

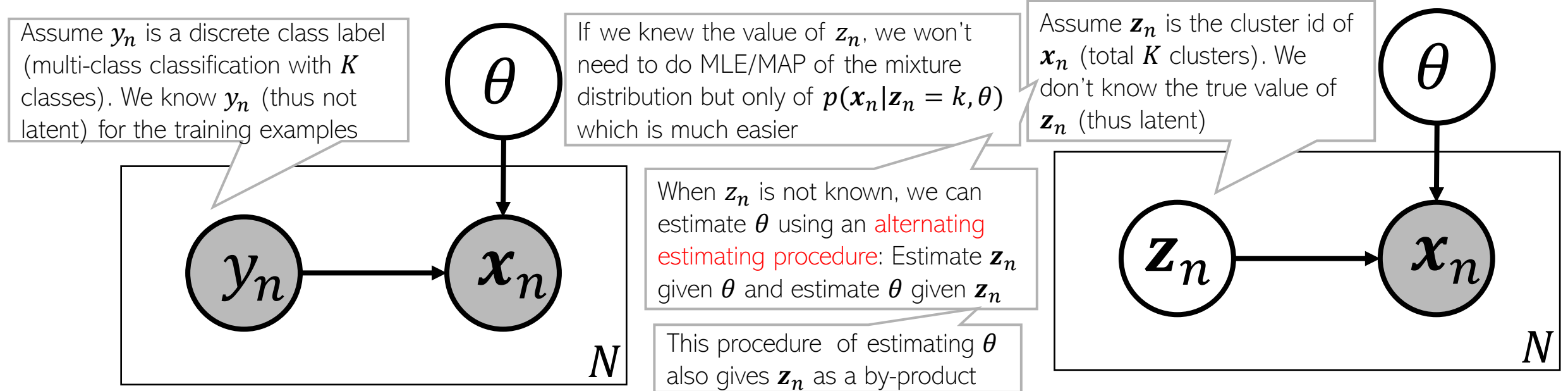
# Latent Variable Models (LVMs)

CS771: Introduction to Machine Learning

Piyush Rai

# Example: Generative Models with Latent Variables<sup>2</sup>

- Two generative models of inputs  $\mathbf{x}_n$  without (left) and with (right) latent variables



- Suppose we wish to estimate (e.g., using MLE/MAP) params  $\theta$  of distribution of  $\mathbf{x}_n$
- For case 1, the distribution is  $p(\mathbf{x}_n | y_n, \theta)$  and MLE/MAP of  $\theta$  easy since  $y_n$  is known
- For case 2, distribution is more complex because true  $z_n$  is not known

$$p(\mathbf{x}_n | \theta) = \sum_{k=1}^K p(\mathbf{x}_n, z_n = k | \theta) = \sum_{k=1}^K p(z_n = k) p(\mathbf{x}_n | z_n = k, \theta)$$

Reason: The functional form of mixture can be messy

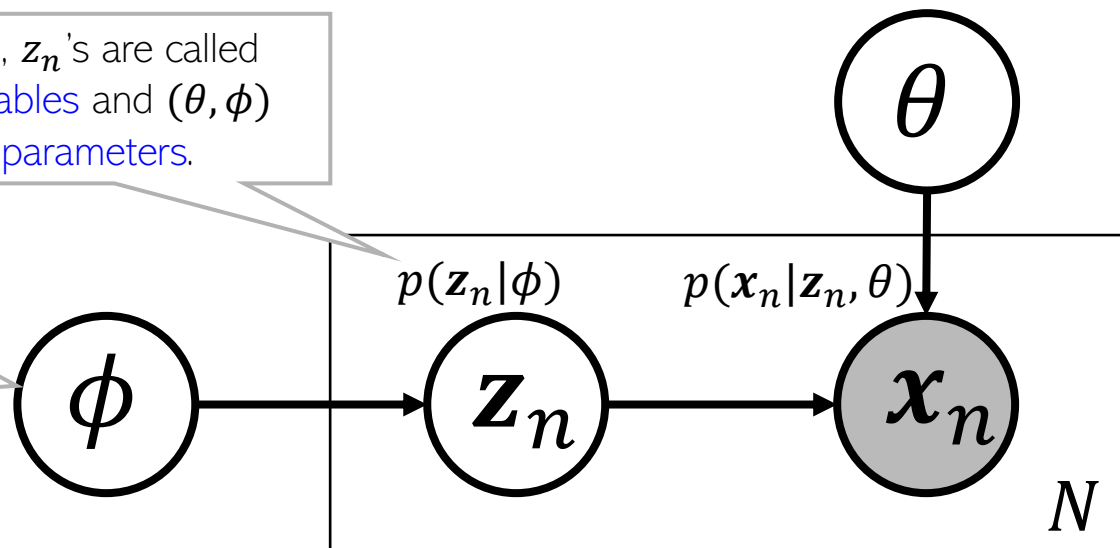
MLE/MAP a bit difficult for this more complex "mixture" of distributions

# Components of an LVM

- Recall that the goal is to estimate  $\theta$  (and  $\mathbf{z}_n$  is also unknown)
- In LVM, we treat  $\mathbf{z}_n$  as a random variable and assume a prior distribution  $p(\mathbf{z}_n|\phi)$

In an LVM,  $\mathbf{z}_n$ 's are called **latent variables** and  $(\theta, \phi)$  are called **parameters**.

Will also need to estimate  $\phi$  in addition to  $\theta$



This prior tells us what the value of  $\mathbf{z}_n$  is before we have seen the input  $\mathbf{x}_n$

Ultimately, we will compute the distribution of  $\mathbf{z}_n$  conditioned on the input  $\mathbf{x}_n$

For example, a  $D$ -dimensional Gaussian if  $\mathbf{x}_n \in \mathbb{R}^D$

- We will also assume a suitable conditional distribution  $p(\mathbf{x}_n|\mathbf{z}_n, \theta)$  for  $\mathbf{x}_n$
- The form of  $p(\mathbf{z}_n|\phi)$  will depend on the nature of  $\mathbf{z}_n$ , e.g.,
  - If  $\mathbf{z}_n$  is discrete with  $K$  possible values,  $p(\mathbf{z}_n|\phi) = \text{multinoulli}(\mathbf{z}_n|\boldsymbol{\pi})$
  - If  $\mathbf{z}_n \in \mathbb{R}^K$ ,  $p(\mathbf{z}_n|\phi) = \mathcal{N}(\mathbf{z}_n|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , a  $K$ -dim Gaussian



# Why Direct MLE/MAP is Hard for LVMs?

- Direct MLE/MAP of parameters  $(\theta, \phi) = \Theta$  without estimating  $\mathbf{z}_n$  is hard
- Reason: Given  $N$  observations  $\mathbf{x}_n, n = 1, 2, \dots, N$ , the MLE problem for  $\Theta$  will be

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

Also note that  $p(\mathbf{x}_n, \mathbf{z}_n | \Theta) = p(\mathbf{z}_n | \phi) p(\mathbf{x}_n | \mathbf{z}_n, \theta)$

Summing over all possible values  $\mathbf{z}_n$  can take (would be an integral instead of sum if  $\mathbf{z}_n$  is continuous)

Gaussian Mixture Model (GMM).

- For a mixture of  $K$  Gaussians,  $p(\mathbf{x}_n | \Theta)$  will be

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

- The MLE problem for GMM would be

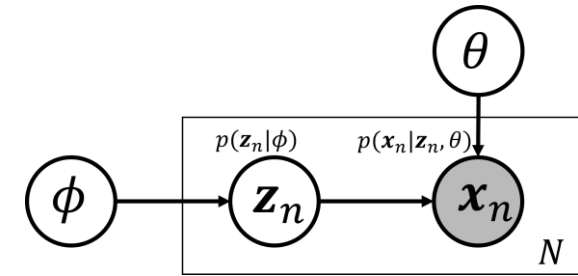
$$\operatorname{argmax}_{\Theta} \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)$$

ALT-OPT or EM makes it simpler by using hard/soft guesses of  $\mathbf{z}_n$ 's

The log of sum doesn't give us a simple expression; MLE can still be done using gradient based methods but updates will be complicated.

# How to Guess $\mathbf{z}_n$ in an LVM?

- Note that  $\mathbf{z}_n$  is a random variable with prior distribution  $p(\mathbf{z}_n|\phi)$
- Can compute its **conditional posterior** (CP) distribution as



Called conditional posterior because it is conditioned on data as well as  $\Theta$  (assuming we have already estimated  $\Theta$ )

$$(\theta, \phi) = \Theta$$

$$p(\mathbf{z}_n|\mathbf{x}_n, \Theta) = \frac{p(\mathbf{z}_n|\Theta)p(\mathbf{x}_n|\mathbf{z}_n, \Theta)}{p(\mathbf{x}_n|\Theta)} = \frac{p(\mathbf{z}_n|\phi)p(\mathbf{x}_n|\mathbf{z}_n, \theta)}{p(\mathbf{x}_n|\Theta)}$$

- If we just want the single best (hard) guess of  $\mathbf{z}_n$  then that can be computed as

Used in ALT-OPT for LVMs

$$\hat{\mathbf{z}}_n = \operatorname{argmax}_{\mathbf{z}_n} p(\mathbf{z}_n|\mathbf{x}_n, \Theta) = \operatorname{argmax}_{\mathbf{z}_n} p(\mathbf{z}_n|\phi)p(\mathbf{x}_n|\mathbf{z}_n, \theta)$$

Used in Expectation-Maximization (EM) algo for LVMs

- Otherwise, we can compute and use CP  $p(\mathbf{z}_n|\mathbf{x}_n, \Theta)$  to get a soft/probabilistic guess
  - Using the CP  $p(\mathbf{z}_n|\mathbf{x}_n, \Theta)$  we can compute quantities such as **expectation** of  $\mathbf{z}_n$
  - If  $p(\mathbf{z}_n|\phi)$  and  $p(\mathbf{x}_n|\mathbf{z}_n, \theta)$  are conjugate to each other then CP  $p(\mathbf{z}_n|\mathbf{x}_n, \Theta)$  is easy to compute
- Computing hard guess is usually easier but ignores the uncertainty in  $\mathbf{z}_n$

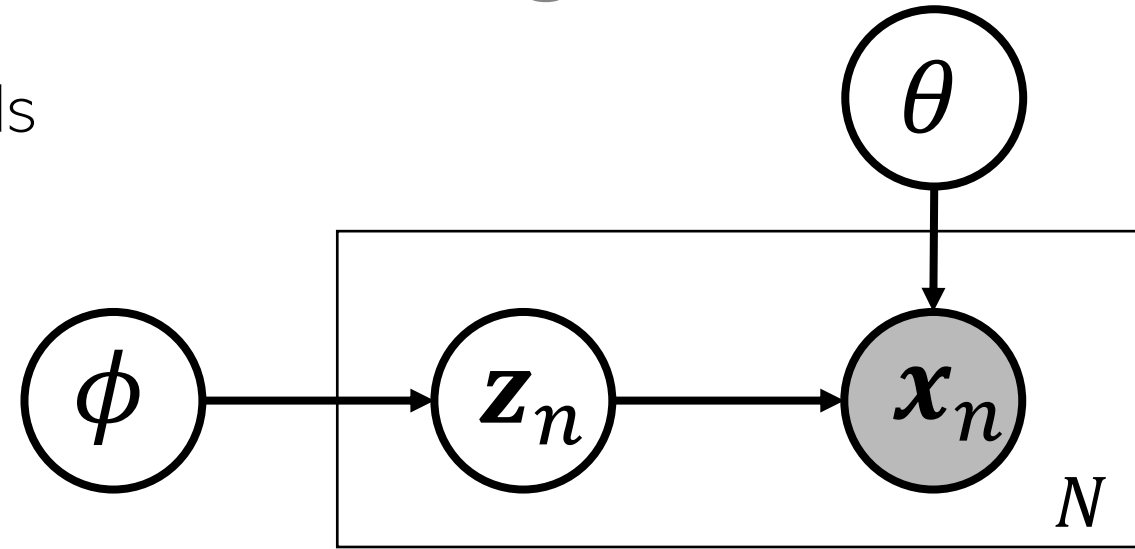


# LVMs: Incomplete vs Complete Data Log Likelihood<sup>6</sup>

- We can define two types of likelihoods for LVMs

- Incomplete data log likelihood (ILL)  $\log p(\mathbf{X}|\Theta)$

- Complete data log likelihood (CLL)  $\log p(\mathbf{X}, \mathbf{Z}|\Theta)$



- Named so because we can think of latent  $\mathbf{Z}$  “completing” the observed data  $\mathbf{X}$
- Since  $\mathbf{Z}$  is never observed (is latent), to estimate  $\Theta$  we must maximize the ILL

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

- But since ILL maximization is hard (log of sum/integral over the unknown  $\mathbf{Z}$ ), we instead maximize the CLL  $p(\mathbf{X}, \mathbf{Z}|\Theta)$  using hard/soft guesses of  $\mathbf{Z}$



# MLE for LVM

- If using a hard guess

$$\Theta_{MLE} = \operatorname{argmax}_{\Theta} \log p(\mathbf{X}, \hat{\mathbf{Z}} | \Theta)$$

- If using a soft (probabilistic) guess

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

- In LVMs, hard and soft guesses of  $\mathbf{Z}$  would depend on  $\Theta$  (since  $\mathbf{Z}$  and  $\Theta$  are coupled)
- Thus we need a procedure which alternates between estimating  $\mathbf{Z}$  and estimating  $\Theta$

Also, we can use this idea to find MAP solution of  $\Theta$  if we want. Assume a prior  $p(\Theta)$  and simply add a  $\log p(\Theta)$  term to these objectives

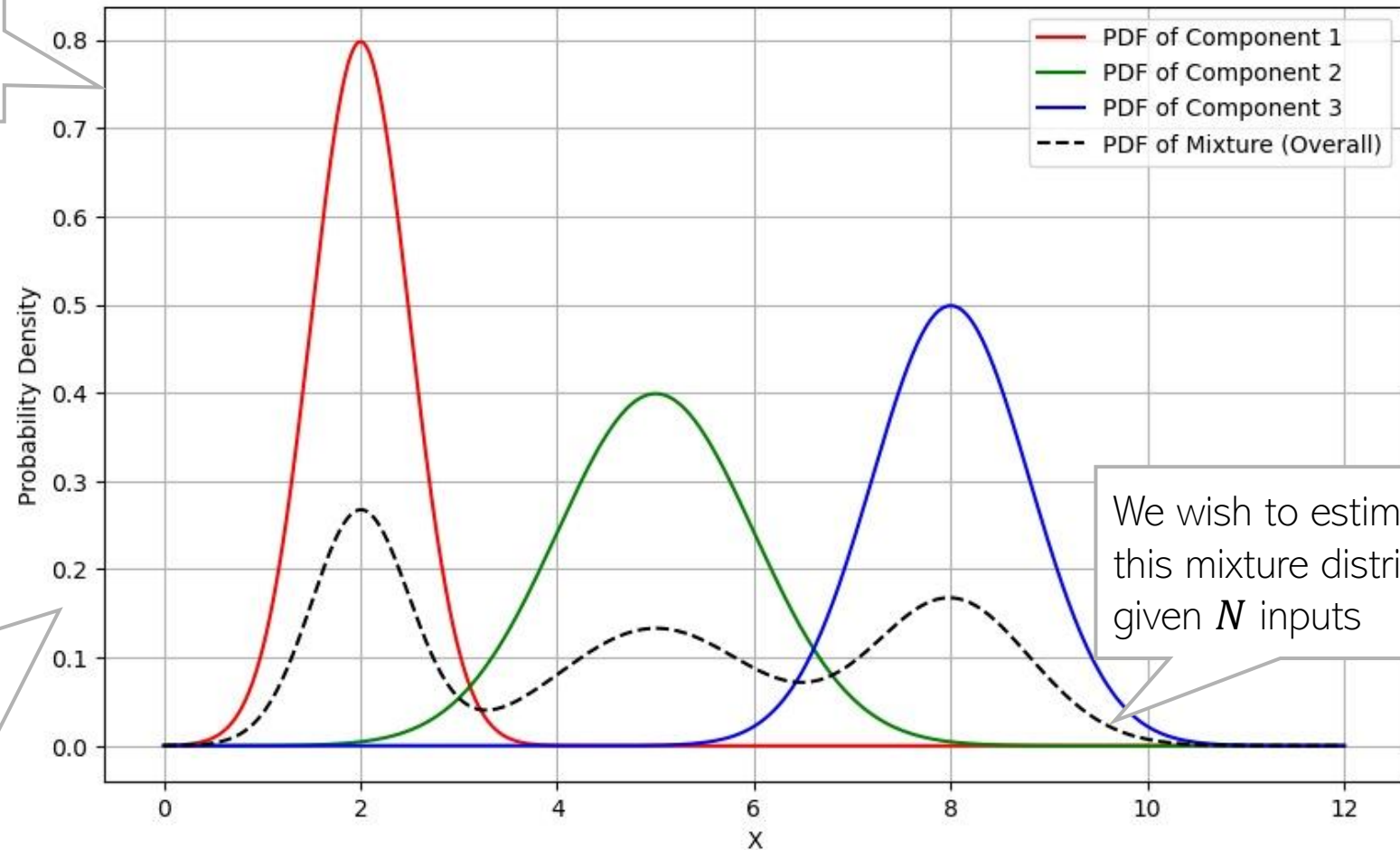
Note that we aren't solving the original MLE problem  $\operatorname{argmax}_{\Theta} \log p(\mathbf{X} | \Theta)$  anymore. However, what we are solving now is still justifiable theoretically (will see later)



# An LVM: Gaussian Mixture Model

Inputs are assumed generated from a mixture of Gaussians. But we don't know which input was generated by which Gaussian

If we knew which input came from which Gaussian (akin to knowing their true labels), the problem is easy – simply estimate each Gaussian using the inputs that came from that Gaussian (just like generative classification)



We wish to estimate this mixture distribution given  $N$  inputs





# 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  $\Sigma_k$
- Can assume label  $\mathbf{z}_n$  to be one-hot and then  $\mathbf{z}_{nk} = 1$  if  $\mathbf{z}_n = k$ , and  $\mathbf{z}_{nk} = 0$ , o/w
  - Note: For each label, using notation  $\mathbf{z}_n$  instead of  $\mathbf{y}_n$
- Assuming class marginal  $p(\mathbf{z}_n = k) = \pi_k$ , the model's params  $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$
- The MLE objective  $\log p(\mathbf{X}, \mathbf{Z}|\Theta)$  is (will provide a note for the proof)

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

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

Same as  $\frac{N_k}{N}$

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

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



# MLE for GMM: Using Guesses of $z_n$

Will have the exact same form for the expression of MLE objective as generative classification with Gaussian class-conditionals (except  $z_n$  is unknown)

- Using a hard guess  $\hat{z}_n = \operatorname{argmax}_{z_n} p(z_n | x_n, \Theta)$ , the MLE problem for GMM

Log likelihood of  $\Theta$  w.r.t. data  $\mathbf{X}$  and hard guesses  $\hat{\mathbf{Z}}$  of cluster ids

Assuming  $x_n$  given  $z_n$  and  $\Theta$  are i.i.d.

$$\Theta_{MLE} = \operatorname{argmax}_{\Theta} \log p(\mathbf{X}, \hat{\mathbf{Z}} | \Theta) = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K \hat{z}_{nk} [\log \pi_k + \log \mathcal{N}(x_n | \mu_k, \Sigma_k)]$$

- Using a soft guess  $\mathbb{E}[z_n]$ , the MLE problem for GMM

Expected log likelihood of  $\Theta$  w.r.t. data  $\mathbf{X}$  and  $\mathbf{Z}$

$z_{nk}$  appears at only one place in the log likelihood expression so easily replaced by expectation of  $z_{nk}$  w.r.t the CP  $p(z_n | x_n, \Theta)$

$$\Theta_{MLE} = \operatorname{argmax}_{\Theta} \mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \Theta)] = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}[z_{nk}] [\log \pi_k + \log \mathcal{N}(x_n | \mu_k, \Sigma_k)]$$

- In both cases, the MLE solution for  $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  will be identical to that of generative classification with Gaussian class cond with  $z_{nk}$  replaced by  $\hat{z}_{nk}$  or  $\mathbb{E}[z_{nk}]$ 
  - Case 1 solved using **ALT-OPT** alternating b/w estimating  $\Theta_{MLE}$  and  $\hat{\mathbf{Z}}$
  - Case 2 solved using **Expectation Maximization (EM)** alternating b/w estimating  $\Theta_{MLE}$  and  $\mathbb{E}[\mathbf{Z}]$



# ALT-OPT for GMM

- We will assume we have a “hard” (most probable) guess of  $\mathbf{z}_n$ , say  $\hat{\mathbf{z}}_n$
- ALT-OPT which maximizes  $\log p(\mathbf{X}, \hat{\mathbf{Z}} | \Theta)$  would look like this
  - Initialize  $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  as  $\hat{\Theta}$
  - Repeat the following until convergence

Proportional to prior prob times likelihood, i.e.,  
 $p(z_n = k | \hat{\Theta}) p(x_n | z_n = k, \hat{\Theta}) = \hat{\pi}_k \mathcal{N}(x_n | \hat{\mu}_k, \hat{\Sigma}_k)$

Posterior probability of point  $\mathbf{x}_n$  belonging to cluster  $k$ , given current  $\Theta$

- For each  $n$ , compute most probable value (our best guess) of  $\mathbf{z}_n$  as

$$\hat{\mathbf{z}}_n = \operatorname{argmax}_{k=1,2,\dots,K} p(z_n = k | \hat{\Theta}, \mathbf{x}_n)$$

- Solve MLE problem for  $\Theta$  using most probable  $\mathbf{z}_n$ 's

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

$$\begin{aligned} \hat{\pi}_k &= \frac{1}{N} \sum_{n=1}^N \hat{z}_{nk} & \hat{\mu}_k &= \frac{1}{N_k} \sum_{n=1}^N \hat{z}_{nk} \mathbf{x}_n \\ \hat{\Sigma}_k &= \frac{1}{N_k} \sum_{n=1}^N \hat{z}_{nk} (\mathbf{x}_n - \hat{\mu}_k)(\mathbf{x}_n - \hat{\mu}_k)^\top \end{aligned}$$

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



# Expectation-Maximization (EM) for GMM

- EM finds  $\Theta_{MLE}$  by maximizing  $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$
- Note: Expectation will be w.r.t. the CP of  $\mathbf{Z}$ , i.e.,  $p(\mathbf{Z}|\mathbf{X}, \Theta)$
- 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 CP  $p(\mathbf{Z}|\mathbf{X}, \hat{\Theta})$  using current estimate of  $\Theta$ . Since obs are i.i.d, compute for each  $n$  (and for  $k = 1, 2, \dots, K$ )

Why w.r.t. this distribution?  
Will see justification in a bit

Note that EM for GMM also gives a soft clustering  $\mathbf{z}_n = [\gamma_{n1}, \gamma_{n2}, \dots, \gamma_{nK}]$  for each input  $\mathbf{x}_n$



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

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

- Update  $\Theta$  by maximizing  $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$

$$\frac{\hat{\pi}_k \mathcal{N}(\mathbf{x}_n | \hat{\mu}_k, \hat{\Sigma}_k)}{\sum_{\ell=1}^K \hat{\pi}_\ell \mathcal{N}(\mathbf{x}_n | \hat{\mu}_\ell, \hat{\Sigma}_\ell)}$$

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

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

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

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

$$\begin{aligned} \mathbb{E}[\mathbf{z}_{nk}] &= \gamma_{nk} = 0 \times p(\mathbf{z}_{nk} = 0 | \mathbf{x}_n, \hat{\Theta}) + 1 \times p(\mathbf{z}_{nk} = 1 | \mathbf{x}_n, \hat{\Theta}) \\ &= p(\mathbf{z}_{nk} = 1 | \mathbf{x}_n, \hat{\Theta}) \\ &\propto \hat{\pi}_k \mathcal{N}(\mathbf{x}_n | \hat{\mu}_k, \hat{\Sigma}_k) \end{aligned}$$

Posterior probability of  $\mathbf{x}_n$  belonging to  $k^{th}$  cluster

# EM for GMM (Contd)

## 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 are more of a heuristic to get soft-clustering, also gave us probabilities but didn't account for cluster shapes or fraction of points in each cluster

Accounts for fraction of points in each cluster

$$\mathbb{E}[z_{nk}^{(t)}] = \gamma_{nk}^{(t)}$$

$$= \frac{\pi_k^{(t-1)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k^{(t-1)}, \boldsymbol{\Sigma}_k^{(t-1)})}{\sum_{\ell=1}^K \pi_{\ell}^{(t-1)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_{\ell}^{(t-1)}, \boldsymbol{\Sigma}_{\ell}^{(t-1)})} \quad \forall n, k$$

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  $\Theta$  via MLE

$$\boldsymbol{\mu}_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk}^{(t)} \mathbf{x}_n$$

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

M-step:

$$\boldsymbol{\Sigma}_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk}^{(t)} (\mathbf{x}_n - \boldsymbol{\mu}_k^{(t)}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{(t)})^T$$

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

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



# What is EM Doing?

14

Maximization of ILL

Assuming  $\mathbf{Z}$  to be discrete, else replace it by an integral

- The MLE problem was  $\Theta_{MLE} = \underset{\Theta}{\operatorname{argmax}} \log p(\mathbf{X}|\Theta) = \underset{\Theta}{\operatorname{argmax}} \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta)$   
As an approximation
- What EM (and ALT-OPT) maximized is expected CLL:  $\Theta_{MLE} = \underset{\Theta}{\operatorname{argmax}} \mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$
- We did not solve the original problem (max of ILL). **Is it okay?**
- Assume  $p_{\mathbf{Z}} = p(\mathbf{Z}|\mathbf{X}, \Theta)$  and  $q(\mathbf{Z})$  to be some prob distribution over  $\mathbf{Z}$ , then

Function of a distribution  $q$  and parameter  $\Theta$ 

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

May verify this identity

- In the above  $\mathcal{L}(q, \Theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \right\}$  and  $KL(q||p_{\mathbf{Z}}) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \Theta)}{q(\mathbf{Z})} \right\}$
- Since KL is always non-negative  $\log p(\mathbf{X}|\Theta) \geq \mathcal{L}(q, \Theta)$ , so  $\mathcal{L}(q, \Theta)$  is a lower-bound on ILL
- Thus if we maximize  $\mathcal{L}(q, \Theta)$ , it will also improve  $\log p(\mathbf{X}|\Theta)$



# What is EM Doing?

- As we saw,  $\mathcal{L}(q, \Theta)$  depends on  $q$  and  $\Theta$
- Let's maximize  $\mathcal{L}(q, \Theta)$  w.r.t.  $q$  with  $\Theta$  fixed at  $\Theta^{\text{old}}$

The posterior distribution of  $\mathbf{Z}$  given older parameters  $\Theta^{\text{old}}$  (will need this posterior to get the expectation of CLL)

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  $\hat{q} = p_z = p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})$

$$\begin{aligned} \Theta^{\text{new}} &= \operatorname{argmax}_{\Theta} \mathcal{L}(\hat{q}, \Theta) = \operatorname{argmax}_{\Theta} \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})} \right\} \\ &= \operatorname{argmax}_{\Theta} \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) \\ &= \operatorname{argmax}_{\Theta} \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})} [\log p(\mathbf{X}, \mathbf{Z}|\Theta)] \\ &= \operatorname{argmax}_{\Theta} Q(\Theta, \Theta^{\text{old}}) \end{aligned}$$

Maximization of expected CLL w.r.t. the posterior distribution of  $\mathbf{Z}$  given older parameters  $\Theta^{\text{old}}$





# Recap: ALT-OPT vs EM

- ALT-OPT does the following

- 1 Initialize  $\Theta = \hat{\Theta}$
- 2 Estimate  $\mathbf{Z}$  as  $\hat{\mathbf{Z}} = \arg \max_{\mathbf{Z}} \log p(\mathbf{Z}|\mathbf{X}, \hat{\Theta})$
- 3 Estimate  $\Theta$  as  $\hat{\Theta} = \arg \max_{\Theta} \log p(\mathbf{X}, \hat{\mathbf{Z}}|\Theta)$
- 4 Go to step 2 if not converged

This step could potentially throw away a lot of information about the latent variable  $\mathbf{Z}$

- EM addresses it using “soft” version of ALT-OPT

- 1 Initialize  $\Theta = \hat{\Theta}$
- 2 Compute the posterior distribution of  $\mathbf{Z}$ , i.e.,  $p(\mathbf{Z}|\mathbf{X}, \hat{\Theta})$
- 3 Estimate  $\Theta$  by maximizing the expected CLL  $\hat{\Theta} = \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \hat{\Theta})} [\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$
- 4 Go to step 2 if not converged

ALT-OPT can be seen as an approximation of EM – the posterior  $p(\mathbf{Z}|\mathbf{X}, \Theta)$  is replaced by a point mass at its mode





# EM: An Illustration

- As we saw, EM maximizes the lower bound  $\mathcal{L}(q, \Theta)$  in two steps
- Step 1 finds the optimal  $q$  (call it  $\hat{q}$ ) by setting it 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

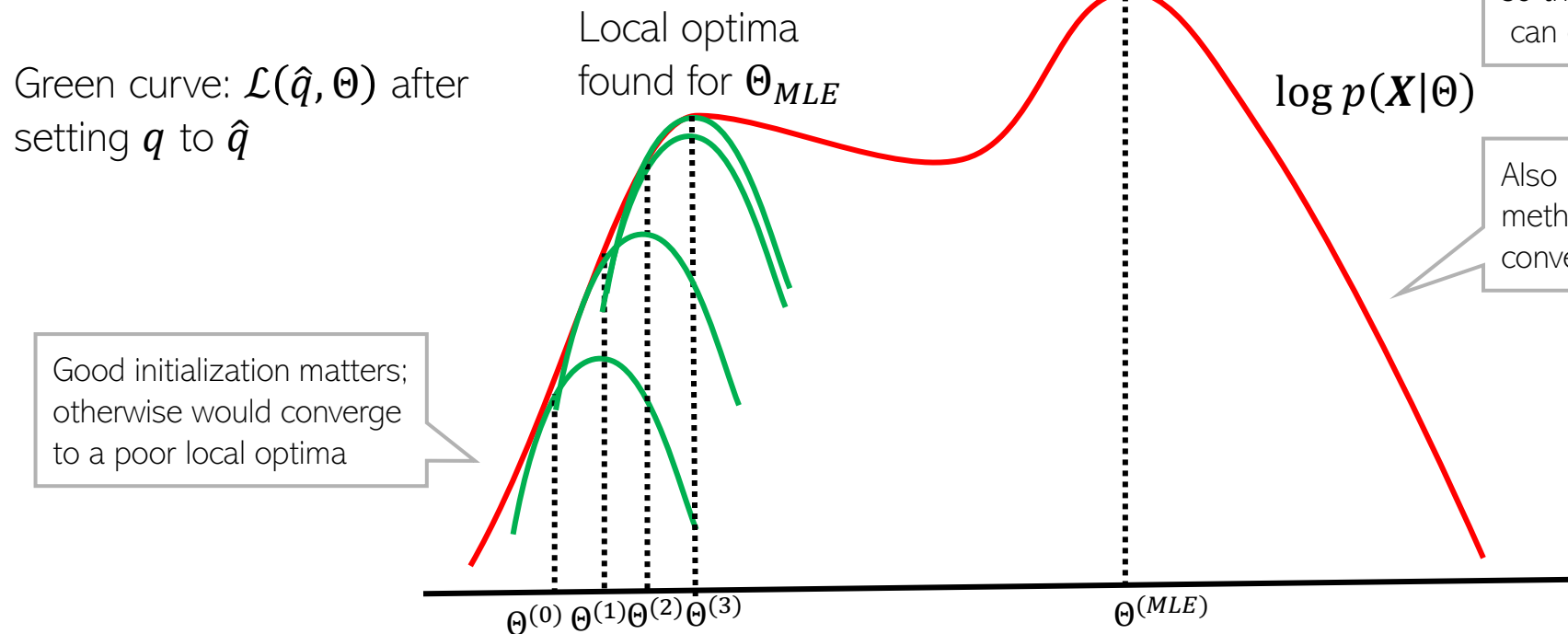
Makes  $\mathcal{L}(q, \Theta)$  equal to  $\log p(\mathbf{X}|\Theta)$ ; thus the curves touch at current  $\Theta$

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

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



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



# The EM Algorithm in its general form..

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

## The EM Algorithm

- Initialize  $\Theta$  as  $\Theta^{(0)}$ , set  $t = 1$
- Step 1: Compute **posterior** of latent variables given current parameters  $\Theta^{(t-1)}$

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

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

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

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

- Note: If we can take the MAP estimate  $\hat{\mathbf{z}}_n$  of  $\mathbf{z}_n$  (not full posterior) in Step 2 and maximize the CLL in Step 3 using that, i.e., do  $\arg \max_{\Theta} \sum_{n=1}^N [\log p(\mathbf{x}_n, \hat{\mathbf{z}}_n^{(t)} | \Theta)]$  this will be ALT-OPT



# The Expected CLL

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

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

Was indeed the case of GMM:  $p(\mathbf{z}_n|\Theta)$  was multinoulli,  $p(\mathbf{x}_n|\mathbf{z}_n, \Theta)$  was Gaussian

- If  $p(\mathbf{z}_n|\Theta)$  and  $p(\mathbf{x}_n|\mathbf{z}_n, \Theta)$  are exponential family distributions, then  $Q(\Theta, \Theta^{old})$  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 (beyond the scope of CS771)



# Detour: Exponential Family

- Exponential Family is a family of prob. distributions that have the form

$$p(x|\theta) = h(x)\exp[\theta^\top T(x) - A(\theta)]$$

Even though their standard form may not look like this, they can be rewritten in this form after some algebra

- Many well-known distribution (Bernoulli, Binomial, multinoulli, Poisson, beta, gamma, Gaussian, etc.) are examples of exponential family distributions
- $\theta$  is called the **natural parameter** of the family
- $h(x)$ ,  $T(x)$ , and  $A(\theta)$  are known functions (specific to the distribution)
- $T(x)$  is called the **sufficient statistics**: estimates of  $\theta$  contain  $x$  in form of suff-stats
- Every exp. family distribution also has a conjugate distribution (often also in exp. family)
- Also, MLE/MAP is usually quite simple since  $\log p(x|\theta)$  will have a simple expression
- Also useful in fully Bayesian inference since they have conjugate priors

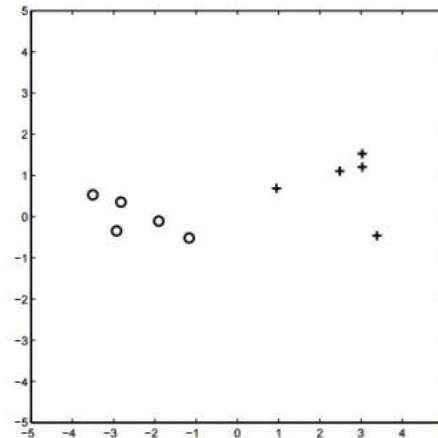
Natural params are a function of the distribution parameters in the standard form

[https://en.wikipedia.org/wiki/Exponential\\_family](https://en.wikipedia.org/wiki/Exponential_family)



# LVM for Semi-supervised Learning

- Unlabeled data can help in supervised learning as well

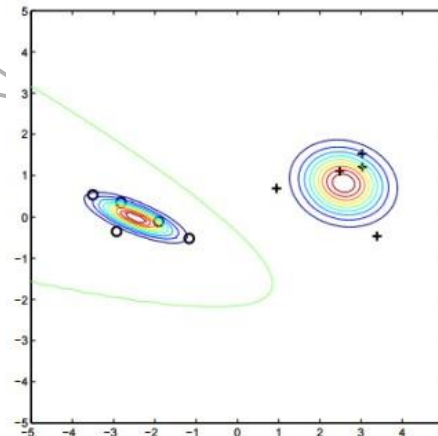


labeled data

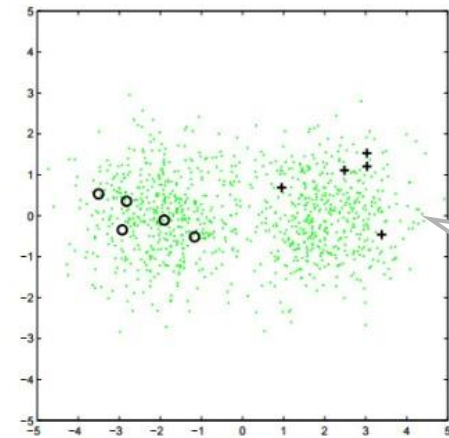
A generative classifier learned using the labelled examples

Decision boundary (green curve)

A small training set may not be able to help learn the true decision boundary



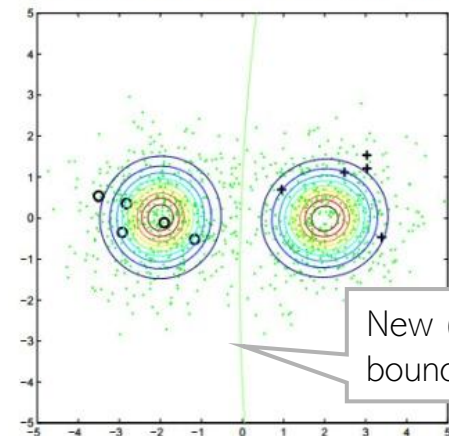
model learned from labeled data



labeled and unlabeled data (small dots)

Added some unlabeled inputs (green points)

We don't know their true labels but can treat the true labels as latent variables and estimate them along with the rest of the parameters of the generative classification model



model learned from labeled and unlabeled data

New (improved) decision boundary



# LVM for Semi-supervised Learning (SSL)

- Suppose we have  $N$  labeled  $\{\mathbf{x}_n, y_n\}_{n=1}^N$  and  $M$  unlabeled examples  $\{\mathbf{x}_n\}_{n=N+1}^{N+M}$
- We wish to learn a classifier using both labelled and unlabeled examples
- We can treat the labels of  $\{\mathbf{x}_n\}_{n=N+1}^{N+M}$  as latent variables and use ALT-OPT or EM
  - $z_n = y_n$  for labeled examples  $n = 1, 2, \dots, N$
  - $z_n$  estimated (hard/soft guess) for unlabeled examples  $n = N + 1, \dots, N + M$
- Assuming generative classification with Gaussian class-conditional ( $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ )

This SSL model is a hybrid of supervised generative classification (with Gaussian class-conditionals) and GMM

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

Assuming we are using EM (soft guess), otherwise ALT-OPT (hard guess) can be used too, as we did in GMM



# Another LVM: Probabilistic PCA (PPCA)

- Assume  $\mathbf{x}_n \in \mathbb{R}^D$  as a linear mapping of a latent var  $\mathbf{z}_n \in \mathbb{R}^K$  + Gaussian noise

A “reverse” generative way of thinking about PCA (low-dim  $\mathbf{z}_n$  generating high-dim  $\mathbf{x}_n$ )

$D \times 1$  offset/bias

$D \times K$  matrix

Drawn from a zero-mean  $D$ -dim Gaussian  $\mathcal{N}(\mathbf{0}, \sigma^2 I_D)$

This linear mapping can be replaced by more powerful nonlinear mapping of the form  $\mathbf{x}_n = \mathbf{f}(\mathbf{z}_n) + \epsilon_n$  where  $\mathbf{f}$  can be modeled using a deep neural net (e.g., models like variational autoencoders)

$$\mathbf{x}_n = \boldsymbol{\mu} + \mathbf{W}\mathbf{z}_n + \boldsymbol{\epsilon}_n$$

- Equivalent to saying  $p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu} + \mathbf{W}\mathbf{z}_n, \sigma^2 I_D)$
- Assume a zero-mean  $K$ -dim Gaussian prior on  $\mathbf{z}_n$ , so  $p(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n | \mathbf{0}, I_K)$
- We would like to do MLE for  $\Theta = (\boldsymbol{\mu}, \mathbf{W}, \sigma^2)$
- ILL for this model  $p(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}, \sigma^2)$  is also a Gaussian (thanks to Gaussian properties)

$$p(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \int p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\mu}, \mathbf{W}, \sigma^2) p(\mathbf{z}_n) d\mathbf{z}_n = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 I_D)$$

- Maximizing ILL w.r.t.  $\Theta = (\boldsymbol{\mu}, \mathbf{W}, \sigma^2)$  is possible but requires solving eig decomp. problem
- We can use ALT-OPT/EM to estimate  $\Theta = (\boldsymbol{\mu}, \mathbf{W}, \sigma^2)$  more efficiently without eig decomp.

PRML 12.2.1



# Learning PPCA using EM

- Instead of maximizing the ILL  $p(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = N(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D)$ , let's use ALT-OPT/EM
- EM will instead maximize expected CLL, with CLL (assume  $\boldsymbol{\mu} = \mathbf{0}$ ) given by

$$\log p(X, Z | W, \sigma^2) = \log \prod_{n=1}^N p(\mathbf{x}_n, \mathbf{z}_n | \mathbf{W}, \sigma^2) = \log \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \sigma^2) p(\mathbf{z}_n)$$

- Using  $p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \sigma^2) = \mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n, \sigma^2 \mathbf{I}_D)$  and  $p(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I}_K)$

$$\text{CLL} = - \sum_{n=1}^N \left\{ \frac{D}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \|\mathbf{x}_n\|^2 - \frac{1}{\sigma^2} \mathbf{z}_n^\top \mathbf{W}^\top \mathbf{x}_n + \frac{1}{2\sigma^2} \text{trace}(\mathbf{z}_n \mathbf{z}_n^\top \mathbf{W}^\top \mathbf{W}) + \frac{1}{2} \text{trace}(\mathbf{z}_n \mathbf{z}_n^\top) \right\}$$

- Expected CLL will require  $\mathbb{E}[\mathbf{z}_n]$  and  $\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top]$  w.r.t. conditional posterior of  $\mathbf{z}_n$

Using the fact that  $p(\mathbf{x}_n | \mathbf{z}_n)$  and  $p(\mathbf{z}_n)$  are Gaussians and the CP is just the reverse conditional  $p(\mathbf{z}_n | \mathbf{x}_n)$  and must also be Gaussian

$$\begin{aligned} p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) &= \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \quad \text{where } \mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K \\ \mathbb{E}[\mathbf{z}_n] &= \mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n \\ \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] &= \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \text{cov}(\mathbf{z}_n) = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \sigma^2 \mathbf{M}^{-1} \end{aligned}$$



# Learning PPCA using EM

- The EM algo for PPCA alternates between two steps
  - Compute CP of  $\mathbf{z}_n$  given parameters  $\Theta = (\mathbf{W}, \sigma^2)$  and required expectatuions

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}) = \mathcal{N}(\mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n, \sigma^2 \mathbf{M}^{-1}) \quad \text{where } \mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K$$

$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^\top \mathbf{x}_n$$

$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \text{cov}(\mathbf{z}_n) = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \sigma^2 \mathbf{M}^{-1}$$

Note: This approach does not assume/ensure that  $\mathbf{W}$  is orthonormal

- Maximize the expected CLL  $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2)]$  w.r.t.  $\mathbf{W}$  and  $\sigma^2$

$$\mathbb{E}[\text{CLL}] = - \sum_{n=1}^N \left\{ \frac{D}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \|\mathbf{x}_n\|^2 - \frac{1}{\sigma^2} \mathbb{E}[\mathbf{z}_n^\top] \mathbf{W}^\top \mathbf{x}_n + \frac{1}{2\sigma^2} \text{trace}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] \mathbf{W}^\top \mathbf{W}) + \frac{1}{2} \text{trace}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top]) \right\}$$

Unlike standard PCA (non-probabilistic), no eigendecomposition needed to estimate  $\mathbf{W}$

$$\mathbf{W}_{\text{new}} = \left[ \sum_{n=1}^N \mathbf{x}_n \mathbb{E}[\mathbf{z}_n]^\top \right] \left[ \sum_{n=1}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] \right]^{-1}$$

$$\sigma_{\text{new}}^2 = \frac{1}{ND} \sum_{n=1}^N \left\{ \|\mathbf{x}_n\|^2 - 2 \mathbb{E}[\mathbf{z}_n]^\top \mathbf{W}_{\text{new}}^\top \mathbf{x}_n + \text{tr} \left( \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] \mathbf{W}_{\text{new}}^\top \mathbf{W}_{\text{new}} \right) \right\}$$

Note: setting  $\sigma^2 = 0$  makes it equivalent to standard PCA without orthonormality constraint, but EM is **more efficient** since no eigendecomposition is needed

- Will get ALT-OPT if we use mode of the CP as  $\hat{\mathbf{z}}_n$  in the CLL

# Generative Models can generate synthetic data!

26

- Once parameters  $\Theta = (\mu, W, \sigma^2)$  are learned, we can even generate new data, e.g.,
  - Generate a random  $\mathbf{z}_n$  from  $\mathcal{N}(\mathbf{0}, I_K)$
  - Generate  $\mathbf{x}_n$  condition on  $\mathbf{z}_n$  from  $\mathcal{N}(\mu + W\mathbf{z}_n, \sigma^2 I_D)$

In addition to, of course, reducing the data dimensionality



(a) Training data



(b) Random samples

Generated using a more sophisticated generative model, not PPCA (but similar in formulation)

Methods such as variational autoencoders, GAN, diffusion models, etc are based on similar ideas



# EM: Some Comments

- Good initialization is important
- The E and M steps may not always be possible to perform exactly. Some reasons
  - CP of latent variables  $p(\mathbf{Z}|\mathbf{X}, \Theta)$  may not be easy to find and may require approx.
  - Even if  $p(\mathbf{Z}|\mathbf{X}, \Theta)$  is easy, expected CLL, i.e.,  $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$  may still not be tractable

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

Monte-Carlo EM

..and may need to be approximated, e.g., using [Monte-Carlo expectation](#)

Gradient methods may still be needed for this step

- 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](#))
- Other advanced probabilistic inference algorithms are based on ideas similar to EM
  - E.g., Variational Bayesian inference a.k.a. [Variational Inference \(VI\)](#)
- EM is also related to non-convex optimization algorithms [Majorization-Maximization \(MM\)](#)
  - MM maximizes a difficult-to-optimize objective function by iteratively constructing surrogate functions that are easier to maximize (in EM, the surrogate function was the CLL)

