

Latent Variable Models and EM

Hybrid Inference (posterior infer. + point est.)

- In many models, we infer posterior on some unknowns and do point est. for others
- We have already seen MLE-II for lin reg. which alternates between
 - Inferring Conditional Posterior (CP) over the main parameter given the point estimates of hyperparams

CP of w : $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \hat{\lambda}, \hat{\beta})$
 - Maximizing the marginal lik. to do point estimation for hyperparams

$\{\hat{\lambda}, \hat{\beta}\} = \operatorname{argmax}_{\lambda, \beta} p(\mathbf{y}|\mathbf{X}, \lambda, \beta)$
- The Expectation-Maximization algorithm (will see today) also does something similar
 - In E step, the CP of latent variables is inferred, given current point-est of params
 - M step maximizes **expected complete data log-lik.** to get point estimates of params
- If we can't (due to computational or other reasons) infer posterior over all unknowns, how to decide which variables to infer posterior on, and for which to do point-est?
- Usual approach: Infer **posterior over local vars** and **point estimates for global vars**
 - Reason: We typically have plenty of data to reliably estimate the global variables so it is okay even if we just do point estimation for those

Nomenclature/Notation Alert

- Why call some unknowns as **parameters** and others as **latent variables**?
- Well, no specific reason. Sort of a convention adopted by some algorithms
 - EM: Unknowns estimated in E step referred to as latent vars; those in M step as params
 - Usually: **Latent vars – (Conditional) posterior computed**; **parameters – point estimation**
- Some algos won't make such distinction and will infer posterior over all unknowns
- Sometimes the “global” or “local” unknown distinction makes it clear
 - Local variables = latent variables, global variables = parameters
- But remember that this nomenclature isn't really cast in stone, no need to be confused so long as you are clear as to what the role of each unknown is, and how we want to estimate it (posterior or point estimate) and using what type of inference algorithm

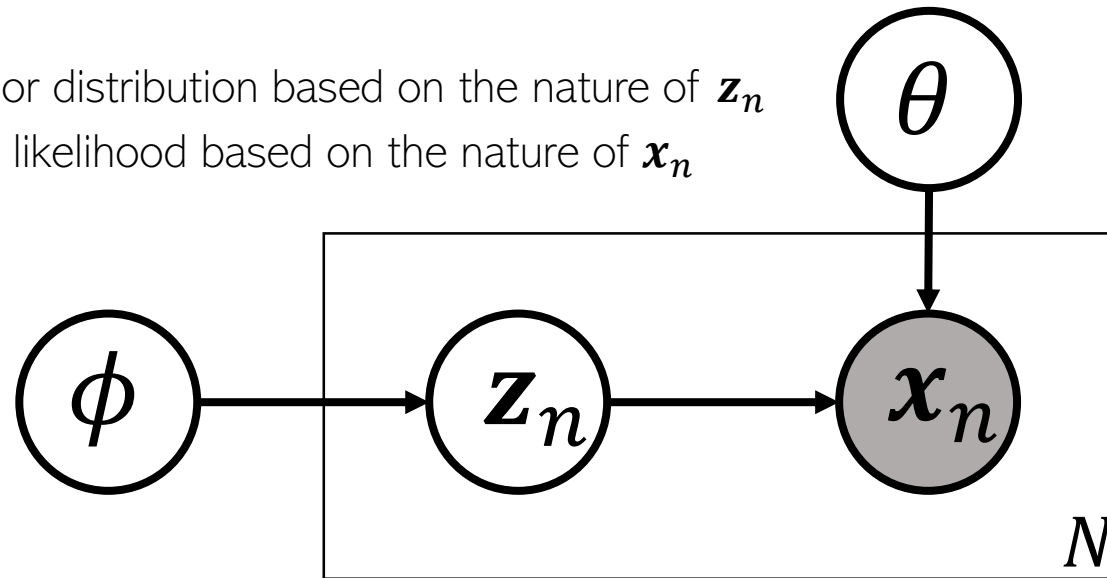
Inference/Parameter Estimation in Latent Variable Models using Expectation-Maximization (EM)

Parameter Estimation in Latent Variable Models

- Assume each observation \mathbf{x}_n to be associated with a “local” latent variable \mathbf{z}_n

$p(\mathbf{z}_n|\phi)$: A suitable prior distribution based on the nature of \mathbf{z}_n

$p(\mathbf{x}_n|\mathbf{z}_n, \theta)$: A suitable likelihood based on the nature of \mathbf{x}_n



- Although we can do fully Bayesian inference for all the unknowns, suppose we only want a point estimate of the “global” parameters $\Theta = (\theta, \phi)$ via MLE/MAP
- Such MLE/MAP problems in LVMs are difficult to solve in a “clean” way
 - Would typically require **gradient based methods** with no closed form updates for Θ
 - However, **EM** gives a clean way to obtain closed form updates for Θ

Why MLE/MAP of Params is Hard for LVMs?

- Suppose we want to estimate $\Theta = (\theta, \phi)$ via MLE. If we knew \mathbf{z}_n , we could solve

$$\Theta_{MLE} = \arg \max_{\Theta} \sum_{n=1}^N \log p(\mathbf{x}_n, \mathbf{z}_n | \Theta) = \arg \max_{\Theta} \sum_{n=1}^N [\log p(\mathbf{z}_n | \phi) + \log p(\mathbf{x}_n | \mathbf{z}_n, \theta)]$$

Easy to solve

In particular, if they are exp-fam distributions

- Easy. Usually closed form if $p(\mathbf{z}_n | \phi)$ and $p(\mathbf{x}_n | \mathbf{z}_n, \theta)$ have simple forms

- However, since in LVMs, \mathbf{z}_n is hidden, the MLE problem for Θ will be the following

Basically, the **marginal likelihood** after integrating out \mathbf{z}_n

$$\Theta_{MLE} = \arg \max_{\Theta} \sum_{n=1}^N \log p(\mathbf{x}_n | \Theta) = \arg \max_{\Theta} \log p(\mathbf{X} | \Theta)$$

- $\log p(\mathbf{x}_n | \Theta)$ will not have a simple expression since $p(\mathbf{x}_n | \Theta)$ requires sum/integral

$$p(\mathbf{x}_n | \Theta) = \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n | \Theta) \quad \dots \text{ or if } \mathbf{z}_n \text{ is continuous: } p(\mathbf{x}_n | \Theta) = \int p(\mathbf{x}_n, \mathbf{z}_n | \Theta) d\mathbf{z}_n$$

- MLE now becomes difficult (basically MLE-II now), no closed form expression for Θ .

- Can we maximize some other quantity instead of $\log p(\mathbf{x}_n | \Theta)$ for this MLE?

MLE-II

Sometimes our model has two layers of unknowns:

- **Parameters / weights** w (local variables).
- **Hyperparameters** λ, β (global parameters controlling prior distributions, variances, etc.).
- Example: **Bayesian linear regression**
 - Data likelihood: $p(y \mid \mathbf{X}, w, \beta)$
 - Prior: $p(w \mid \lambda)$
 - Unknowns: w, β, λ
- If we were fully Bayesian, we'd put priors on everything and integrate

MLE-II

In **MLE-II**

- Integrate out the local parameter w to compute the **marginal likelihood** of hyperparameters:

$$p(y \mid \mathbf{X}, \lambda, \beta) = \int p(y \mid \mathbf{X}, w, \beta) p(w \mid \lambda) dw$$

- Then maximize w.r.t. hyperparameters:

$$(\hat{\lambda}, \hat{\beta}) = \operatorname{argmax}_{\lambda, \beta} p(y \mid \mathbf{X}, \lambda, \beta)$$

- This is **MLE-II** (Maximum Likelihood Estimation, stage II).

MLE-I: fit parameters (w) by maximizing likelihood of data.

MLE-II: fit hyperparameters (λ, β) by maximizing the marginal likelihood (a.k.a. “evidence”) after integrating out parameters w .

An Important Identity

- Assume $p_z = p(\mathbf{Z}|\mathbf{X}, \Theta)$ and $q(\mathbf{Z})$ to be some prob distribution over \mathbf{Z} , then

$$\log p(\mathbf{X}|\Theta) = \mathcal{L}(q, \Theta) + KL(q||p_z)$$

Assume \mathbf{Z} discrete

Verify the identity

- In the above $\mathcal{L}(q, \Theta) = \sum_Z q(Z) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(Z)} \right\}$

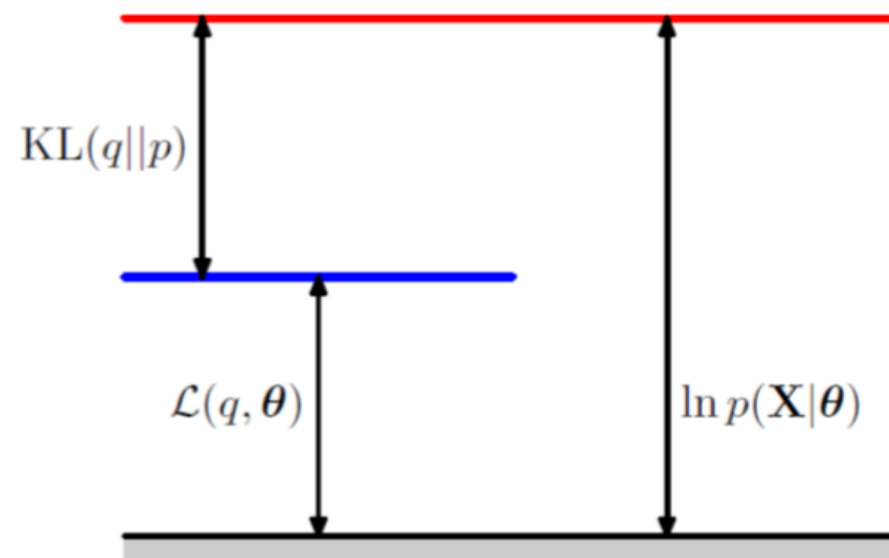
- $KL(q||p_z) = -\sum_Z q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \Theta)}{q(\mathbf{Z})} \right\}$

- KL is always non-negative, so $\log p(\mathbf{X}|\Theta) \geq \mathcal{L}(q, \Theta)$

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

- Thus if we maximize $\mathcal{L}(q, \Theta)$, it will also improve $\log p(\mathbf{X}|\Theta)$

- Also, as we'll see, it's easier to maximize $\mathcal{L}(q, \Theta)$



Lower bound

KL divergence

Derivation

$$\text{KL}(q(Z) \parallel p(Z|X, \Theta)) = \sum_Z q(Z) \log \frac{q(Z)}{p(Z|X, \Theta)}.$$

Intuition

KL is a measure of how far our guess distribution $q(Z)$ is from the true posterior $p(Z|X, \Theta)$. If they're equal, $\text{KL} = 0$.

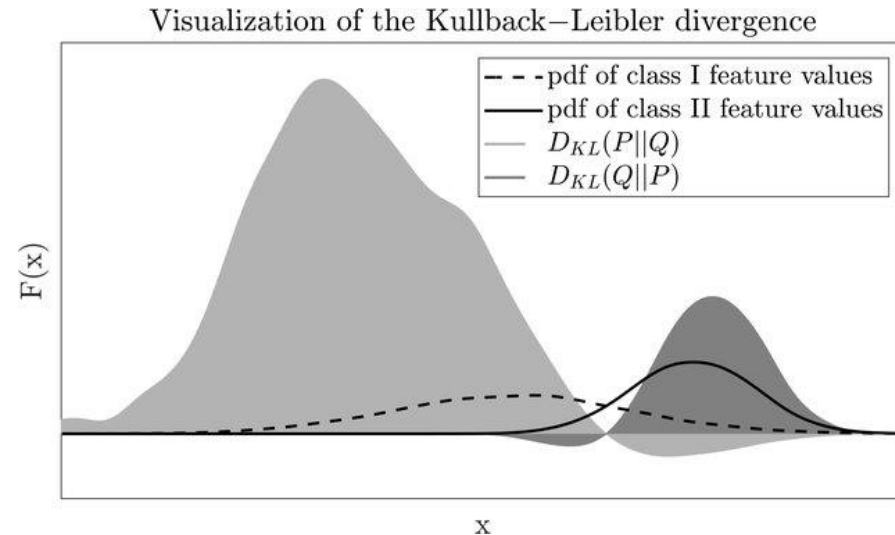
Deriving Lower bound

KL divergence

$$\text{KL}(q(Z) \parallel p(Z|X, \Theta)) = \sum_Z q(Z) \log \frac{q(Z)}{p(Z|X, \Theta)}.$$

Intuition

KL is a measure of how far our guess distribution $q(Z)$ is from the true posterior $p(Z|X, \Theta)$. If they're equal, $\text{KL} = 0$.



Deriving Lower bound

Substitute Bayes' rule for the posterior

$$p(Z|X, \Theta) = \frac{p(X, Z|\Theta)}{p(X|\Theta)}.$$

So,

$$\text{KL}(q||p) = \sum_Z q(Z) \log \frac{q(Z) p(X|\Theta)}{p(X, Z|\Theta)}$$

We replace the “posterior” with things we know:

- joint $p(X, Z | \Theta)$ and
- evidence (marginal likelihood) $p(X | \Theta)$

Deriving Lower bound

Expand the terms

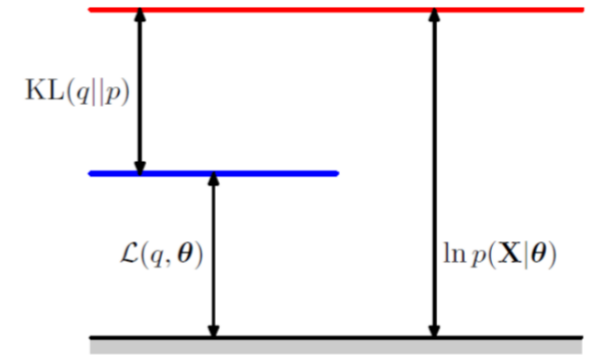
$$\text{KL}(q||p) = \sum_Z q(Z) \log q(Z) - \sum_Z q(Z) \log p(X, Z|\Theta) + \log p(X|\Theta)$$

Intuition : We now have three pieces

- an entropy-like term $\sum q \log q$
- a data-fit term $\sum q \log p(X, Z | \Theta)$
- Constant log-likelihood of the data $\log p(X | \Theta)$

Deriving Lower bound

Rearrange to isolate $\log p(X | \Theta)$



$$\log p(X|\Theta) = \underbrace{\sum_Z q(Z) [\log p(X, Z|\Theta) - \log q(Z)]}_{\mathcal{L}(q, \Theta)} + \text{KL}(q||p)$$

Intuition : We've decomposed the log-likelihood into two parts:

- 1.ELBO ($\mathcal{L}(q, \Theta)$)** – how well our chosen $q(Z)$ explains the joint data + latent variables, adjusted by its own complexity.
- 2.KL gap** – the “price” we pay because $q(Z)$ isn't the true posterior.

Maximizing $\mathcal{L}(q, \Theta)$

Basically, log of marginal likelihood w.r.t. Θ with \mathbf{Z} integrated out

$\log p(\mathbf{X}|\Theta)$ is called **Incomplete-Data Log Likelihood (ILL)**



- $\mathcal{L}(q, \Theta)$ depends on q and Θ . We'll use ALT-OPT to maximize it
- Let's maximize $\mathcal{L}(q, \Theta)$ w.r.t. q with Θ fixed at some Θ^{old}

Since $\log p(\mathbf{X}|\Theta) = \mathcal{L}(q, \Theta) + KL(q||p_z)$ is constant when Θ is held fixed at Θ^{old}

$$\hat{q} = \operatorname{argmax}_q \mathcal{L}(q, \Theta^{\text{old}}) = \operatorname{argmin}_q KL(q||p_z) = p_z = p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})$$

- Now let's maximize $\mathcal{L}(q, \Theta)$ w.r.t. Θ with q fixed at $\hat{q} = p_z = p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})$

The posterior distribution of \mathbf{Z} given current parameters Θ^{old}

$$\Theta^{\text{new}} = \operatorname{argmax}_{\Theta} \mathcal{L}(\hat{q}, \Theta) = \operatorname{argmax}_{\Theta} \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})} \right\}$$

Maximization of **expected CLL** where the expectation is w.r.t. the posterior distribution of \mathbf{Z} given current parameters Θ^{old}

$$= \operatorname{argmax}_{\Theta} \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

$$= \operatorname{argmax}_{\Theta} \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})} [\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$$

Complete-Data Log Likelihood (CLL)

$$= \operatorname{argmax}_{\Theta} Q(\Theta, \Theta^{\text{old}})$$

Much easier than maximizing ILL since CLL will have simple expressions (since it is akin to knowing \mathbf{Z})

The Expectation-Maximization (EM) Algorithm

- ALT-OPT of $\mathcal{L}(q, \Theta)$ w.r.t. q and Θ gives the EM algorithm (Dempster, Laird, Rubin, 1977)

The EM Algorithm

Primarily designed for doing point estimation of the parameters Θ but also gives (CP of) latent variables z_n

Usually computing CP + expected CLL is referred to as the **E step**, and max. of exp-CLL w.r.t. Θ as the **M step**



① Initialize Θ as $\Theta^{(0)}$, set $t = 1$

② Step 1: Compute **posterior** of latent variables given current parameters $\Theta^{(t-1)}$

Conditional posterior of each latent variable z_n

Latent variables also assumed indep. a priori

$$p(z_n^{(t)} | x_n, \Theta^{(t-1)}) = \frac{p(z_n^{(t)} | \Theta^{(t-1)}) p(x_n | z_n^{(t)}, \Theta^{(t-1)})}{p(x_n | \Theta^{(t-1)})} \propto \text{prior} \times \text{likelihood}$$

③ Step 2: Now maximize the **expected complete data log-likelihood** w.r.t. Θ

Assuming the (expected) CLL $\mathbb{E}_{p(\mathbf{Z} | \mathbf{X}, \Theta^{\text{old}})} [\log p(\mathbf{X}, \mathbf{Z} | \Theta)]$ factorizes over all observations

$$\Theta^{(t)} = \arg \max_{\Theta} \mathcal{Q}(\Theta, \Theta^{(t-1)}) = \arg \max_{\Theta} \sum_{n=1}^N \mathbb{E}_{p(z_n^{(t)} | x_n, \Theta^{(t-1)})} [\log p(x_n, z_n^{(t)} | \Theta)]$$

④ If not yet converged, set $t = t + 1$ and go to step 2.

- Note: If we can take the MAP estimate \hat{z}_n of z_n (not full posterior) in Step 1 and maximize the CLL in Step 2 using that, i.e., do $\arg \max_{\Theta} \sum_{n=1}^N [\log p(x_n, \hat{z}_n^{(t)} | \Theta)]$ this will be ALT-OPT

The Expected CLL

- Expected CLL in EM is given by (assume observations are i.i.d.)

$$\begin{aligned} \mathcal{Q}(\Theta, \Theta^{old}) &= \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n | \mathbf{x}_n, \Theta^{old})} [\log p(\mathbf{x}_n, \mathbf{z}_n | \Theta)] \\ &= \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n | \mathbf{x}_n, \Theta^{old})} [\log p(\mathbf{x}_n | \mathbf{z}_n, \Theta) + \log p(\mathbf{z}_n | \Theta)] \end{aligned}$$

- If $p(\mathbf{z}_n | \Theta)$ and $p(\mathbf{x}_n | \mathbf{z}_n, \Theta)$ are exp-family distributions, $\mathcal{Q}(\Theta, \Theta^{old})$ has a very simple form
- In resulting expressions, replace terms containing \mathbf{z}_n 's by their respective expectations, e.g.,
 - \mathbf{z}_n replaced by $\mathbb{E}_{p(\mathbf{z}_n | \mathbf{x}_n, \hat{\Theta})} [\mathbf{z}_n]$
 - $\mathbf{z}_n \mathbf{z}_n^\top$ replaced by $\mathbb{E}_{p(\mathbf{z}_n | \mathbf{x}_n, \hat{\Theta})} [\mathbf{z}_n \mathbf{z}_n^\top]$
- However, in some LVMs, these expectations are intractable to compute and need to be approximated (will see some examples later)

What's Going On?

- As we saw, the maximization of lower bound $\mathcal{L}(q, \Theta)$ had two steps
- Step 1 finds the optimal q (call it \hat{q}) by setting it as the posterior of \mathbf{Z} given current Θ
- Step 2 maximizes $\mathcal{L}(\hat{q}, \Theta)$ w.r.t. Θ which gives a new Θ .

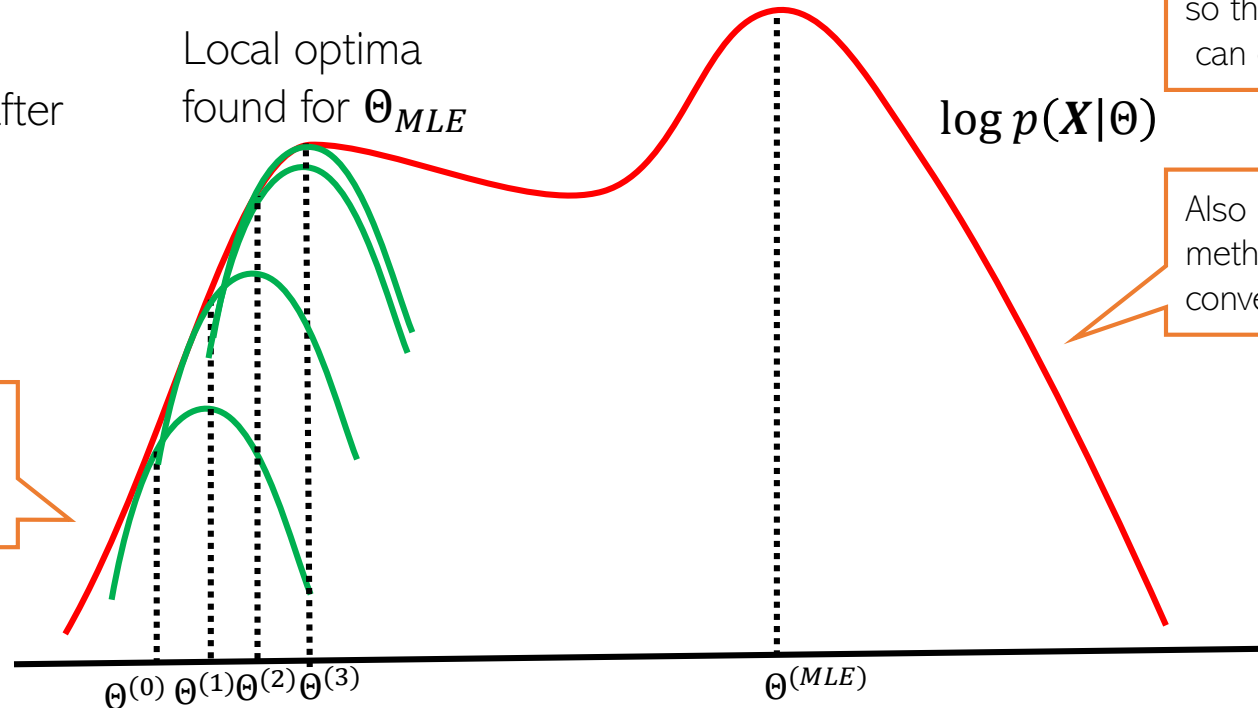
Alternating between them until convergence to some local optima

KL becomes zero and $\mathcal{L}(q, \Theta)$ becomes equal to $\log p(\mathbf{X}|\Theta)$; thus their curves touch at current Θ

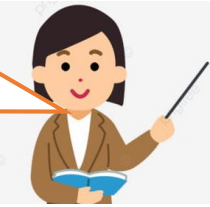
Green curve: $\mathcal{L}(\hat{q}, \Theta)$ after setting q to \hat{q}

Local optima found for Θ_{MLE}

Good initialization matters; otherwise would converge to a poor local optima



Note that Θ only changes in Step 2 so the objective $\log p(\mathbf{X}|\Theta)$ can only change in Step 2



Also kind of similar to Newton's method (and has second order like convergence behavior in some cases)

Unlike Newton's method, we don't construct and optimize a quadratic approximation, but a lower bound

Even though original MLE problem $\text{argmax}_{\Theta} \log p(\mathbf{X}|\Theta)$ could be solved using gradient methods, EM often works faster and has cleaner updates

EM vs Gradient-based Methods

- Can also estimate params using gradient-based optimization instead of EM
 - We can usually explicitly sum over or integrate out the latent variables \mathbf{Z} , e.g.,

$$\mathcal{L}(\Theta) = \log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta)$$

- Now we can optimize $\mathcal{L}(\Theta)$ using first/second order optimization to find the optimal Θ
- EM is usually preferred over this approach because
 - The M step has often simple closed-form updates for the parameters Θ
 - Often constraints (e.g., PSD matrices) are automatically satisfied due to form of updates
 - In some cases[†], EM usually converges faster (and often like second-order methods)
 - E.g., Example: Mixture of Gaussians with when the data is reasonably well-clustered
 - EM applies even when the explicit summing over/integrating out is expensive/intractable
 - EM also provides the conditional posterior over the latent variables \mathbf{Z} (from E step)

[†]Optimization with EM and Expectation-Conjugate-Gradient (Salakhutdinov et al, 2003), On Convergence Properties of the EM Algorithm for Gaussian Mixtures (Xu and Jordan, 1996), Statistical guarantees for the EM algorithm: From population to sample-based analysis (Balakrishnan et al, 2017)

Some Applications of EM

- Mixture Models and Dimensionality Reduction/Representation Learning
 - Mixture Models: Mixture of Gaussians, Mixture of Experts, etc
 - Dim. Reduction/Representation Learning: Probabilistic PCA, Variational Autoencoders
- Problems with missing features or missing labels (which are treated as latent variables)
 - $\hat{\Theta} = \operatorname{argmax}_{\Theta} \log p(\mathbf{x}^{obs} | \Theta) = \operatorname{argmax}_{\Theta} \log \int p([\mathbf{x}^{obs}, \mathbf{x}^{miss}] | \Theta) d\mathbf{x}^{miss}$
 - $\hat{\Theta} = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \log p(x_n, y_n | \Theta) + \sum_{n=N+1}^{N+M} \log \sum_{c=1}^K p(x_n, y_n = c | \Theta)$
- Hyperparameter estimation in probabilistic models (an alternative to MLE-II)
 - MLE-II estimates hyperparams by maximizing the marginal likelihood, e.g.,

$$\{\hat{\lambda}, \hat{\beta}\} = \operatorname{argmax}_{\lambda, \beta} p(\mathbf{y} | \mathbf{X}, \lambda, \beta) = \operatorname{argmax}_{\lambda, \beta} \int p(\mathbf{y} | \mathbf{w}, \mathbf{X}, \beta) p(\mathbf{w} | \lambda) d\mathbf{w}$$

For a Bayesian linear regression model

- With EM, can treat \mathbf{w} as latent var, and λ, β as “parameters”
 - E step will estimate the CP of \mathbf{w} given current estimates of λ, β
 - M step will re-estimate λ, β by maximizing the expected CLL

$$\mathbb{E}[\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}, \beta, \lambda)] = \mathbb{E}[\log p(\mathbf{y} | \mathbf{w}, \mathbf{X}, \beta) + \log p(\mathbf{w} | \lambda)]$$

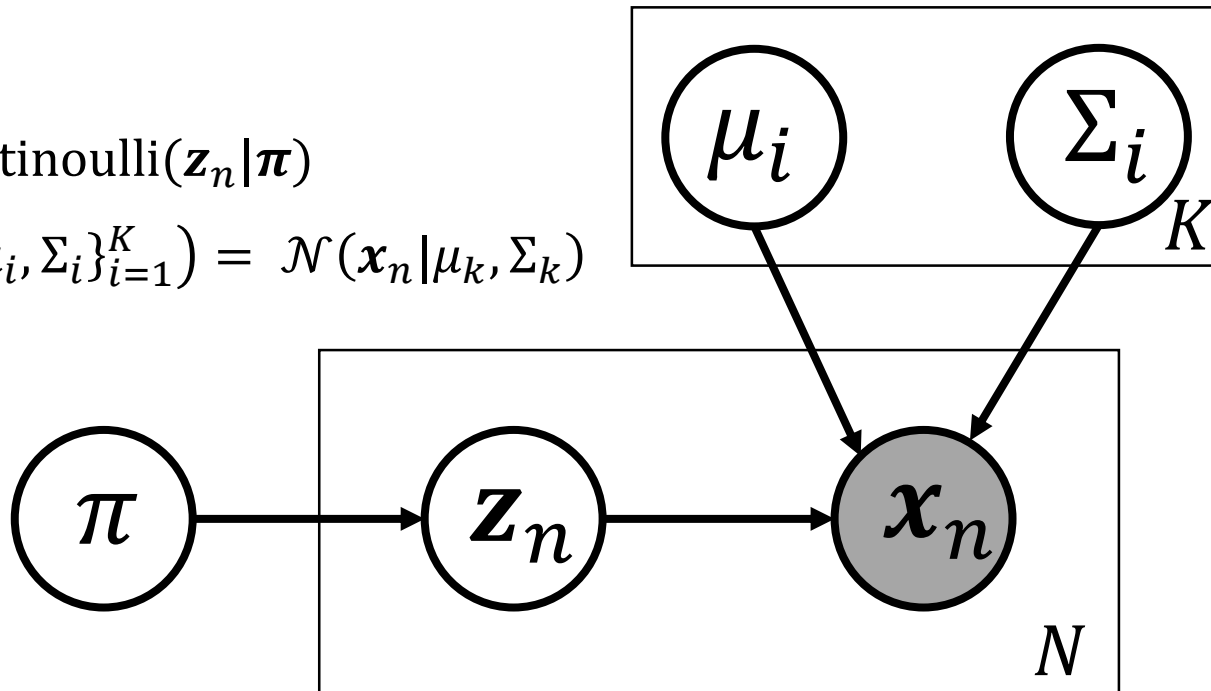
Expectations w.r.t.
the CP of \mathbf{w}

Gaussian Mixture Model (GMM)

- N observations $\{\mathbf{x}_n\}_{n=1}^N$ each from one of the K Gaussians $\{\mathcal{N}(\mu_i, \Sigma_i)\}_{i=1}^K$
- We don't know which Gaussian each observation \mathbf{x}_n comes from
- Assume $\mathbf{z}_n \in \{1, 2, \dots, K\}$ denotes which Gaussian generated \mathbf{x}_n
- Suppose we want to do point estimation for the parameters $\{\mu_i, \Sigma_i\}_{i=1}^K$

$$p(\mathbf{z}_n | \boldsymbol{\pi}) = \text{multinoulli}(\mathbf{z}_n | \boldsymbol{\pi})$$

$$p(\mathbf{x}_n | \mathbf{z}_n = k, \{\mu_i, \Sigma_i\}_{i=1}^K) = \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)$$



$$p(\mathbf{x}_n | \{\pi_i, \mu_i, \Sigma_i\}_{i=1}^K)$$

$$= \sum_{i=1}^K \pi_i \mathcal{N}(\mathbf{x}_n | \mu_i, \Sigma_i)$$

$$\log p(\mathbf{x}_n | \boldsymbol{\Theta}) = \log \sum_{i=1}^K \pi_i \mathcal{N}(\mathbf{x}_n | \mu_i, \Sigma_i)$$

Can use gradient based optimization for MLE of $\boldsymbol{\Theta}$ but the update equations are a bit complicated

EM would give simpler updates

Detour: MLE for GMM when \mathbf{Z} is known

GMM then is just like generative classification with Gaussian class conditionals

- Derivation of the MLE solution for $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ when \mathbf{Z} is known

$$\begin{aligned}\hat{\Theta} &= \operatorname{argmax}_{\Theta} p(\mathbf{X}, \mathbf{Z} | \Theta) = \operatorname{argmax}_{\Theta} \prod_{n=1}^N p(\mathbf{x}_n, \mathbf{z}_n | \Theta) \\ &= \operatorname{argmax}_{\Theta} \prod_{n=1}^N \underbrace{p(\mathbf{z}_n | \Theta)}_{\text{multinoulli}} \underbrace{p(\mathbf{x}_n | \mathbf{z}_n, \Theta)}_{\text{Gaussian}} \\ &= \operatorname{argmax}_{\Theta} \prod_{n=1}^N p(\mathbf{z}_n | \Theta) p(\mathbf{x}_n | \mathbf{z}_n, \Theta) \\ &= \operatorname{argmax}_{\Theta} \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \prod_{k=1}^K p(\mathbf{x}_n | \mathbf{z}_n = k, \Theta)^{z_{nk}} \\ &= \operatorname{argmax}_{\Theta} \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\mathbf{x}_n | \mathbf{z}_n = k, \Theta)]^{z_{nk}} \\ &= \operatorname{argmax}_{\Theta} \log \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\mathbf{x}_n | \mathbf{z}_n = k, \Theta)]^{z_{nk}} \\ &= \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K z_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)]\end{aligned}$$

In general, in models with probability distributions from the **exponential family**, the MLE problem will usually have a simple analytic form

Also, due to the form of the likelihood (Gaussian) and prior (multinoulli), the MLE problem had a nice separable structure after taking the log

Can see that, when estimating the parameters of the k^{th} Gaussian (π_k, μ_k, Σ_k) , we only will only need training examples from the k^{th} class, i.e., examples for which $z_{nk} = 1$

EM for Gaussian Mixture Model (GMM)

1. Initialize $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ as $\Theta^{(0)}$. Set $t = 1$
2. Set CP $q^{(t)} = p(\mathbf{Z}|\mathbf{X}, \Theta^{(t-1)})$. Assuming i.i.d. data, this means computing $\forall n, k$

Probability of data point n
belonging to the k -th Gaussian

"Soft-clustering"

Same as writing $z_n = k$

$$p(\mathbf{z}_{nk} = 1 | \mathbf{x}_n, \Theta^{(t-1)}) \propto p(\mathbf{z}_{nk} = 1 | \Theta^{(t-1)}) p(\mathbf{x}_n | \mathbf{z}_{nk} = 1, \Theta^{(t-1)})$$

$$= \pi_k^{(t-1)} \mathcal{N}(\mathbf{x}_n | \mu_k^{(t-1)}, \Sigma_k^{(t-1)})$$

3. Set $\Theta^{(t)} = \operatorname{argmax}_{\Theta} \mathbb{E}_{q^{(t)}} [\log p(\mathbf{X}, \mathbf{Z} | \Theta)] = \operatorname{argmax}_{\Theta} Q(\Theta, \Theta^{(t-1)})$

This only required expectation for EM
for GMM is $\mathbb{E}[\mathbf{z}_{nk}]$ which can be
computed easily using the CP of \mathbf{z}_n

EM for GMM does **two**
things: soft-clustering
and estimating the
density $p(\mathbf{X} | \Theta)$



$$\Theta^{(t)} = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n | \mathbf{x}_n, \Theta^{(t-1)})} [\log p(\mathbf{x}_n, \mathbf{z}_n | \Theta)]$$

$$= \operatorname{argmax}_{\Theta} \mathbb{E} \left[\sum_{n=1}^N \sum_{k=1}^K \mathbf{z}_{nk} \left[\log \pi_k^{(t-1)} + \log \mathcal{N}(\mathbf{x}_n | \mu_k^{(t-1)}, \Sigma_k^{(t-1)}) \right] \right]$$

$$= \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}[\mathbf{z}_{nk}] [\log \pi_k^{(t-1)} + \log \mathcal{N}(\mathbf{x}_n | \mu_k^{(t-1)}, \Sigma_k^{(t-1)})]$$

$N_k = \sum_{n=1}^N \mathbb{E}[\mathbf{z}_{nk}]$
denotes the effective
number of points
from k -th Gaussian

$$\pi_k^{(t)} = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[\mathbf{z}_{nk}]$$

$$\mu_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \mathbb{E}[\mathbf{z}_{nk}] \mathbf{x}_n$$

$$\Sigma_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \mathbb{E}[\mathbf{z}_{nk}] (\mathbf{x}_n - \mu_k^{(t)}) (\mathbf{x}_n - \mu_k^{(t)})^\top$$

4. Go to step 2 if not converged

EM for GMM: The Full Algorithm

- The EM algo for GMM required $\mathbb{E}[z_{nk}]$. Note $z_{nk} \in \{0,1\}$

Reason: $\sum_{k=1}^K \gamma_{nk} = 1$

Need to normalize: $\mathbb{E}[z_{nk}] = \frac{\hat{\pi}_k \mathcal{N}(x_n | \hat{\mu}_k, \hat{\Sigma}_k)}{\sum_{\ell=1}^K \hat{\pi}_\ell \mathcal{N}(x_n | \hat{\mu}_\ell, \hat{\Sigma}_\ell)}$

$$\mathbb{E}[z_{nk}] = \gamma_{nk} = 0 \times p(z_{nk} = 0 | x_n, \hat{\Theta}) + 1 \times p(z_{nk} = 1 | x_n, \hat{\Theta}) = p(z_{nk} = 1 | x_n, \hat{\Theta}) \propto \hat{\pi}_k \mathcal{N}(x_n | \hat{\mu}_k, \hat{\Sigma}_k)$$

EM for Gaussian Mixture Model

① Initialize $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ as $\Theta^{(0)}$, set $t = 1$

② E step: compute the expectation of each z_n (we need it in M step)

Soft K -means, which is more of a heuristic to get soft-clustering, also gave us probabilities but doesn't account for cluster shapes or fraction of points in each cluster

Accounts for fraction of points in each cluster

$$\mathbb{E}[z_{nk}^{(t)}] = \gamma_{nk}^{(t)} = \frac{\pi_k^{(t-1)} \mathcal{N}(x_n | \mu_k^{(t-1)}, \Sigma_k^{(t-1)})}{\sum_{\ell=1}^K \pi_\ell^{(t-1)} \mathcal{N}(x_n | \mu_\ell^{(t-1)}, \Sigma_\ell^{(t-1)})} \quad \forall n, k$$

Accounts for cluster shapes (since each cluster is a Gaussian)

③ Given “responsibilities” $\gamma_{nk} = \mathbb{E}[z_{nk}]$, and $N_k = \sum_{n=1}^N \gamma_{nk}$, re-estimate Θ via MLE

$$\mu_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk}^{(t)} x_n$$

Effective number of points in the k^{th} cluster

M-step:

$$\Sigma_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk}^{(t)} (x_n - \mu_k^{(t)})(x_n - \mu_k^{(t)})^\top$$

$$\pi_k^{(t)} = \frac{N_k}{N}$$

④ Set $t = t + 1$ and go to step 2 if not yet converged

Bayesian Linear Regression (Revisited)

$N \times D$ input matrix

$N \times 1$ responses

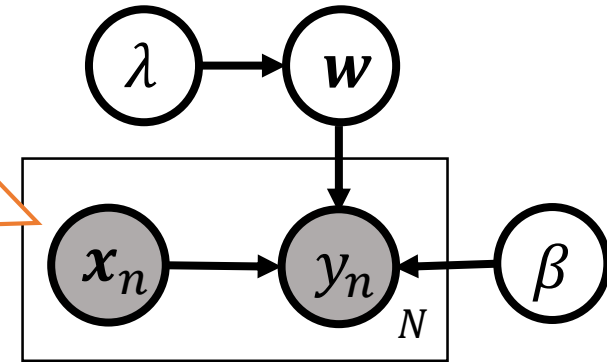
- N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$ from a lin-reg model with weights \mathbf{w}
- Suppose the hyperparameters are also unknown, so need to estimate $\mathbf{w}, \beta, \lambda$

$$p(y_n | \mathbf{x}_n, \mathbf{w}, \beta) = \mathcal{N}(y_n | \mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) \quad p(\mathbf{w} | \lambda) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda^{-1} \mathbf{I})$$

CP of \mathbf{w} : $p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$\boldsymbol{\Sigma} = (\beta \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \quad \boldsymbol{\mu} = \beta \mathbf{A}^{-1} \mathbf{X}^\top \mathbf{y}$$

In this latent variable model, there are no local variables. $\mathbf{w}, \beta, \lambda$ are all "global"



Many ways to optimize the marginal likelihood in MLE-II, e.g., gradient descent

MLE-II $(\hat{\beta}, \hat{\lambda}) = \operatorname{argmax}_{\beta, \lambda} \log p(\mathbf{y} | \mathbf{X}, \beta, \lambda)$

EM solves the MLE-II problem by optimizing a lower bound on the log marginal likelihood, and gives simple update equations for β, λ

EM

Expected CLL

$$\begin{aligned} (\hat{\beta}, \hat{\lambda}) &= \operatorname{argmax}_{\beta, \lambda} \mathbb{E}_{p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta^{(t-1)}, \lambda^{(t-1)})} [\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}, \beta, \lambda)] \\ &= \operatorname{argmax}_{\beta, \lambda} \mathbb{E}_{p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta^{(t-1)}, \lambda^{(t-1)})} [\log p(\mathbf{y} | \mathbf{w}, \mathbf{X}, \beta) + \log p(\mathbf{w} | \lambda)] \end{aligned}$$

Data

\mathbf{w} treated as latent variable here

EM for Bayesian Linear Regression

$$(\beta^{(t)}, \lambda^{(t)}) = \operatorname{argmax}_{\beta, \lambda} \mathbb{E}[\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}, \beta^{(t-1)}, \lambda^{(t-1)})]$$

1. Initialize β as $\beta^{(0)}$ and λ as $\lambda^{(0)}$. Set $t = 1$

2. Update the CP of \mathbf{w} as

$$p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta^{(t-1)}, \lambda^{(t-1)}) = \mathcal{N}(\boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)})$$

$$\boldsymbol{\Sigma}^{(t)} = (\beta^{(t-1)} \mathbf{X}^\top \mathbf{X} + \lambda^{(t-1)} \mathbf{I})^{-1} \quad \boldsymbol{\mu}^{(t)} = \beta^{(t-1)} \boldsymbol{\Sigma}^{(t)} \mathbf{X}^\top \mathbf{y}$$

3. Update β and λ as

$$\lambda^{(t)} = \frac{D}{\mathbb{E}[\mathbf{w}^\top \mathbf{w}]} = \frac{D}{\boldsymbol{\mu}^{(t)\top} \boldsymbol{\mu}^{(t)} + \operatorname{trace}(\boldsymbol{\Sigma}^{(t)})}$$

Intuition: If weights are large, λ decreases (weaker prior). If weights are small, λ increases (stronger shrinkage).

$$\beta^{(t)} = \frac{N}{\|\mathbf{y} - \mathbf{X} \boldsymbol{\mu}^{(t)}\|^2 + \operatorname{trace}(\mathbf{X}^\top \boldsymbol{\Sigma}^{(t)} \mathbf{X})}$$

Intuition: β (noise precision) increases when residuals are small (model fits well), decreases when residuals are large (data is noisy).

4. If not converged, set $t = t + 1$ and go to step 2

EM: Some other examples

- Problems with missing features (which are treated as latent variables)
 - Suppose each input \mathbf{x}_n has two parts - observed and missing: $\mathbf{x}_n = [\mathbf{x}_n^{obs}, \mathbf{x}_n^{miss}]$
 - For such problems, MLE for a model $p(\mathbf{X}|\Theta)$, assuming i.i.d. data, would have the form

$$\hat{\Theta} = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \log p(\mathbf{x}_n^{obs} | \Theta)$$

Suppose we are estimating the mean/covariance of a multivariate Gaussian given N input, with some inputs observations may have missing features

$$= \operatorname{argmax}_{\Theta} \sum_{n=1}^N \log \int p([\mathbf{x}_n^{obs}, \mathbf{x}_n^{miss}] | \Theta) d\mathbf{x}_n^{miss}$$

- Here \mathbf{x}_n^{miss} can be treated as a latent variable
- The CP will be $p(\mathbf{x}_n^{miss} | \mathbf{x}_n^{obs}, \Theta)$
- Using the CP, compute expected CLL and maximize it w.r.t. Θ
- Problems with missing labels (which are treated as latent variables)

An example of semi-supervised learning

This part is like GMM, thus EM can be used

$$\hat{\Theta} = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \log p(x_n, y_n | \Theta) + \sum_{n=N+1}^{N+M} \log \sum_{c=1}^K p(x_n, y_n = c | \Theta)$$

EM when CP and/or expectation is intractable

- EM solves the following step for estimating Θ

$$\Theta^{(t)} = \operatorname{argmax}_{\Theta} \mathbb{E}_{q^{(t)}}[\log p(\mathcal{D}, \mathbf{Z}|\Theta)] = \operatorname{argmax}_{\Theta} \int \log p(\mathcal{D}, \mathbf{Z}|\Theta) p(\mathbf{Z}|\Theta^{(t-1)}, \mathcal{D}) d\mathbf{Z}$$

- The above problem may be difficult to solve if one/both of the following is true

1. CP $p(\mathbf{Z}|\Theta^{(t-1)}, \mathcal{D})$ can't be computed exactly (Solution: Need to approximate the CP)
2. Integral for the expectation is intractable (Solution: Use Monte Carlo approximation)

- Draw M i.i.d. samples of \mathbf{Z} from the current (exact/approximate) CP $p(\mathbf{Z}|\Theta^{(t-1)}, \mathcal{D})$

$$\{\mathbf{Z}^{(i)}\}_{i=1}^M \sim p(\mathbf{Z}|\Theta^{(t-1)}, \mathcal{D})$$

- Use these samples to get a Monte-Carlo approximation of expected CLL and maximize

$$\Theta^{(t)} = \operatorname{argmax}_{\Theta} \frac{1}{M} \sum_{i=1}^M \log p(\mathcal{D}, \mathbf{Z}^{(i)}|\Theta)$$

- Monte-Carlo approximation is commonly used in such problems

EM: Some Final Comments

- The E and M steps may not always be possible to perform exactly. Some reasons

- The conditional posterior of latent variables $p(\mathbf{Z}|\mathbf{X}, \Theta)$ may not be easy to compute
 - Will need to approximate $p(\mathbf{Z}|\mathbf{X}, \Theta)$ using methods such as MCMC or variational inference

- Even if $p(\mathbf{Z}|\mathbf{X}, \Theta)$ is easy, the expected CLL may not be easy to compute

$$\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] = \int \log p(\mathbf{X}, \mathbf{Z}|\Theta) p(\mathbf{Z}|\mathbf{X}, \Theta) d\mathbf{Z}$$

Can often be approximated by Monte-Carlo using sample from the CP of \mathbf{Z}

Results in Monte-Carlo EM

- Maximization of the expected CLL may not be possible in closed form
- EM works even if the M step is only solved approximately (Generalized EM)
- If M step has multiple parameters whose updates depend on each other, they are updated in an alternating fashion - called Expectation Conditional Maximization (ECM)
- Other advanced probabilistic inference algos are based on ideas similar to EM
 - E.g., Variational EM, Variational Bayes (VB) inference, a.k.a. Variational Inference (VI)