# Latent Variable Models and EM Algorithm

CS772A: Probabilistic Machine Learning

Piyush Rai

# Conditional Posterior

- Consider a model with $K$ unknown params/hyperparams $\Theta = (\theta_1, \theta_2, \ldots, \theta_K)$

Joint posterior

$$p(\Theta|\boldsymbol{X}) = \frac{p(\Theta)p(\boldsymbol{X}|\Theta)}{p(\boldsymbol{X})} = \frac{p(\Theta)p(\boldsymbol{X}|\Theta)}{\int p(\Theta)p(\boldsymbol{X}|\Theta)d\theta_1 d\theta_2 \ldots d\theta_K}$$

Usually intractable integral so the full posterior can't be computed exactly

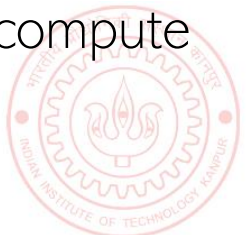- We can however compute conditional posteriors (CP) which for each $\theta_i$ looks like

$$p(\theta_i | \text{whatever } \theta_i \text{ depends on})$$

Can be data and/or other params/hyperparams given their fixed values (or current estimates)

- To compute each CP, look at the joint distribution $p(\boldsymbol{X}, \Theta)$

$$p(\boldsymbol{X}, \Theta) = p(\boldsymbol{X}, \theta_1, \theta_2, \ldots, \theta_K) = p(\boldsymbol{X}|\theta_1, \theta_2, \ldots, \theta_K)p(\theta_1|\theta_2, \ldots, \theta_K)p(\theta_2|\theta_3, \ldots, \theta_K) \ldots p(\theta_K)$$

- CP of $\theta_i$ will be proportional to the product of all the terms involving $\theta_i$
  - If those terms are conjugate to each other, it is called local conjugacy. CP is then easy to compute

- Many algorithms for computing point estimate/full posterior use the CPs
  - Expectation Maximization, Variational Inference , MCMC (especially Gibbs sampling)
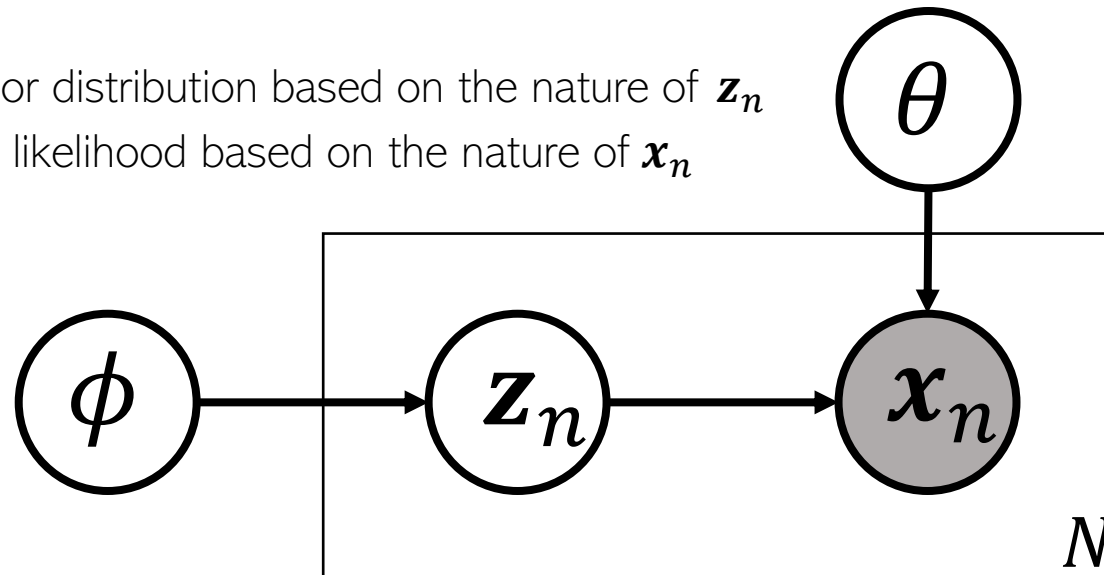
# Latent Variable Models

- Application 1: Can use latent variables to learn latent properties/features of data, e.g.,
  - Cluster assignment of each observation (in mixture models)
  - Low-dim rep. or "code" of each observation (e.g., prob. PCA, variational autoencoders, etc)

Plate notation of a generic LVM

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

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



- In such apps, latent variables ($\boldsymbol{z_n}$'s) are called "local variables" (specific to individual obs.) and other unknown parameters/hyperparams ($\boldsymbol{\theta}, \boldsymbol{\phi}$ above) are called "global var"

# Latent Variable Models

- Application 2: Sometimes, <u>augmenting</u> a model by latent variables simplifies inference
  - These latent variables aren't part of the original model definition
- Some of the popular examples of such augmentation include
  - In Probit regression for binary classification, we can model each label $y_n \in \{0,1\}$ as

$$y_n = \mathbb{I}[z_n > 0] \qquad \text{where} \qquad z_n \sim \mathcal{N}(\boldsymbol{w}^\top \boldsymbol{x}_n, 1) \quad \text{is an auxiliary latent variable}$$

  .. and use EM etc, to infer the unknowns $\boldsymbol{w}$ and $z_n$'s (PML-2, Sec 15.4)

  - Many sparse priors on weights can be thought of as Gaussian "scale-mixtures"

$$\text{Laplace}(w_d|0, 1/\gamma) = \frac{\gamma}{2}\exp(-\gamma|w_d|) = \int \mathcal{N}(w_d|0, \tau_d^2)\text{Gamma}(\tau_d^2|1, \gamma^2/2)d\tau_d^2$$

  .. where $\tau_d$'s are latent vars. Can use EM to infer $\boldsymbol{w}$, $\boldsymbol{\tau}$ (MLAPP 13.4.4 - EM for LASSO)

- Such augmentations can often make a non-conjugate model a locally conjugate one
  - Conditional posteriors of the unknowns often have closed form in such cases

# 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

# 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

  CP of $w$: $p(\boldsymbol{w}|\boldsymbol{X}, \boldsymbol{y}, \hat{\lambda}, \hat{\beta})$

  - Inferring CP over the main parameter given the point estimates of hyperparams
  - Maximizing the marginal lik. to do point estimation for hyperparams

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

- The Expectation-Maximization algorithm (will see today) also does something similar
  - In E step, the CP of latent variables is inferred, given <u>current</u> point-est of params
  - M step maximizes <span style="color:red">expected complete data log-lik</span>. 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 <span style="color:blue">posterior over local vars</span> and <span style="color:blue">point estimates for global vars</span>
  - 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

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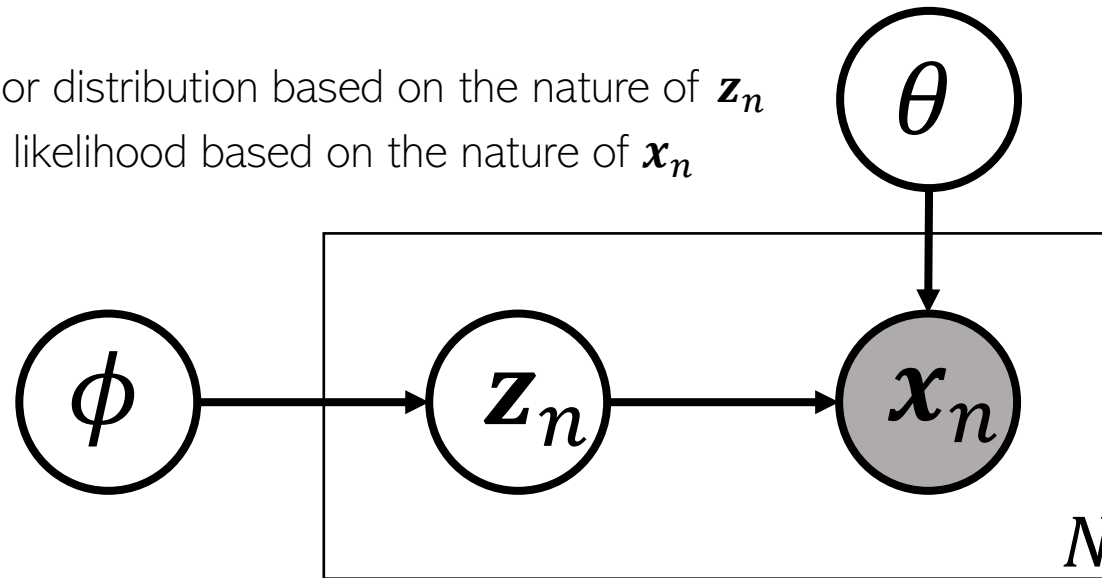
# Parameter Estimation in Latent Variable Models

■ Assume each observation $\boldsymbol{x}_n$ to be associated with a "local" latent variable $\boldsymbol{z}_n$

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

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



■ Although we can do fully Bayesian inference for all the unknowns, suppose we only want a point estimate of the "global" parameters $\Theta = (\theta, \boldsymbol{\phi})$ 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 $\Theta$

■ However, EM gives a clean way to obtain closed form updates for $\Theta$

# 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 $\Theta$ 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

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

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

- Can we maximize some other quantity instead of $\log p(x_n | \Theta)$ for this MLE?

# 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

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

Verify the identity

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

Assume **Z** discrete

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

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



$KL(q||p)$

$\mathcal{L}(q, \boldsymbol{\theta})$

$\ln p(\mathbf{X}|\boldsymbol{\theta})$

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

Basically, log of marginal likelihood w.r.t. $\Theta$ with $Z$ integrated out

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

- $\mathcal{L}(q, \Theta)$ depends on $q$ and $\Theta$. We'll use ALT-OPT to maximize it

- Let's maximize $\mathcal{L}(q, \Theta)$ w.r.t. $q$ with $\Theta$ fixed at some $\Theta^{\text{old}}$

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

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

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

The posterior distribution of $Z$ given current parameters $\Theta^{\text{old}}$

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

Maximization of expected CLL where the expectation is w.r.t. the posterior distribution of $Z$ given current parameters $\Theta^{\text{old}}$

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

Complete-Data Log Likelihood (CLL)

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

Much easier than maximizing ILL since CLL will have simple expressions (since it is akin to knowing $Z$)

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

# The Expectation-Maximization (EM) Algorithm

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

### The EM Algorithm

Primarily designed for doing point estimation of the parameters $\Theta$ 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. $\Theta$ as the M step

① Initialize $\Theta$ 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

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

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

③ Step 2: Now maximize the expected complete data log-likelihood w.r.t. $\Theta$

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

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 $\arg\max_{\Theta} \sum_{n=1}^{N}\left[\log p\left(x_n, \hat{z}_n^{(t)}|\Theta\right)\right]$ this will be ALT-OPT

# The Expected CLL

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

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

- If $p(\boldsymbol{z}_n|\Theta)$ and $p(\boldsymbol{x}_n|\boldsymbol{z}_n, \Theta)$ are exp-family distributions, $\mathcal{Q}(\Theta, \Theta^{\mathbf{old}})$ has a very simple form

- In resulting expressions, replace terms containing $\boldsymbol{z}_n$'s by their respective expectations, e.g.,
  - $\boldsymbol{z}_n$ replaced by $\mathbb{E}_{p(\boldsymbol{z}_n|\boldsymbol{x}_n,\widehat{\Theta})}[\boldsymbol{z}_n]$
  - $\boldsymbol{z}_n\boldsymbol{z}_n^{\top}$ replaced by $\mathbb{E}_{p(\boldsymbol{z}_n|\boldsymbol{x}_n,\widehat{\Theta})}[\boldsymbol{z}_n\boldsymbol{z}_n^{\top}]$

- 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 $\Theta$

- Step 2 maximizes $\mathcal{L}(\hat{q}, \Theta)$ w.r.t. $\Theta$ which gives a new $\Theta$.

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 $\Theta$

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

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

Local optima found for $\Theta_{MLE}$
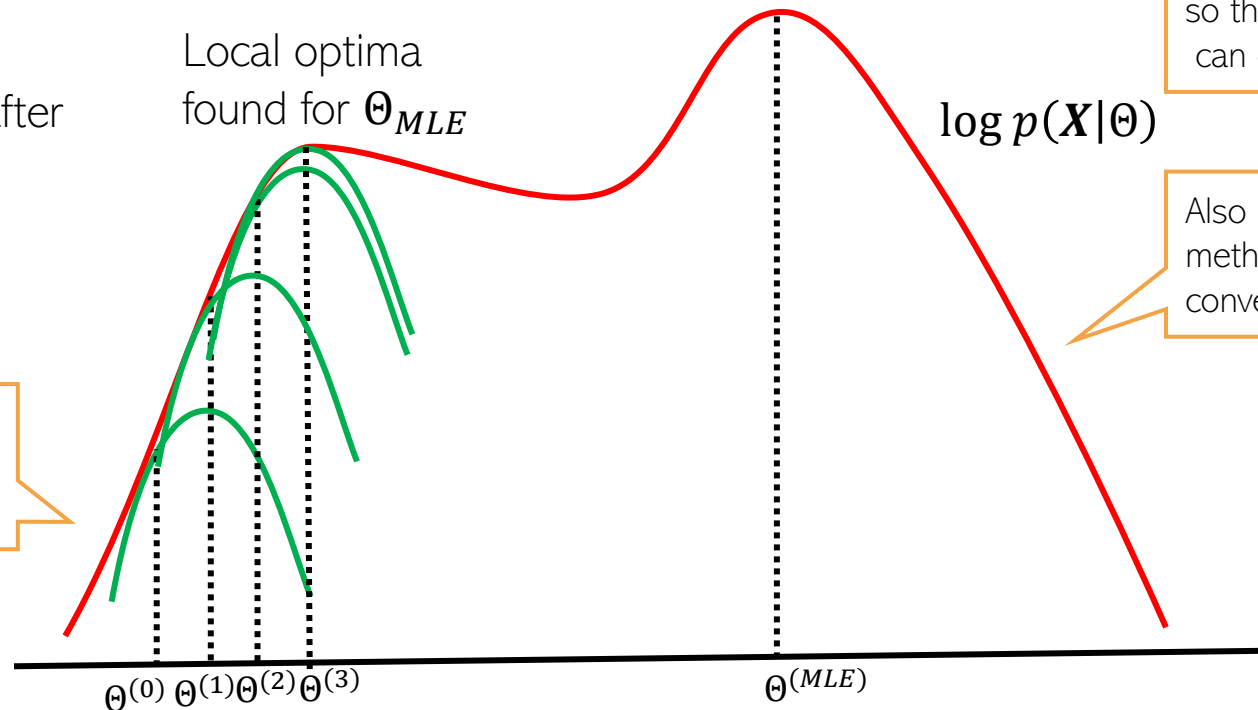
$\log p(\mathbf{X}|\Theta)$

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

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

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

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

$\Theta^{(0)}$ $\Theta^{(1)}$ $\Theta^{(2)}$ $\Theta^{(3)}$          $\Theta^{(MLE)}$

# 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 $\Theta$

- EM is usually preferred over this approach because
  - The M step has often simple closed-form updates for the parameters $\Theta$
  - 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 Z (from E step)

# 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)
    - $\widehat{\Theta} = \text{argmax}_\Theta \log p(x^{obs}|\Theta) = \text{argmax}_\Theta \log \int p([x^{obs}, x^{miss}]|\Theta) dx^{miss}$
    - $\widehat{\Theta} = \text{argmax}_\Theta \sum_{n=1}^{N} \log p(x_n, y_n|\Theta) + \sum_{n=N+1}^{N+M} \log \sum_{c=1}^{K} p(x_n, y_n = c|\Theta)$

- Hyperparameter estimation in probabilistic models (an alternative to MLE-II)
    - MLE-II estimates hyperparams by maximizing the marginal likelihood, e.g.,

$$\{\hat{\lambda}, \hat{\beta}\} = \text{argmax}_{\lambda,\beta} \, p(y|X, \lambda, \beta) = \text{argmax}_{\lambda,\beta} \int p(y|w, X, \beta) p(w|\lambda) dw$$

> For a Bayesian linear regression model

    - With EM, can treat $w$ as latent var, and $\lambda, \beta$ as "parameters"
        - E step will estimate the CP of $w$ given current estimates of $\lambda, \beta$
        - M step will re-estimate $\lambda, \beta$ by maximizing the expected CLL

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

> Expectations w.r.t. the CP of $w$

# An Example: Mixture Models

■ Assume $K$ probability distributions (e.g., Gaussians), one for each cluster

$p(z_n \mid \phi) = \text{multinoulli}(\boldsymbol{\pi})$
(also means $p(z_n = k \mid \phi) = \pi_k$)

Parameters of the $K$ distributions, e.g., $\theta = \{\mu_k, \Sigma_k\}_{k=1}^{K}$

$p(x)$ is a Gaussian mixture model (GMM)

Discrete latent variable (with $K$ possible values) or a one-hot vector of length $K$. Modeled by a multinoulli distribution as prior
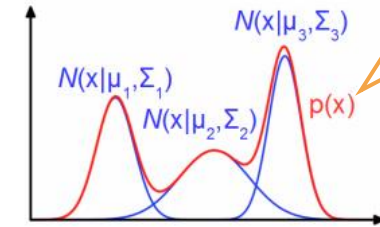
Assumed generated from one of the $K$ distributions depending on the true (but unknown) value of $z_n$ (which clustering will find))

The likelihood distributions

The parameter vector $\boldsymbol{\pi} = [\pi_1, \pi_2, \dots, \pi_K]$ of the multinoulli distribution



$$p(x_n \mid z_n = k, \theta) = \mathcal{N}(\mu_k, \Sigma_k)$$

■ The log-likelihood will be

MLE on this objective won't give closed form solution for the parameters

$$\log p(\boldsymbol{x}_n \mid \Theta) = \log \sum_{k=1}^{K} p(\boldsymbol{x}_n, \boldsymbol{z}_n = k \mid \Theta)$$

$$= \log \sum_{k=1}^{K} p(\boldsymbol{z}_n = k \mid \phi) p(\boldsymbol{x}_n \mid \boldsymbol{z}_n = k, \theta) = \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n \mid \mu_k, \Sigma_k)$$

# Detour: MLE for Generative Classification

- Assume a $K$ class generative classification model with Gaussian class-conditionals
- Assume class $k = 1, 2, \ldots, K$ is modeled by a Gaussian with mean $\mu_k$ and cov matrix $\Sigma_k$
- The labels $\mathbf{z}_n$ (known) are one-hot vecs. Also, $z_{nk} = 1$ if $\mathbf{z}_n = k$, and $\mathbf{z}_{nk} = 0$, o/w
- Assuming class prior as $p(\mathbf{z}_n = k) = \pi_k$, the model has params $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$
- Given training data $\{\boldsymbol{x}_n, \boldsymbol{z}_n\}_{n=1}^{N}$, the MLE solution will be

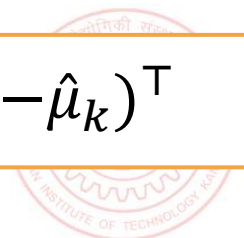$$\hat{\pi}_k = \frac{1}{N} \sum_{n=1}^{N} z_{nk}$$

Same as $\frac{N_k}{N}$ where $N_k$ is # of training ex. for which $y_n = k$

$$\hat{\mu}_k = \frac{1}{N_k} \sum_{n=1}^{N} z_{nk} \boldsymbol{x}_n$$

Same as $\frac{1}{N_k} \sum_{n:\boldsymbol{z}_n=k}^{N} \boldsymbol{x}_n$

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^{N} z_{nk} (\boldsymbol{x}_n - \hat{\mu}_k)(\boldsymbol{x}_n - \hat{\mu}_k)^{\mathsf{T}}$$

Same as $\frac{1}{N_k} \sum_{n:\boldsymbol{z}_n=k}^{N} (\boldsymbol{x}_n - \hat{\mu}_k)(\boldsymbol{x}_n - \hat{\mu}_k)^{\mathsf{T}}$

# Detour: MLE for Generative Classification

- Here is a formal derivation of the MLE solution for $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$

$$\hat{\Theta} = \text{argmax}_{\Theta}\ p(X, Z|\Theta) = \text{argmax}_{\Theta}\ \prod_{n=1}^{N} p(x_n, z_n|\Theta)$$

multinoulli

Gaussian

$$= \text{argmax}_{\Theta}\ \prod_{n=1}^{N} p(z_n|\Theta)\ p(x_n|z_n, \Theta)$$

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

$$= \text{argmax}_{\Theta}\ \prod_{n=1}^{N} \prod_{k=1}^{K} \pi_k^{z_{nk}} \prod_{k=1}^{K} p(x_n|z_n = k, \Theta)^{z_{nk}}$$

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

$$= \text{argmax}_{\Theta}\ \prod_{n=1}^{N} \prod_{k=1}^{K} [\pi_k p(x_n|z_n = k, \Theta)]^{z_{nk}}$$

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

$$= \text{argmax}_{\Theta}\ \log \prod_{n=1}^{N} \prod_{k=1}^{K} [\pi_k p(x_n|z_n = k, \Theta)]^{z_{nk}}$$

$$= \text{argmax}_{\Theta}\ \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk}[\log \pi_k + \log \mathcal{N}(x_n|\mu_k, \Sigma_k)]$$

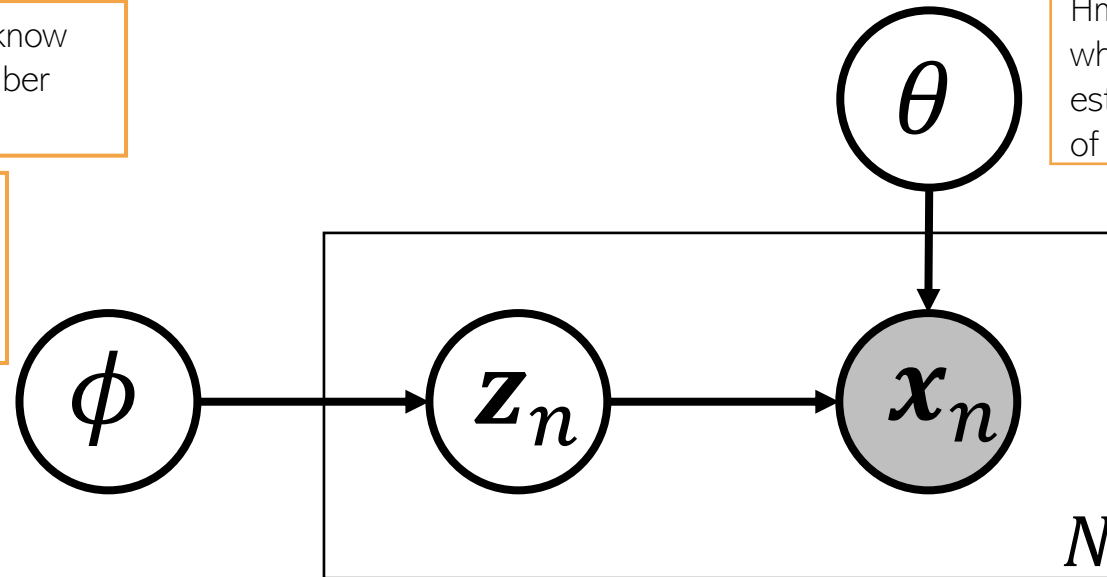The form of this expression is important; will encounter this in GMM too

# EM for Mixture Models

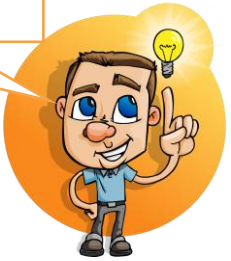- So how do we estimate the parameters of a GMM where $z_n$'s are <u>unknown</u>?

Well, you kind of already know how to do this. ☺ Remember generative classification?

Hmmmm.. So can we make a guess what the value of each $z_n$ and then estimate $\theta$ and $\phi$ as we do in case of generative classification??

Yes, exactly. ☺ However, just like in gen-class, you will need to repeat the guess and estimate them a few times until you converge

$\theta$

$\phi$ → $z_n$ → $x_n$

$N$

- The guess about $z_n$ can be in one of the two forms
  - A "hard" guess – a single best value $\hat{z}_n$ (some "optimal" value of the random variable $z_n$)
  - The "expected" value $\mathbb{E}[z_n]$ of the random variable $z_n$

EM is pretty much like ALT-OPT but with soft/expected values of the latent variables

- Using the hard guess $\hat{z}_n$ of $z_n$ will result in an ALT-OPT like algorithm
- Using the expected value of $z_n$ will give the so-called Expectation-Maximization (EM) algo

# EM for Gaussian Mixture Model (GMM)

- EM finds $\Theta_{MLE}$ by maximizing $\mathbb{E}[\log p(\boldsymbol{X}, \boldsymbol{Z}|\Theta)]$ rather than $\log p(\boldsymbol{X}, \widehat{\boldsymbol{Z}}|\Theta)$

Expectation of CLL

- Note: Expectation will be w.r.t. the <u>conditional</u> posterior distribution of $\boldsymbol{Z}$, i.e., $p(\boldsymbol{Z}|\boldsymbol{X}, \Theta)$

It is "conditional" posterior because it is also conditioned on $\Theta$, not just data $X$

Requires knowing $\Theta$

- The EM algorithm for GMM operates as follows
    - Initialize $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ as $\widehat{\Theta}$
    - Repeat until convergence

    Needed to get the expected CLL

        - Compute conditional posterior $p(\boldsymbol{Z}|\boldsymbol{X}, \widehat{\Theta})$. Since obs are i.i.d, compute separately for each $\boldsymbol{n}$ (and for $k = 1, 2, ..K$)

Same as $p(z_{nk} = 1|\boldsymbol{x}_n, \widehat{\Theta})$, just a different notation

$$p(\boldsymbol{z}_n = k|\boldsymbol{x}_n, \widehat{\Theta}) \propto p(\boldsymbol{z}_n = k|\widehat{\Theta}) p(\boldsymbol{x}_n|\boldsymbol{z}_n = k, \widehat{\Theta}) = \hat{\pi}_k \mathcal{N}(x_n|\hat{\mu}_k, \hat{\Sigma}_k)$$

        - Update $\Theta$ by maximizing the <u>expected</u> complete data log-likelihood
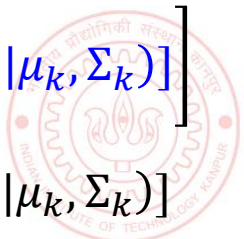
Solution has a similar form as ALT-OPT (or gen. class.), except we now have the **expectation** of $z_{nk}$ being used

$$\widehat{\Theta} = \text{argmax}_\Theta \mathbb{E}_{p(\boldsymbol{Z}|\boldsymbol{X},\widehat{\Theta})}[\log p(\boldsymbol{X}, \boldsymbol{Z}|\Theta)] = \sum_{n=1}^N \mathbb{E}_{p(\boldsymbol{z}_n|\boldsymbol{x}_n,\widehat{\Theta})}[\log p(\boldsymbol{x}_n, \boldsymbol{z}_n|\Theta)]$$

$$\hat{\pi}_k = \frac{1}{N}\sum_{n=1}^N \mathbb{E}[z_{nk}] \quad \hat{\mu}_k = \frac{1}{N_k}\sum_{n=1}^N \mathbb{E}[z_{nk}]\boldsymbol{x}_n$$

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

$N_k$ : Effective number of points in cluster k

$$\hat{\Sigma}_k = \frac{1}{N_k}\sum_{n=1}^N \mathbb{E}[z_{nk}](\boldsymbol{x}_n - \hat{\mu}_k)(\boldsymbol{x}_n - \hat{\mu}_k)^\top$$

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

# EM for GMM (Contd)

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

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

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

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

> Accounts for fraction of points in each cluster

> 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 cluster shapes (since each cluster is a Gaussian

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

❸ Given "responsibilities" $\gamma_{nk} = \mathbb{E}[z_{nk}]$, and $N_k = \sum_{n=1}^{N} \gamma_{nk}$, re-estimate $\Theta$ via MLE

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

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

M-step:

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

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

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

# 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 $p(Z|X, \Theta)$ may not be easy to compute
    - Will need to approximate $p(Z|X, \Theta)$ using methods such as MCMC or variational inference

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

  Results in
  Monte-Carlo EM

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

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

  - 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)