

# Latent Variable Models and the EM Algorithm (Contd)

CS772A: Probabilistic Machine Learning

Piyush Rai

# Announcement

- Mid-sem exam on Sept 19, 1800-2000
- Venue: L17, ERES seating scheme
- Syllabus: Up to today's lecture
- Closed book exam
  - Necessary formulae/equations will be provided in the question paper



# Plan Today

- Expectation Maximization (EM) algorithm for param-est/inference in LVMs
- An example of EM
  - Parameter estimation for Gaussian mixture models

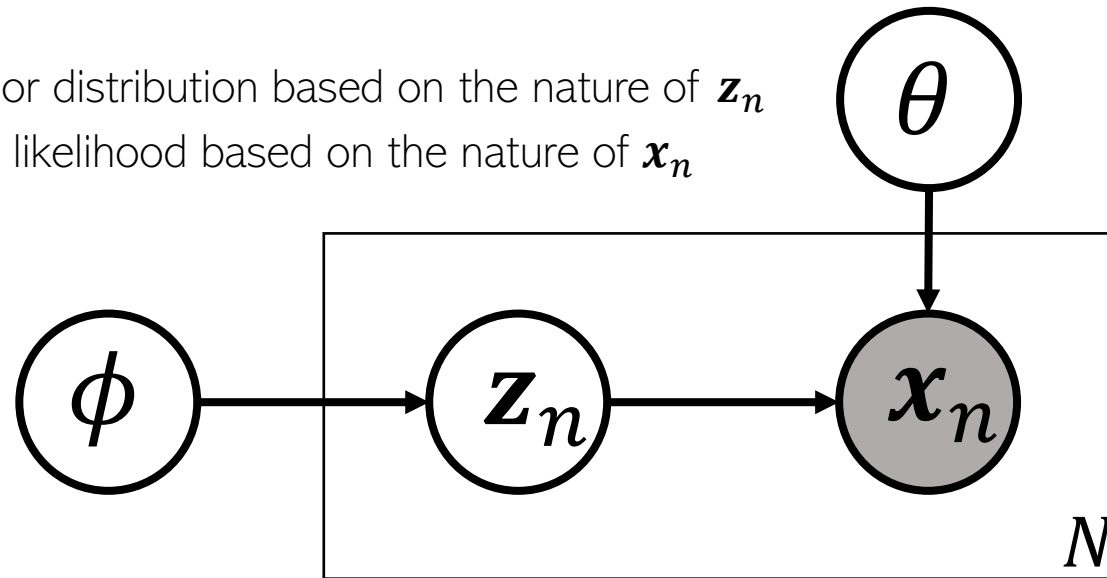


# Parameter Estimation in Latent Variable Models

- 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
  - Can do **gradient based opt on log-lik** but usually won’t get closed form updates for  $\Theta$
  - However, **EM** algo gives a clean way to obtain closed form updates for  $\Theta$  for such LVMs



# Why MLE/MAP of Params is Hard for LVMs?

- Suppose we want to estimate parameters  $\Theta$  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

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

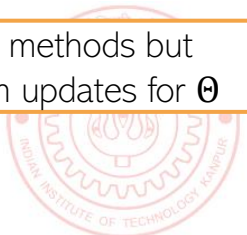
- $\log p(\mathbf{x}_n | \Theta)$  will not have a simple expression since  $p(\mathbf{x}_n | \Theta)$  requires sum/integral

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

- MLE now becomes difficult, no closed form expression for  $\Theta$

Can still use gradient based methods but won't get clean, closed-form updates for  $\Theta$

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

Assume  $\mathbf{z}$  discrete

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

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

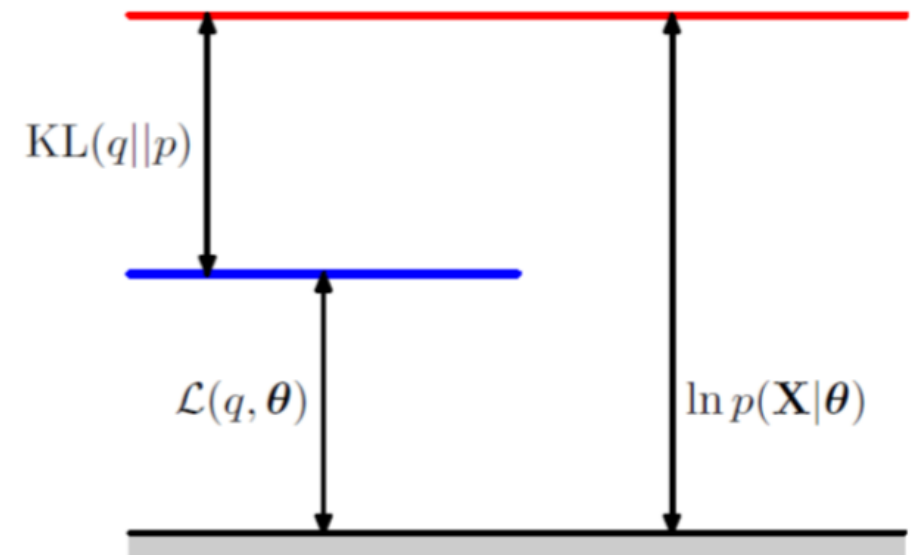
- $KL(q||p_z) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \Theta)}{q(\mathbf{Z})} \right\}$

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



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

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

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

- Now let's maximize  $\mathcal{L}(q, \Theta)$  w.r.t.  $\Theta$  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  $\Theta^{\text{old}}$

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

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

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

**Complete-Data Log Likelihood (CLL)**

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

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

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



# 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

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

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

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

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

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 the CP) in Step 1 and maximize the CLL in Step 2 using that, i.e., do  $\arg \max_{\Theta} \sum_{n=1}^N [\log p(\mathbf{x}_n, \hat{z}_n^{(t)} | \Theta)]$  then this will be std ALT-OPT





# The Expected CLL

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

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

- If  $p(\mathbf{z}_n|\Theta)$  and  $p(\mathbf{x}_n|\mathbf{z}_n, \Theta)$  are exp-family distributions,  $Q(\Theta, \Theta^{old})$  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  $\mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \hat{\Theta})}[\mathbf{z}_n]$
  - $\mathbf{z}_n\mathbf{z}_n^T$  replaced by  $\mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \hat{\Theta})}[\mathbf{z}_n\mathbf{z}_n^T]$
- However, in some LVMs, these expectations are intractable to compute and need to be approximated (will see some examples later), e.g., using MC integral approximation



# What's Going On In EM?

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  $\Theta$

- 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  $\Theta$ .

Note that  $\Theta$  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

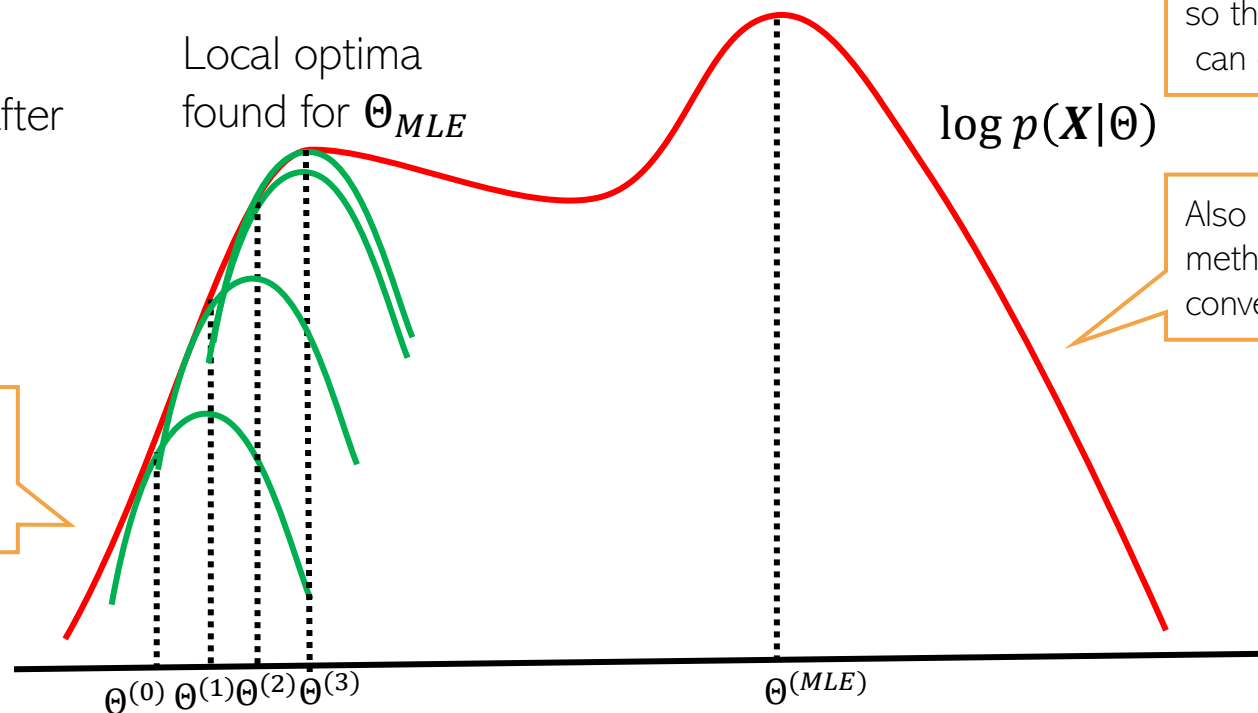
Even though original MLE problem  $\operatorname{argmax}_{\Theta} \log p(\mathbf{X}|\Theta)$  could be solved using gradient methods, EM often works faster and has cleaner updates

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

Local optima found for  $\Theta_{MLE}$

$\log p(\mathbf{X}|\Theta)$

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



# Some Applications of EM

- Mixture Models (each data-point comes from one of  $K$  mixture components)
  - Examples: Mixture of Gaussians, Mixture of Experts, etc
- Latent variable models for dimensionality reduction or representation learning
  - Probabilistic PCA, Factor Analysis, Variational Autoencoders, etc
- Problems models with missing features/labels (treated as latent variables)
  - An example of problem with missing labels: [Semi-supervised learning](#)
- [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

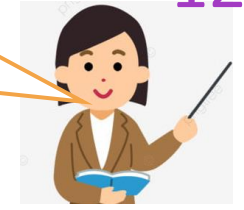
- With EM, can treat  $\mathbf{w}$  as latent var, and  $\lambda, \beta$  as “parameters”
  - E step will estimate the CP of  $\mathbf{w}$  given current estimates of  $\lambda, \beta$
  - M step will re-estimate  $\lambda, \beta$  by maximizing the expected CLL

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

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



# An Example: Mixture Models



If the  $\mathbf{z}_n$  were known, it just becomes **generative classification**, for which which we know how to estimate  $\theta$  and  $\phi$ , given training data

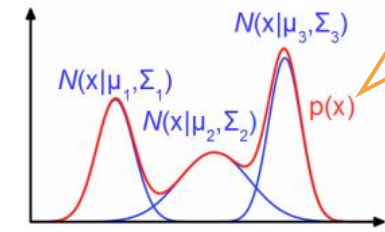
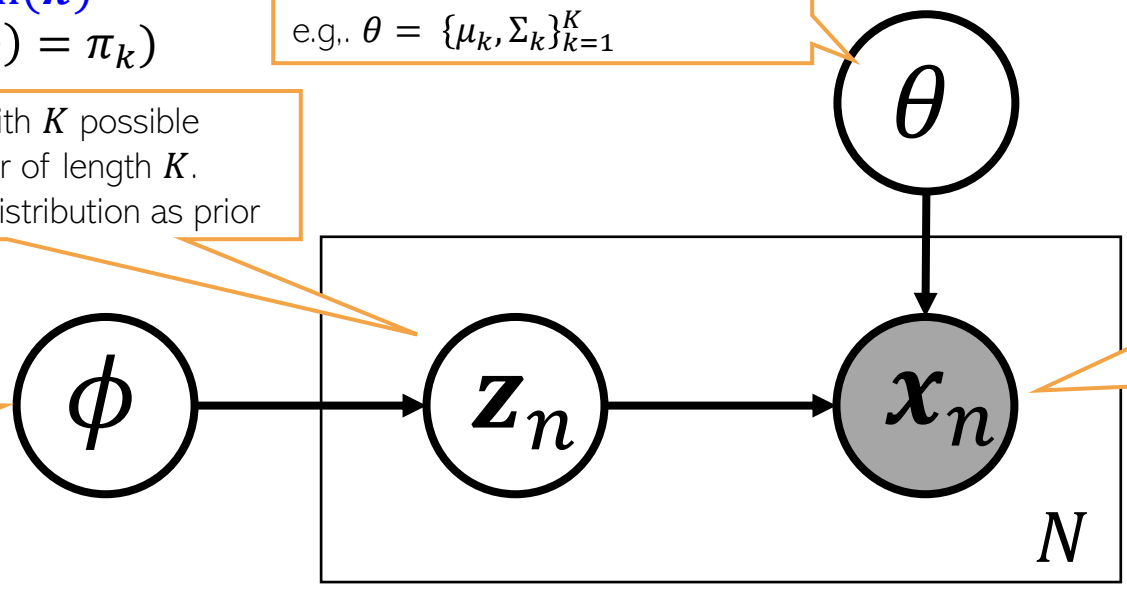
- Assume  $K$  probability distributions (e.g., Gaussians), one for each cluster

$p(\mathbf{z}_n | \phi) = \text{multinoulli}(\boldsymbol{\pi})$   
 (also means  $p(\mathbf{z}_n = k | \phi) = \pi_k$ )

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

Discrete latent variable (with  $K$  possible values) or a one-hot vector of length  $K$ . Modeled by a multinoulli distribution as prior

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



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

Assumed generated from one of the  $K$  distributions depending on the true (but unknown) value of  $\mathbf{z}_n$  (which clustering will find))

The likelihood distributions

$p(\mathbf{x}_n | \mathbf{z}_n = k, \theta) = \mathcal{N}(\mu_k, \Sigma_k)$

- The log-likelihood will be

$$\begin{aligned} \log p(\mathbf{x}_n | \Theta) &= \log \sum_{k=1}^K p(\mathbf{x}_n, \mathbf{z}_n = k | \Theta) \\ &= \log \sum_{k=1}^K p(\mathbf{z}_n = k | \phi) p(\mathbf{x}_n | \mathbf{z}_n = k, \theta) = \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \end{aligned}$$

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

# Detour: MLE for Generative Classification

- Assume a  $K$  class generative classification model with Gaussian class-conditionals
- Assume class  $k = 1, 2, \dots, K$  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  $z_{nk} = 0$ , 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  $\{\mathbf{x}_n, \mathbf{z}_n\}_{n=1}^N$ , the MLE solution will be

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

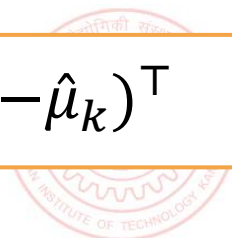
Same as  $\frac{N_k}{N}$  where  $N_k$  is # of training ex. for which  $y_n = k$

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

Same as  $\frac{1}{N_k} \sum_{n:z_n=k} \mathbf{x}_n$

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

Same as  $\frac{1}{N_k} \sum_{n:z_n=k} (\mathbf{x}_n - \hat{\mu}_k)(\mathbf{x}_n - \hat{\mu}_k)^\top$



# Detour: MLE for Generative Classification

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

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

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

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



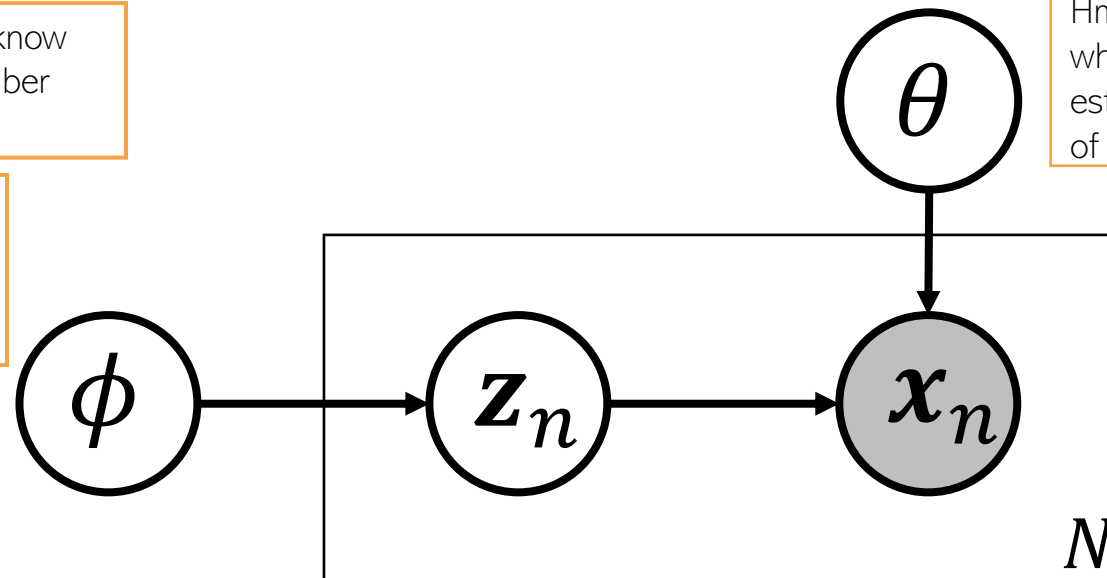
# EM for Mixture Models

- So how do we estimate the parameters of a GMM where  $\mathbf{z}_n$ 's are unknown?

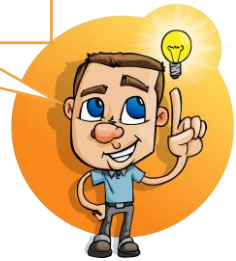


Well, you kind of already know how to do this. 😊 Remember generative classification?

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



Hmmm.. So can we make a guess what the value of each  $\mathbf{z}_n$  and then estimate  $\theta$  and  $\phi$  as we do in case of generative classification??



- The guess about  $\mathbf{z}_n$  can be in one of the two forms
  - A “hard” guess – a single best value  $\hat{\mathbf{z}}_n$  (some “optimal” value of the random variable  $\mathbf{z}_n$ )
  - The “expected” value  $\mathbb{E}[\mathbf{z}_n]$  of the random variable  $\mathbf{z}_n$
- Using the hard guess  $\hat{\mathbf{z}}_n$  of  $\mathbf{z}_n$  will result in an ALT-OPT like algorithm
- Using the expected value of  $\mathbf{z}_n$  will give the so-called Expectation-Maximization (EM) algo

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

# EM for Gaussian Mixture Model (GMM)

- EM finds  $\Theta_{MLE}$  by maximizing  $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$  rather than  $\log p(\mathbf{X}, \hat{\mathbf{Z}}|\Theta)$ 
  - Note: Expectation will be w.r.t. the conditional posterior distribution of  $\mathbf{Z}$ , i.e.,  $p(\mathbf{Z}|\mathbf{X}, \Theta)$ 
    - Expectation of CLL
    - It is "conditional" posterior because it is also conditioned on  $\Theta$ , not just data  $\mathbf{X}$
    - Requires knowing  $\Theta$
  - The EM algorithm for GMM operates as follows
    - Initialize  $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  as  $\hat{\Theta}$
    - Repeat until convergence
      - Compute conditional posterior  $p(\mathbf{Z}|\mathbf{X}, \hat{\Theta})$ . Since obs are i.i.d, compute separately for each  $n$  (and for  $k = 1, 2, \dots, K$ )
        - Needed to get the expected CLL

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

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

- Update  $\Theta$  by maximizing the expected complete data log-likelihood

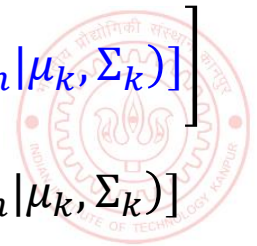
Solution has a similar form as ALT-OPT (or gen. class.), except we now have the expectation of  $z_{nk}$  being used

$N_k$  : Effective number of points in cluster  $k$

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

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

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





# EM for GMM (Contd)

- 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)

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)

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

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

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

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

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



# 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<sup>†</sup>, 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 also provides the conditional posterior over the latent variables  $Z$  (from E step)

<sup>†</sup>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)



# 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 \(VI\)](#)

