\pi_k^* = \arg \max \langle l_c(\boldsymbol{\theta}) \rangle, \quad \Rightarrow \quad \frac{\partial}{\partial \pi_k} \langle l_c(\boldsymbol{\theta}) \rangle = 0, \forall k, \quad \text{s.t.} \quad \sum_k \pi_k = 1$$
$$\Rightarrow \quad \pi_k^* = \frac{\sum_n \langle z_n^k \rangle_{q^{(t)}}}{N} = \frac{\sum_n \tau_n^{k(t)}}{N} = \frac{\langle n_k \rangle}{N}$$

$$\mu_k^* = \arg \max \langle l(\boldsymbol{\theta}) \rangle, \quad \Rightarrow \quad \mu_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} x_n}{\sum_n \tau_n^{k(t)}}$$

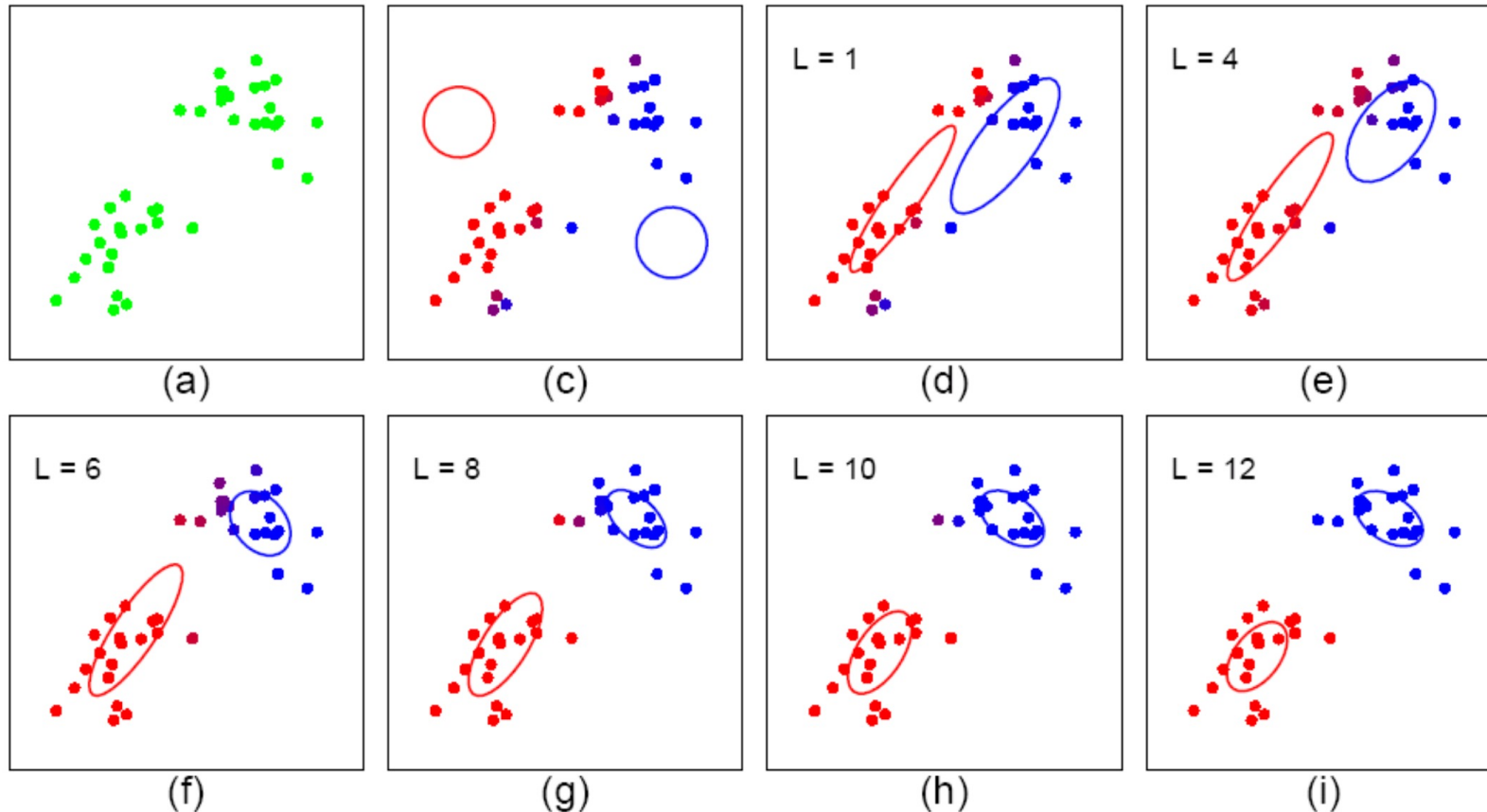
$$\Sigma_k^* = \arg \max \langle l(\boldsymbol{\theta}) \rangle, \quad \Rightarrow \quad \Sigma_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} (x_n - \mu_k^{(t+1)})(x_n - \mu_k^{(t+1)})^T}{\sum_n \tau_n^{k(t)}}$$

Fact:

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

Example: Gaussian Mixture Models (GMMs)

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

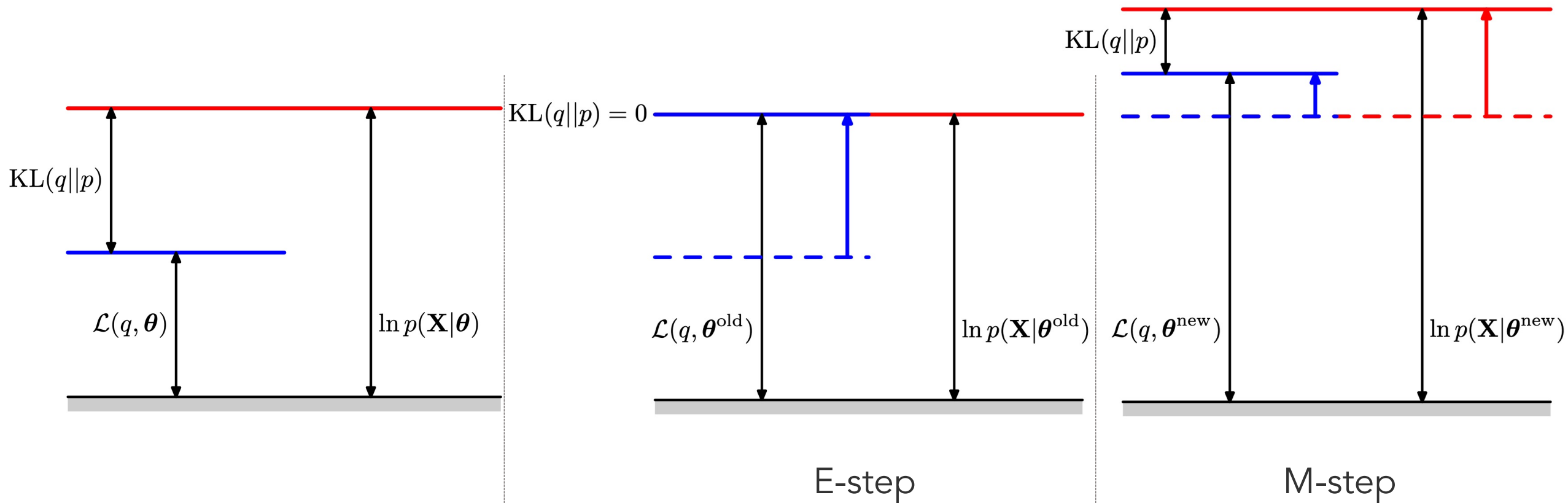


Summary: EM Algorithm

- 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.
- Alternate between filling in the latent variables using the best guess (posterior) and updating the parameters based on this guess:
 - E-step: $q^{t+1} = \arg \min_q F(q, \theta^t)$
 - M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta)$

Each EM iteration guarantees to improve the likelihood

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



EM Variants

- Sparse EM
 - Do not re-compute exactly the posterior probability on each data point under all models, because it is almost zero.
 - Instead keep an “active list” which you update every once in a while.
- Generalized (Incomplete) EM:
 - It might be hard to find the ML parameters in the M-step, even given the completed data. We can still make progress by doing an M-step that improves the likelihood a bit (e.g. gradient step).

Summary

- Supervised Learning
 - Maximum likelihood estimation (MLE)
 - Duality between MLE and Maximum Entropy Principle
- Unsupervised learning
 - Maximum likelihood estimation (MLE) with latent variables
 - EM algorithm for MLE

Questions?