

Module 12

Latent Variabl Model

Kentaro Nakamura

GOV 2003

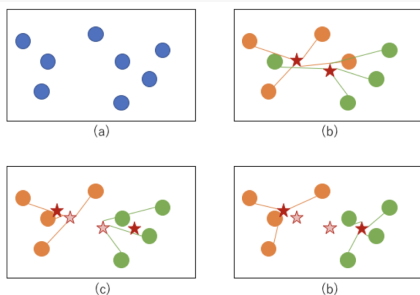
May 4th, 2026

Agenda

- k-means algorithm
- Gaussian Mixture Models (GMM)
- EM Algorithm in general
 - How to derive E / M-step
 - Why EM algorithm works
 - Uncertainty quantification under EM
- Useful resource: Bishop (2006) Chapter 9
- For future reference (not for final exam, but in appendix)
 - Variational Inference
 - MCMC Methods

k-means algorithm (Overview)

- **Goal:** discover the latent cluster that separates the data well
 - Input: \mathbf{X}_i (observed data), K (number of clusters)
- **Assumption:** Each observation belongs to one cluster (**hard clustering**)
- **Algorithm (Lloyd's algorithm)**
 - Randomly assign each observation to one of K clusters.
 - Compute the centroid of the observations assigned to each cluster.
 - Reassign each observation to the cluster with the nearest centroid.
 - Repeat steps 2 and 3 until the cluster assignments no longer change.



k-means algorithm: optimization (1)

- Formally, the objective function is written as

$$\{\mathbf{Z}, \boldsymbol{\mu}\} := \arg \min_{\mathbf{Z}, \boldsymbol{\mu}} \sum_{i=1}^n \sum_{k=1}^K \mathbf{1}\{Z_i = k\} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2$$

- Let's see why Lloyd's algorithm solves this objective function.
- Now, let's fix the assignment \mathbf{Z} . Then, the derivative w.r.t. μ_k is

$$\frac{\partial}{\partial \mu_k} \sum_{i:Z_i=k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2 = -2 \cdot \sum_{i:Z_i=k} (\mathbf{x}_i - \boldsymbol{\mu}_k)$$

which means that setting $\boldsymbol{\mu}_k = \frac{1}{n_k} \sum_{i:Z_i=k} \mathbf{x}_i$ (equivalent to STEP 2 in Lloyd's algorithm) achieves the first order condition

k-means algorithm: optimization (2)

- Now, think about optimizing cluster assignment given the centroid

$$\{\mathbf{Z}^{(t+1)}\} := \arg \min_{\mathbf{Z}} \sum_{i=1}^n \sum_{k=1}^K \mathbf{1}\{Z_i = k\} \|\mathbf{X}_i - \boldsymbol{\mu}_k^{(t)}\|_2^2$$

- Note that the objective function is separable across i :

$$= \sum_{i=1}^n \sum_{k=1}^K \mathbf{1}\{Z_i = k\} \|\mathbf{X}_i - \boldsymbol{\mu}_k^{(t)}\|_2^2 = \sum_{i=1}^n \min_{k \in \{1, \dots, K\}} \|\mathbf{X}_i - \boldsymbol{\mu}_k^{(t)}\|_2^2$$

- Therefore, for each observation i , we can optimize independently:

$$Z_i^{(t+1)} = \arg \min_{k \in \{1, \dots, K\}} \|\mathbf{X}_i - \boldsymbol{\mu}_k^{(t)}\|_2^2$$

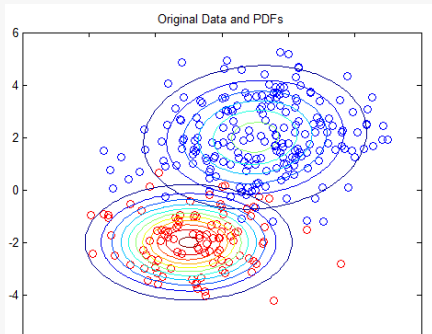
- This is exactly STEP 3 in Lloyd's algorithm (assign to nearest centroid)
 - Note that this might yield local optimal depending on initialization, so typically repeat the procedure with multiple initialization and pick up the one that minimizes the objective function the most.

Gaussian (Finite) Mixture Models

- We want to extend k-means by allowing each observations to belong to multiple cluster
- **Gaussian Mixture Models** (GMM): The entire data generating process is

$$\mathbf{X}_i \mid Z_i, \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K \sim \mathcal{N}(\boldsymbol{\mu}_{Z_i}, \boldsymbol{\Sigma}_{Z_i})$$
$$Z_i \sim \text{Categorical}(\boldsymbol{\pi})$$

where $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]^\top$, $\sum_{k=1}^K \pi_k = 1$ and $0 \leq \pi_k \leq 1$.



Likelihood of Latent Variable Model (1)

- Let's start with the likelihood of GMM as if we observe latent variable.

$$\mathcal{L}^{(\text{full})}(\mathbf{X}, \mathbf{Z}; \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{i=1}^n \prod_{k=1}^K (\pi_k \mathcal{N}(X_i; \mu_k, \Sigma_k))^{\mathbf{1}\{Z_i=k\}}$$

- Thus, the log-likelihood is given by

$$\ell^{(\text{full})}(\mathbf{X}, \mathbf{Z}; \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^n \sum_{k=1}^K \mathbf{1}\{Z_i = k\} (\log \pi_k + \log \mathcal{N}(X_i; \mu_k, \Sigma_k))$$

- Thus, if we observe Z_i , then we get MLE

$$\hat{\pi}_k = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{Z_i = k\}, \quad \hat{\mu}_k = \frac{\sum_{i=1}^n \mathbf{1}\{Z_i = k\} X_i}{\sum_{i=1}^n \mathbf{1}\{Z_i = k\}}$$

$$\hat{\Sigma}_k = \frac{1}{\sum_{i=1}^n \mathbf{1}\{Z_i = k\}} \sum_{i=1}^n \mathbf{1}\{Z_i = k\} (X_i - \hat{\mu}_k)(X_i - \hat{\mu}_k)^\top$$

- So, if we observe Z_i , then the problem is easy

Likelihood of Latent Variable Model (2)

- However, Z_i is latent. We cannot observe. What we observe is the observed likelihood:

$$\begin{aligned}\mathcal{L}^{(\text{obs})}(\mathbf{X}; \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) &= \sum_{\mathbf{z}} \mathcal{L}^{(\text{full})}(\mathbf{X}, \mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) \\ &= \sum_{z_1, \dots, z_n} \prod_{i=1}^n \prod_{k=1}^K (\pi_k \mathcal{N}(X_i; \mu_k, \Sigma_k)) \mathbf{1}_{\{Z_i=k\}} \\ &= \sum_{z_1, \dots, z_n} \prod_{i=1}^n (\pi_{Z_i} \mathcal{N}(X_i; \mu_{Z_i}, \Sigma_{Z_i})) \\ &= \prod_{i=1}^n \sum_{z_i} (\pi_{z_i} \mathcal{N}(X_i; \mu_{z_i}, \Sigma_{z_i})) \\ &= \prod_{i=1}^n \sum_k \pi_k \mathcal{N}(X_i; \mu_k, \Sigma_k)\end{aligned}$$

where the last equality is by re-indexing.

- Recall that $\Pr(X = x) = \sum_y \Pr(X_i = x, Y_i = y)$ (this is why you need to integrate out the latent)

Problem in Latent Variable Model

- Therefore, the observed data log likelihood is given by

$$\log p(\mathbf{X} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{i=1}^n \log \left\{ \sum_k \pi_k \mathcal{N}(X_i; \mu_k, \Sigma_k) \right\}$$

- This is hard to maximize because of the structure $\log(\text{sum})$
 - The log of sum over convex function is generally not convex
 - To see why, consider the example

$$f(x) = \log\{\exp(-[x - 2]^2) + \exp(-([x + 2]^2))\}$$

where you see each term is convex

- When $x = 2 \rightarrow$ first term dominates (peak)
 - When $x = -2 \rightarrow$ second term dominates (peak)
- Similarly, GMM suffers from non-convexity problem
 - There are K latent cluster, but each cluster label does not have meaning
 - **Label switching**: There are $K!$ equivalent solutions

EM Algorithm for GMM (1)

- **Goal:** Maximize the observed data likelihood $\ln p(\mathbf{X} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})$
- Now, recall that we have

$$\ln p(\mathbf{X} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{i=1}^n \log \left\{ \sum_k \pi_k \mathcal{N}(\mathbf{X}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

- Now, check the first order condition for each parameter. The derivative w.r.t. $\boldsymbol{\mu}_k$ is given by

$$0 = - \sum_{i=1}^n \frac{\pi_k \mathcal{N}(\mathbf{X}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\underbrace{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{X}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}_{:=\gamma_{i,k}(=\Pr(Z_i=k|\mathbf{X}_i))}} \boldsymbol{\Sigma}_k^{-1} (\mathbf{X}_i - \boldsymbol{\mu}_k)$$

where $\gamma_{i,k}$ is called responsibility parameter.

- Similarly, the derivative with respect to $\boldsymbol{\mu}_k$ yields

$$\boldsymbol{\mu}_k = \frac{\sum_{i=1}^n \gamma_{i,k} \mathbf{X}_i}{\sum_{i=1}^n \gamma_{i,k}}$$

EM Algorithm for GMM (2)

- Next, let's derive the MLE for Σ_k .
- The first order condition with respect to Σ_k^{-1} is

$$\begin{aligned} 0 &= \sum_{i=1}^n \gamma_{i,k} \frac{\partial}{\partial \Sigma_k^{-1}} \log \mathcal{N}(\mathbf{X}_i; \boldsymbol{\mu}_k, \Sigma_k) \\ &= \frac{1}{2} \sum_{i=1}^n \gamma_{i,k} \left[\Sigma_k - (\mathbf{X}_i - \boldsymbol{\mu}_k)(\mathbf{X}_i - \boldsymbol{\mu}_k)^\top \right]. \end{aligned}$$

- Therefore,

$$\sum_{i=1}^n \gamma_{i,k} \Sigma_k = \sum_{i=1}^n \gamma_{i,k} (\mathbf{X}_i - \boldsymbol{\mu}_k)(\mathbf{X}_i - \boldsymbol{\mu}_k)^\top.$$

- This gives

$$\Sigma_k = \frac{\sum_{i=1}^n \gamma_{i,k} (\mathbf{X}_i - \boldsymbol{\mu}_k)(\mathbf{X}_i - \boldsymbol{\mu}_k)^\top}{\sum_{i=1}^n \gamma_{i,k}}.$$

EM Algorithm for GMM (3)

- For the membership probability π_k , we need to incorporate constraints.
 - Using Lagrange multiplier and maximizing

$$\log p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) + \lambda \left(\sum_{k=1}^K -1 \right)$$

which yields

$$0 = \sum_{i=1}^n \frac{\mathcal{N}(\mathbf{X}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{X}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} + \lambda$$

- Now, multiplying both sides by π_k , we get

$$0 = \sum_{i=1}^n \gamma_{i,k} + \lambda \pi_k$$

- Since $\sum \pi_k = 1$, $\sum_k \{ \sum_{i=1}^n \gamma_{i,k} + \lambda \pi_k \} = \sum_k \sum_{i=1}^n \gamma_{i,k} + \lambda = 0$.
 - I.e., $\lambda = -N$
 - This gives us $\pi_k = \sum_{i=1}^n \gamma_{i,k} / N$.

EM Algorithm for GMM (4)

- However, we do not observe $\gamma_{i,k}$
 - If we observe $\gamma_{i,k} = \Pr(Z_i = k \mid \mathbf{X}_i)$, then we can maximize the observed data log-likelihood
 - This suggests the iterative procedure
- **E-Step:** Calculate the responsibility parameter given μ_k , Σ_k , and π :

$$\gamma_{i,k} = \frac{\pi_k \mathcal{N}(\mathbf{X}_i; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{X}_i; \mu_j, \Sigma_j)}$$

- **M-Step:** Maximize the log-likelihood. I.e., obtain the MLE estimator given $\gamma_{i,k}$ by

$$\mu_k = \frac{\sum_{i=1}^n \gamma_{i,k} \mathbf{X}_i}{\sum_{i=1}^n \gamma_{i,k}}, \quad \Sigma_k = \frac{\sum_{i=1}^n \gamma_{i,k} (\mathbf{X}_i - \mu_k)(\mathbf{X}_i - \mu_k)^\top}{\sum_{i=1}^n \gamma_{i,k}}$$
$$\pi_k = \frac{\sum_{i=1}^n \gamma_{i,k}}{N}$$

EM Algorithm in general (1)

- Let's understand the general logic and formulation for EM algorithm
- **Setup:** Let \mathbf{x} be the observed data, \mathbf{z} be latent, and θ be the parameter of the model
 - Then, the full-data log likelihood is written as $\log p(\mathbf{x}, \mathbf{z}, \theta)$ and observed data log likelihood (known as evidence) is $\log p(\mathbf{x}; \theta)$
 - Our goal is to maximize the observed log-likelihood $\log p(\mathbf{x}; \theta)$, but suppose that its optimization is hard while the optimization of full data likelihood is significantly easier
- Understand the difference between latent variable and parameter
 - Latent variable: random, varied by each unit
 - Parameter: fixed, shared across units

EM Algorithm in general (2)

- Now, notice that

$$\underbrace{\log p(\mathbf{x}; \theta)}_{\text{observed}} = \log \left(\sum_{\mathbf{z}} \underbrace{p(\mathbf{x}, \mathbf{z}, \theta)}_{\text{complete}} \right)$$

- As with the case of GMM, the observed log-likelihood is log of sum over complete likelihood, which is hard to optimize
- **Algorithm:**
 - **E-step** (Expectation step):

$$Q(\theta | \theta^{(t)}) = \mathbb{E}_{\mathbf{z} | \mathbf{x}, \theta^{(t)}} [\log p(\mathbf{x}, \mathbf{z}; \theta)]$$

- Take expectation of the **complete-data log-likelihood** with respect to the posterior of latent variables under current parameter $\theta^{(t)}$
 - **M-step** (Maximization step):

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta | \theta^{(t)})$$

- **Intuition:** The difficulty comes from $\log p(\mathbf{x}; \theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}; \theta)$
 - EM avoids this by:
 - **E-step:** “fill in” missing data using current model (soft assignment)
 - **M-step:** maximize as if data were complete

Theory of EM Algorithm (1)

- Let's understand why EM algorithm works
- Now, notice that

$$\begin{aligned}\log p(x; \theta) &= \log \left(\sum_z p(x, z; \theta) \right) \\ &= \log \left(\sum_z \frac{p(x, z; \theta)}{q(z)} q(z) \right) \\ &\geq \sum_z \log \frac{p(x, z; \theta)}{q(z)} q(z) \quad (\because \text{Jensen's inequality}) \\ &= \mathbb{E}_z[\log p(x, z; \theta)] - \mathbb{E}_z[\log q(z)] := \mathcal{L}(q, \theta)\end{aligned}$$

where the last term $\mathcal{L}(q, \theta)$ is called **Evidence Lower Bound (ELBO)**

Theory of EM Algorithm (2)

- Notice that ELBO is written as

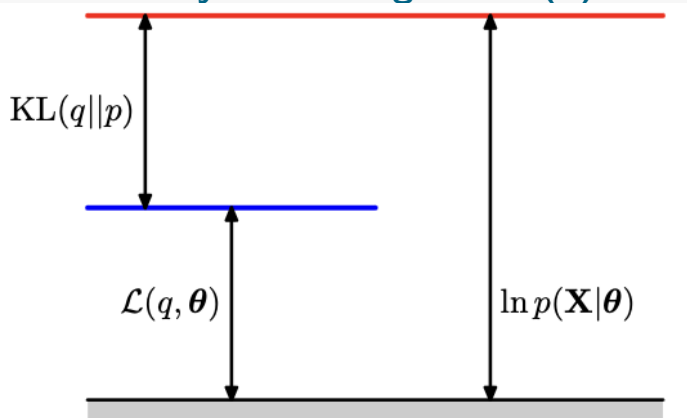
$$\begin{aligned}\mathcal{L}(q, \theta) &= \mathbb{E}_z[\log p(z | x; \theta)p(x; \theta)] - \mathbb{E}_z[\log q(z)] \\ &= \mathbb{E}_z[\log p(z | x; \theta)] - \mathbb{E}_z[\log q(z)] + \log p(x; \theta) \\ &= \underbrace{\log p(x; \theta)}_{\text{Evidence}} - \text{KL}(q(z) || p(z | x; \theta))\end{aligned}$$

- So, ELBO becomes lower bound when KL divergence is zero
- In other words,

$$\log p(x; \theta) = \text{KL}(q(z) || p(z | x; \theta)) + \mathcal{L}(q, \theta)$$

- Recall that KL-divergence is non-negative and takes 0 when $q(z) = p(z | x; \theta)$

Theory of EM Algorithm (3)

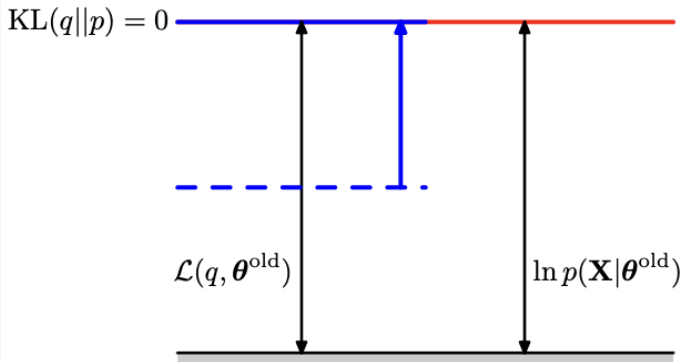


- Recall that our goal is to maximize the observed log-likelihood $\log p(x; \theta)$
 - This corresponds to maximize the lower bound $\mathcal{L}(q, \theta)$
- Notice that we have two unknowns: q and θ
 - We iteratively update them

Theory of EM Algorithm (4)

- **E-step:** Maximize ELBO $\mathcal{L}(q, \theta)$ with respect to q while fixing $\theta = \theta^{\text{old}}$.
 - Notice that once θ is fixed, we can maximize ELBO by choosing q makes KL divergence to be 0 (recall the previous derivation)
 - Such q is given by

$$q(z) = p(z | x; \theta)$$



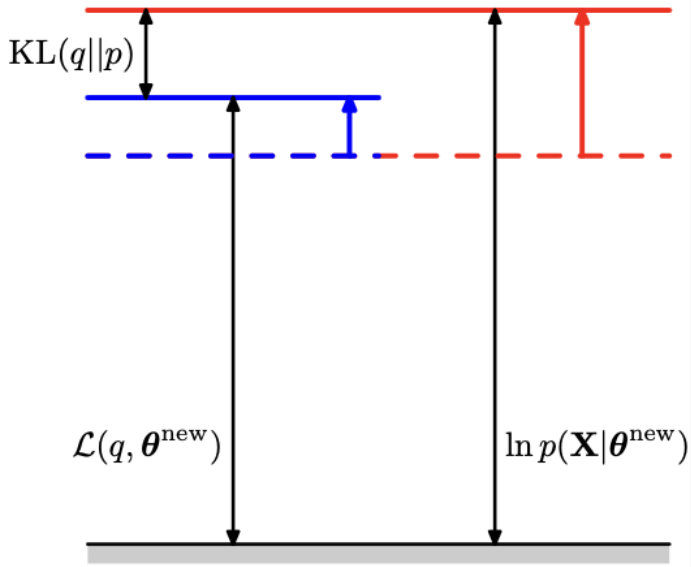
Theory of EM Algorithm (5)

- **M-step:** Maximize ELBO $\mathcal{L}(q, \theta)$ with respect to θ while fixing q
- Notice that after E-step (setting $q(z) = p(z | x; \theta^{\text{old}})$), ELBO is

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_z \log p(x, z; \theta) p(z | x; \theta^{\text{old}}) \\ &\quad - \sum_z p(z | x; \theta^{\text{old}}) \log p(z | x; \theta^{\text{old}}) \\ &= \underbrace{\sum_z \log p(x, z; \theta) p(z | x; \theta^{\text{old}})}_{\text{Q-function}} + \text{constant}\end{aligned}$$

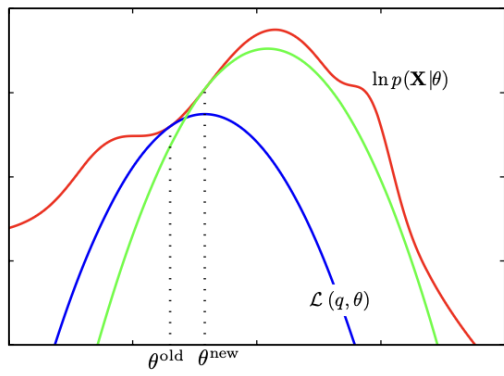
- This is the justification of maximizing Q-function

Theory of EM Algorithm (6)



Theory of EM Algorithm (7)

- EM algorithm has monotone convergence property



- E-step: Make ELBO $\mathcal{L}(q, \theta)$ tangential to log-likelihood (maximization of ELBO w.r.t. q under fixed θ)
- M-step: Make ELBO $\mathcal{L}(q, \theta)$ larger w.r.t. θ

Theory of EM Algorithm (8): Monotone convergence

- Let $q^{(t)}(\mathbf{z}) = p(\mathbf{z} | \mathbf{x}; \theta^{(t)})$

- From the ELBO decomposition,

$$\log p(\mathbf{x}; \theta) = \mathcal{L}(q, \theta) + \text{KL}\{q(\mathbf{z}) || p(\mathbf{z} | \mathbf{x}; \theta)\}.$$

- In the E-step, we set $q^{(t)}(\mathbf{z}) = p(\mathbf{z} | \mathbf{x}; \theta^{(t)})$, so

$$\text{KL}\{q^{(t)}(\mathbf{z}) || p(\mathbf{z} | \mathbf{x}; \theta^{(t)})\} = 0.$$

- Therefore,

$$\log p(\mathbf{x}; \theta^{(t)}) = \mathcal{L}(q^{(t)}, \theta^{(t)}).$$

Theory of EM Algorithm (9): Monotone convergence

- In the M-step, we choose

$$\theta^{(t+1)} = \arg \max_{\theta} \mathcal{L}(q^{(t)}, \theta).$$

- Therefore,

$$\mathcal{L}(q^{(t)}, \theta^{(t+1)}) \geq \mathcal{L}(q^{(t)}, \theta^{(t)}).$$

- Also, because KL divergence is non-negative,

$$\log p(\mathbf{x}; \theta^{(t+1)}) \geq \mathcal{L}(q^{(t)}, \theta^{(t+1)}).$$

- Combining the inequalities,

$$\log p(\mathbf{x}; \theta^{(t+1)}) \geq \mathcal{L}(q^{(t)}, \theta^{(t+1)}) \geq \mathcal{L}(q^{(t)}, \theta^{(t)}) = \log p(\mathbf{x}; \theta^{(t)}).$$

- Thus, the observed-data log-likelihood never decreases.

Uncertainty quantification of EM

- EM algorithm gives you the point estimates $\hat{\theta}$, but you cannot directly obtain the uncertainty estimates
 - Bootstrap is the easiest, but it is computationally too expensive
 - If you can work with the observed log-likelihood, then the asymptotic theory of likelihood tells us that

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, I^{-1}), \quad I(\theta_0) = \mathbb{E}[s_i(\theta_0)s_i(\theta_0)^\top]$$

where $s_i(\theta_0) = \nabla \log p(x_i; \theta)$ is the score.

- However, we often use EM because the observed log-likelihood is intractable and it is hard to compute Hessian
- Key tools: **Missing Information Principle**
 - Formally,

$$\underbrace{-\ell''(\theta | \mathbf{X})}_{\text{Observed}} = \mathbb{E}[\underbrace{-\ell''(\theta | \mathbf{X}, \mathbf{Z})}_{\text{Complete}} | \mathbf{X}] - \mathbb{E}[\underbrace{-\nabla^2 \log f(\mathbf{Z} | \mathbf{X}, \theta)}_{\text{Missing}} | \mathbf{X}]$$

- Implication: We can easily obtain the complete data information (similar to E-step), so we only need to obtain missing data part
 - Importantly, this avoids differentiating the observed log-likelihood directly (which involves log-sum structure)

Derivation of Missing-Data Principle

- Now, by factorization,

$$f(\mathbf{X}, \mathbf{Z} | \theta) = f(\mathbf{Z} | \mathbf{X}, \theta) f(\mathbf{X} | \theta).$$

- So taking logs gives

$$\ell(\theta | \mathbf{X}, \mathbf{Z}) = \log f(\mathbf{X}, \mathbf{Z} | \theta) = \log f(\mathbf{Z} | \mathbf{X}, \theta) + \log f(\mathbf{X} | \theta).$$

- Now take the Hessian with respect to θ , we get

$$\ell''(\theta | \mathbf{X}) = \ell''(\theta | \mathbf{X}, \mathbf{Z}) - \nabla^2 \log f(\mathbf{Z} | \mathbf{X}, \theta).$$

- Finally, take conditional expectation given \mathbf{X} (notice that the first term only depends on \mathbf{X})

$$-\ell''(\theta | \mathbf{X}) = \mathbb{E}[-\ell''(\theta | \mathbf{X}, \mathbf{Z}) | \mathbf{X}] - \mathbb{E}\left[-\nabla^2 \log f(\mathbf{Z} | \mathbf{X}, \theta) | \mathbf{X}\right].$$

Extra: Approximation Methods

- **Problem of EM algorithm:** We need to evaluate the expectation of the complete-data log-likelihood with respect to the posterior distribution of latent variable
 - In many applications, however, this is not possible.
- In such situation, we need to use approximation methods
 - Stochastic: MCMC (Metropolis Hasting, Gibbs Sampling, Hamiltonian Monte Carlo)
 - Gives accurate estimates with uncertainty quantification, but not scalable for large and complex data
 - Deterministic: Variational Inference
 - Gives approximate estimates, but scalable and fast
- This is not required for final exam, but it is the main technique you need to use in future
 - Variational Inference: Ch.10 of Bishop (2006), Blei et al. (2017 JASA)
 - MCMC: Ch.11 of Bishop (2006)

Extra: Variational Inference (1)

- Now, we are in the bayesian framework. So, the parameter θ is a random variable
 - Thus, we denote the set of all parameters and latent variable as \mathbf{Z}
 - Our goal is to find the distribution that is the closest to the exact posterior $p(\mathbf{Z} | \mathbf{X})$, which is minimizing the KL-divergence $\text{KL}(q(z) | p(z | x))$ over some candidate set of distributions \mathcal{Q} .
- Importantly, the minimization of KL divergence is equivalent to the maximization of ELBO:

$$\begin{aligned}q^*(z) &= \arg \min_{q(z) \in \mathcal{Q}} \text{KL}(q(z) || \underbrace{p(z | x)}_{\text{Posterior}}) \\ &= \arg \min_{q(z) \in \mathcal{Q}} \left(\underbrace{\log p(x)}_{\text{constant}} - \text{ELBO}(q) \right) \\ &= \arg \max_{q(z) \in \mathcal{Q}} \text{ELBO}(q)\end{aligned}$$

Extra: Variational Inference (2)

- What kind of Q should we use?
 - **Mean Field Approximation:** Choose Q to restrict that latent variables are mutually independent; i.e.,

$$q(z) = \prod_{j=1}^m q_j(z_j)$$

- It is crazy assumption, but often works well in practice
- **Pro:** Mean-field family is expressive because it can capture any marginal density of the latent variable
- **Con:** It cannot capture the correlation between latent variables

Extra: Variational Inference (3)

- Now, let's focus on j -th latent $q_j(z_j)$.

$$\begin{aligned} \text{ELBO}(q_j) &= \mathbb{E}_{z_j \sim q_j} \left[\underbrace{\mathbb{E}_{z_{-j}} \left\{ \log p(x, z_j, z_{-j}) \right\}}_{:=f(z_j)} \right] - \mathbb{E}_{z_j \sim q_j} [\log q_j(z_j)] + \text{const} \\ &= -\mathbb{E}_{z_j \sim q_j} \left[\log \frac{q_j(z_j)}{\exp(f(z_j))} \right] + \text{const} \\ &= -\text{KL}(q_j(z_j) \parallel \exp(f(z_j))) + \text{const} \end{aligned}$$

- Thus, when ELBO is maximized,

$$q_j^*(z_j) \propto \exp \left\{ \mathbb{E}_{z_{-j}} \left\{ \log p(x, z_j, z_{-j}) \right\} \right\}$$

Extra: Variational Inference (4)

Coordinate ascent mean-field variational inference (CAVI)

- **Input:** A model $p(x, z)$, data $p(z)$
- **Output:** Variational density $q(z)$
- **Algorithm**
 1. for each $j \in 1, \dots, m$, set

$$q_j(z_j) \propto \exp \left\{ \mathbb{E}_{z_{-j}} \left\{ \log p(x, z_j, z_{-j}) \right\} \right\}$$

- Each coordinate density is conditional on all the other coordinates
 - Similar to Gibbs sampling
 - 2. Compute ELBO(q) and repeat until convergence
 - **Caveat:** ELBO is a non-convex objective function, which is only guaranteed to converge to a local optimum
- Sensitive to the initialization

Extra: From VI to MCMC

- In both VI and MCMC, the goal is often to approximate a complicated distribution
- In Bayesian latent variable models, this target is usually the posterior

$$p^*(\mathbf{z}) = p(\mathbf{z} \mid \mathbf{X}).$$

- But for MCMC theory, we do not need to write the observed data \mathbf{X} explicitly
- So, from now on, we simply write the target distribution as

$$p^*(\mathbf{z}).$$

- **VI**: approximate $p^*(\mathbf{z})$ by optimizing over $q(\mathbf{z})$
- **MCMC**: generate samples from $p^*(\mathbf{z})$ by constructing a Markov chain whose stationary distribution is $p^*(\mathbf{z})$

Extra: Theory of MCMC (1)

- A first-order Markov chain is defined as

$$p(\mathbf{z}^{m+1} \mid \mathbf{z}^1, \dots, \mathbf{z}^m) = p(\mathbf{z}^{m+1} \mid \mathbf{z}^m)$$

- Thus, by the law of total probability,

$$p(\mathbf{z}^{m+1}) = \sum_{\mathbf{z}^m} \underbrace{p(\mathbf{z}^{m+1} \mid \mathbf{z}^m)}_{\text{Transition Prob}} p(\mathbf{z}^m)$$

- **Detailed Balancing:** With transition probabilities $T(\mathbf{z}, \mathbf{z}')$,

$$p^*(\mathbf{z})T(\mathbf{z}, \mathbf{z}') = p^*(\mathbf{z}')T(\mathbf{z}', \mathbf{z})$$

- Detailed balancing is a *sufficient* condition to ensure that the required distribution $p(\mathbf{z})$ is stationary

Extra: Theory of MCMC (2)

- The detailed balancing implies stationarity because of the following:

$$\begin{aligned} p(z) &= \sum_{z'} p^{prev}(z') T(z', z) \\ &= \sum_{z'} p^{prev}(z) T(z, z') \quad (\text{detailed balancing}) \\ &= p^{prev}(z) \sum_{z'} T(z, z') \\ &= p^{prev}(z) \underbrace{\sum_{z'} p(z' | z)}_{=1} \\ &= p^{prev}(z) \end{aligned}$$

Extra: Theory of MCMC (3)

- We want our chains to converge (in the limit) to the same stationary state, regardless of starting distribution.
 - Such chains are called **ergodic**.
- For finite state spaces, our chains converge to equilibrium under two relatively weak conditions
 1. **Irreducibility**: Any set of states can be reached from any other state in a finite number of moves → Stationary distribution is unique
 2. **Aperiodicity**: The chain does not get trapped in cycles
- Rarely justified theoretically
 - Instead, we can check it empirically by running the estimation multiple times (chains) and checking the convergence (e.g., trace plot / \hat{R})

Extra: Metropolis-Hasting (M-H) algorithm

1. **Initialization:** Choose an arbitrary point z'
2. For each iteration at step m ,
 - 2.1 **Propose** a candidate z^* for the next sample by picking from a proposal distribution $q(z^* | z^m)$
 - 2.2 **Accept** the sample if $u \leq \alpha(z^*, z^m)$, where $u \sim \text{Unif}(0, 1)$ and

$$\alpha(z^*, z^m) = \min\left(1, \frac{\tilde{p}(z^*)q(z^m | z^*)}{\tilde{p}(z^m)q(z^* | z^m)}\right)$$

- This works because it satisfies the detailed balancing:

$$\begin{aligned} p(\mathbf{z})q(\mathbf{z} | \mathbf{z}')\alpha(\mathbf{z}', \mathbf{z}) &= \min\left\{p(\mathbf{z})q(\mathbf{z} | \mathbf{z}'), p(\mathbf{z}')q(\mathbf{z}' | \mathbf{z})\right\} \\ &= \min\left\{p(\mathbf{z}')q(\mathbf{z}' | \mathbf{z}), p(\mathbf{z})q(\mathbf{z} | \mathbf{z}')\right\} = p(\mathbf{z}')q(\mathbf{z}' | \mathbf{z})\alpha(\mathbf{z}, \mathbf{z}') \end{aligned}$$

Extra: Metropolis algorithm

- Suppose that your proposal distribution is symmetric: i.e.,

$$q(y | x) = q(x | y)$$

- Example: $q(x | y) = q((x - y)^2)$
- Then, your ratio is

$$\alpha(z^*, z^m) = \min\left(1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^m)}\right)$$

- This is called **Metropolis algorithm**
- **Problem (both M and M-H)**: The acceptance rate may be low, sensitive to the choice of $q(\cdot | \cdot)$, does not work well for high-dimensional settings.

Extra: Gibbs Sampling

- **Setting:** Suppose we want to get $p(z_1, \dots, z_k)$
- **Algorithm:** After initialization, for each step m , keep sampling

$$z_1^{m+1} \sim p(z_1 \mid z_2^m, \dots, z_k^m)$$

$$z_2^{m+1} \sim p(z_2 \mid z_1^{m+1}, z_3^m, \dots, z_k^m)$$

⋮

$$z_k^{m+1} \sim p(z_k \mid z_1^{m+1}, z_2^{m+1}, \dots, z_{k-1}^{m+1})$$

- **AD:** Always accept the sample
- **DA:** This approach works only when full conditional distributions are easy to sample

Extra: Gibbs Sampling as Special Case of M-H

- If we use the proposal $q(z^* | z^m) = p(z_i^* | z_{-i}^m)$ and $q(z^m | z^*) = p(z_i^m | z_{-i}^*)$, then Gibbs Sampling can be seen as a special case of M-H
- In Gibbs Sampling, we always accept the proposal because

$$\begin{aligned}\alpha(z^*, z^m) &= \min\left(1, \frac{\tilde{p}(z^*)q(z^m | z^*)}{\tilde{p}(z^m)q(z^* | z^m)}\right) \quad (\text{Definition}) \\ &= \min\left(1, \frac{\tilde{p}(z_{-i}^*)\tilde{p}(z_i^* | z_{-i}^*)p(z_i^m | z_{-i}^*)}{p(z_{-i}^m)\tilde{p}(z_i^m | z_{-i}^m)p(z_i^* | z_{-i}^*)}\right) \\ &= \min(1, 1) = 1\end{aligned}$$

where the last equality is by $z_{-i}^* = z_{-i}^m$ (as only i changes)

Extra: Metropolis within Gibbs

- Gibbs sampling is sometimes hard to implement when full-conditional is not known
 - M-H can solve this problem, but does not work well for high-dimensional setting
 - This is because M-H considers the update of all the parameters at the same time as marginal
- We instead consider doing the M-H for each full conditional in Gibbs sampling
- This is called **Metropolis within Gibbs**

Extra: Collapsed Gibbs Sampling

- We can speed up the Gibbs sampler by marginalizing some of the latent variables.
 - This is called **Collapsed Gibbs Sampler**
 - **AD**: Faster convergence (b/c fewer parameters to be sampled)
- **Setup**: let's split the latent variable into
 1. z : variables to sample
 2. θ : variables to collapse
- Then, instead of sample from $p(z, \theta)$, we can integrate out θ :

$$p(z | x) = \int p(z, \theta | x) d\theta$$

Then, we run Gibbs sampler updates like:

$$z_i^{(t+1)} \sim p(z_i | z_{-i}, x)$$

under the marginal posterior

Extra: Hamiltonian/Hybrid Monte Carlo (HMC)

- **Idea:** To use momentum to update the states to mitigate the random walk behavior of GS and MH:
 - x : position variable / u : momentum variable
 - Consider

$$\underbrace{H(x_i, u_i)}_{\text{Hamiltonian}} = \underbrace{U(x_i)}_{\text{Potential Energy}} + \underbrace{K(u_i)}_{\text{Kinetic Energy}} = -\log p(x) + \frac{1}{2}u^T M^{-1}u$$

which satisfies Hamiltonian Equations:

$$\frac{\partial x}{\partial t} = \frac{\partial H}{\partial u}, \quad \frac{\partial u}{\partial t} = -\frac{\partial H}{\partial x}$$

Extra: Hamiltonian/Hybrid Monte Carlo (HMC)

- For each iteration at step m ,
 1. Updating momentum distribution $u_0 \sim N(0, M)$
 2. For $l = 1, \dots, L - 1$, updating (x, u) by

$$u^{l+1/2} \leftarrow u^l + \frac{1}{2}\epsilon \frac{d \log(x | y)}{dx}$$
$$x^{l+1} \leftarrow x^l + \epsilon M^{-1} u^l$$
$$u^{l+1} \leftarrow u^{l+1/2} + \frac{1}{2}\epsilon \frac{d \log(x | y)}{dx}$$

3. Accept (u^L, x^L) with probability

$$\min\left(\frac{p(x^* | y)p(u^*)}{p(x^{t-1} | y)p(u^{t-1})}, 1\right)$$

- **Takeaway:** HMC is one type of M-H methods!
 - Notice that we accept or reject in the same way as M-H