# Latent Variable Models

CS771: Introduction to Machine Learning

# Dimensionality Reduction: Out-of-sample Embedding

- Some dim-red methods can only compute the embedding of the training data

- Given $N$ training samples $\{x_1, x_2, \ldots, x_N\}$ they will give their embedding $\{z_1, z_2, \ldots, z_N\}$

- However, given a new point $x_*$ (not in the training samples), they can't produce its embedding $z_*$ easily
  - Thus no easy way of getting "out-of-sample" embedding

- Some of the nonlinear dim-red methods like LLE, SNE, KPCA, etc have this limitation
  - Reason: They don't learn an explicit encoder and directly optimize for $\{z_n\}_{n=1}^N$ given $\{x_n\}_{n=1}^N$
  - To get "out-of-sample" embeddings, these methods require some modifications*

- But many other methods do explicitly learn a mapping $z = f(x)$ in form of an "encoder" $f$ that can give $z_*$ for any new $x_*$ as well (such methods are more useful)
  - For PCA, the $D \times K$ projection matrix $W_K$ is this encoder function and $z_* = W_K^\mathsf{T} x_*$
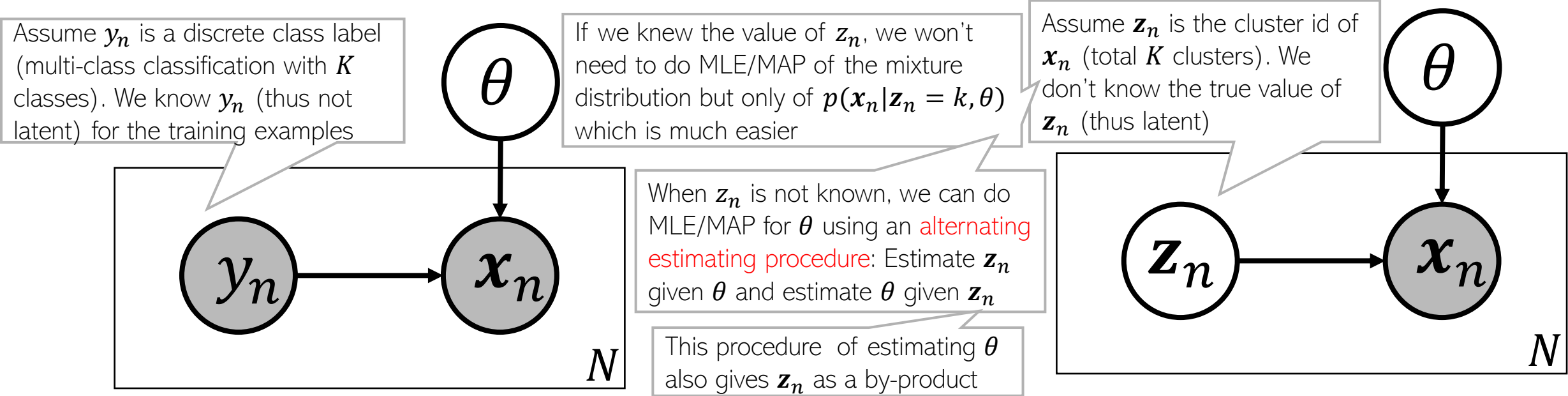  - Neural network based autoencoders can also do this (will see them later)

*Out-of-Sample Extensions for LLE, Isomap, MDS, Eigenmaps, and Spectral Clustering (Bengio et al, 2003)

# Latent Variable Models

# Example: Generative Models with Latent Variables

- Two generative models of inputs $x_n$ without (left) and with (right) latent variables

Assume $y_n$ is a discrete class label (multi-class classification with $K$ classes). We know $y_n$ (thus not latent) for the training examples

If we knew the value of $z_n$, we won't need to do MLE/MAP of the mixture distribution but only of $p(x_n|z_n = k, \theta)$ which is much easier

Assume $z_n$ is the cluster id of $x_n$ (total $K$ clusters). We don't know the true value of $z_n$ (thus latent)

When $z_n$ is not known, we can do MLE/MAP for $\theta$ using an alternating estimating procedure: Estimate $z_n$ given $\theta$ and estimate $\theta$ given $z_n$

This procedure of estimating $\theta$ also gives $z_n$ as a by-product

- Suppose we wish to estimate (e.g., using MLE/MAP) params $\theta$ of distribution of $x_n$
- For case 1, the distribution is $p(x_n|y_n, \theta)$ 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

Reason: The functional form of mixture can be messy

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

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

# 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 $\mathbf{x}_n \in \mathbb{R}^D$

$$p(\mathbf{z}_n|\phi) \qquad p(\mathbf{x}_n|\mathbf{z}_n, \theta)$$

$\theta$

$\phi \longrightarrow \mathbf{z}_n \longrightarrow \mathbf{x}_n$

$N$

- 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 $\boldsymbol{z}_n$ is hard

- Reason: Given $N$ observations $x_n, n = 1, 2, \ldots, N$, the MLE problem for $\Theta$ will be

$$\arg\max_{\Theta} \sum_{n=1}^{N} \log p(\boldsymbol{x}_n | \Theta) = \arg\max_{\Theta} \sum_{n=1}^{N} \log \sum_{\boldsymbol{z}_n} p(\boldsymbol{x}_n, \boldsymbol{z}_n | \Theta)$$

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

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

Gaussian Mixture Model (GMM).

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

$$p(\boldsymbol{x}_n | \Theta) = \sum_{k=1}^{K} p(\boldsymbol{x}_n, \boldsymbol{z}_n = k | \Theta) = \sum_{k=1}^{K} p(\boldsymbol{z}_n = k | \phi) p(\boldsymbol{x}_n | \boldsymbol{z}_n = k, \theta) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \mu_k, \Sigma_k)$$
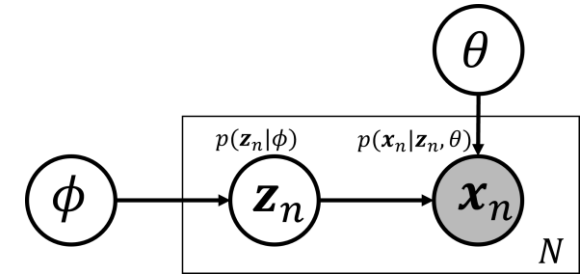
- The MLE problem for GMM would be

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

$$\arg\max_{\Theta} \sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \mu_k, \Sigma_k)$$

# How to Guess $z_n$ in an LVM?

- Note that $\mathbf{z}_n$ is a random variable with prior distribution $p(\mathbf{z}_n|\phi)$



- Can compute its <span style="color:red">conditional posterior</span> (CP) distribution as

> Called conditional posterior because it is conditioned on data as well as Θ (assuming we have already estimated Θ)

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

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

- 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

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

> 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 <span style="color:red">expectation</span> 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

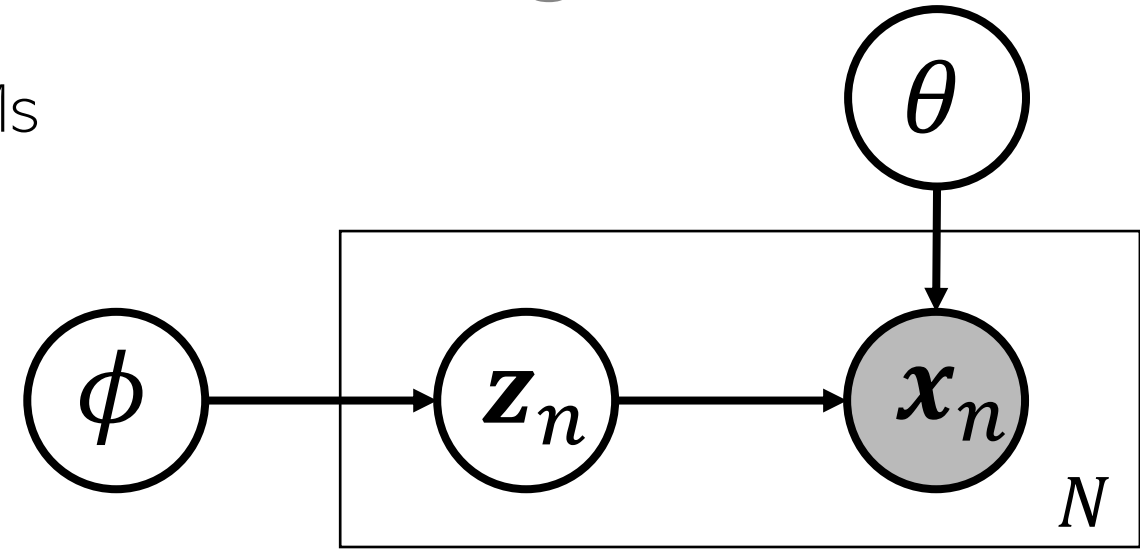- Computing hard guess is usually easier but ignores the uncertainty in $\mathbf{z}_n$

- We can define two types of likelihoods for LVMs

  - <span style="color:red">Incomplete data log likelihood (ILL)</span> $\log p(\boldsymbol{X}|\Theta)$

  - <span style="color:red">Complete data log likelihood (CLL)</span> $\log p(\boldsymbol{X}, \boldsymbol{Z}|\Theta)$

- Named so because we can think of latent $\boldsymbol{Z}$ "completing" the observed data $\boldsymbol{X}$

- Since $\boldsymbol{Z}$ is never observed (is latent), <span style="color:red">to estimate $\Theta$ we must maximize the ILL</span>

$$\operatorname*{argmax}_{\Theta} \log \textcolor{red}{p(\boldsymbol{X}|\Theta)} = \operatorname*{argmax}_{\Theta} \log \textcolor{red}{\sum_{\boldsymbol{Z}} p(\boldsymbol{X}, \boldsymbol{Z}|\Theta)}$$

- But since ILL maximization is hard (log of sum/integral over the unknown $\boldsymbol{Z}$), <span style="color:blue">we instead maximize the CLL $p(\boldsymbol{X}, \boldsymbol{Z}|\Theta)$ using hard/soft guesses of $\boldsymbol{Z}$</span>

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
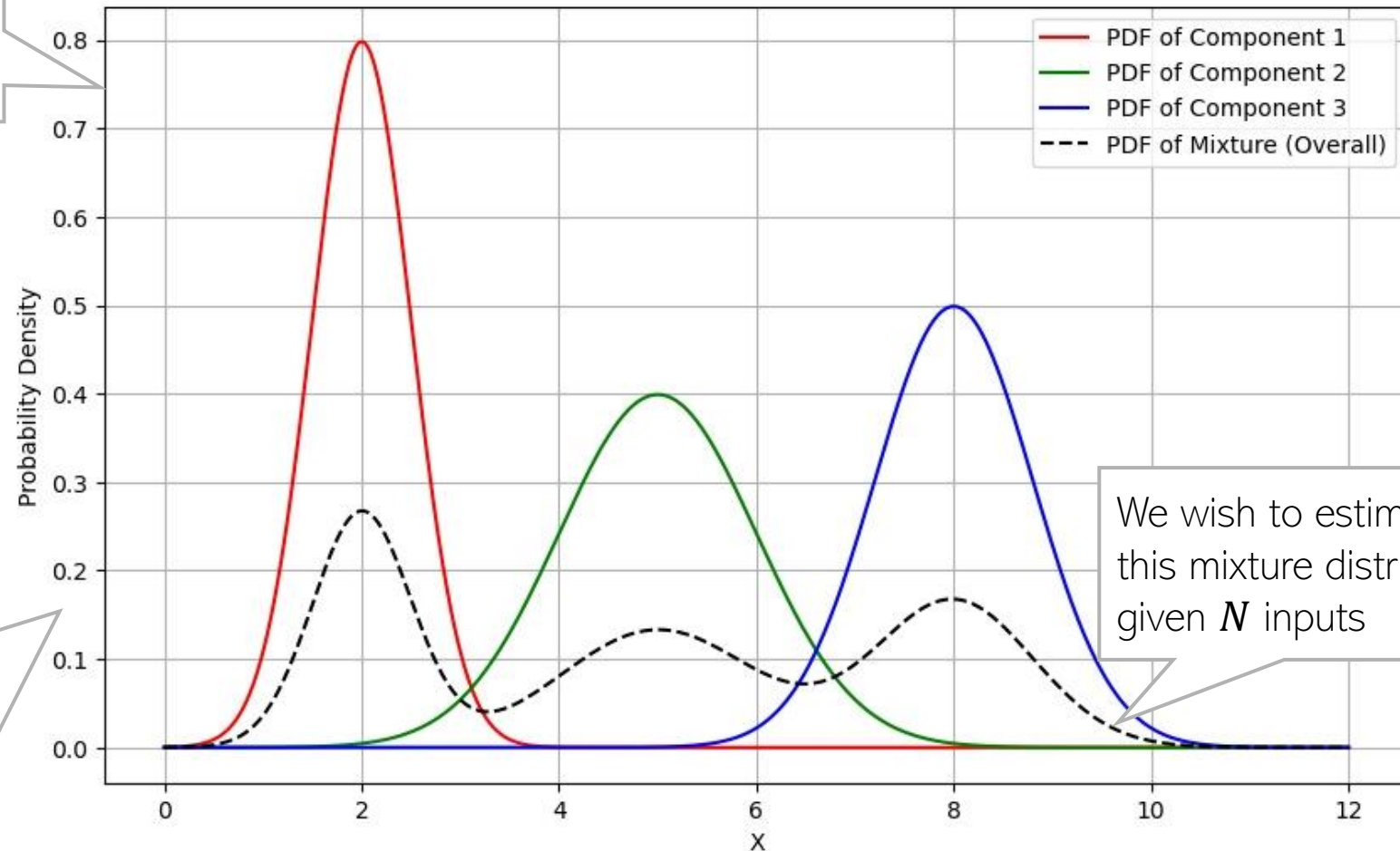
Note that we aren't solving the original MLE problem $\underset{\Theta}{\text{argmax}} \log p(X|\Theta)$ anymore.

However, what we are solving now is still justifiable theoretically (will see later)

- If using a hard guess

$$\Theta_{MLE} = \underset{\Theta}{\text{argmax}} \log p(X, \widehat{\mathbf{Z}} \,|\Theta)$$

- If using a soft (probabilistic) guess

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

- 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 $\Theta$

# 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 $k = 1, 2, \ldots, K$ is modeled by a Gaussian with mean $\mu_k$ and cov matrix $\Sigma_k$
- Can assume label $z_n$ to be one-hot and then $z_{nk} = 1$ if $z_n = k$, and $z_{nk} = 0$, o/w
  - Note: For each label, using notation $z_n$ instead of $y_n$
- Assuming class marginal $p(z_n = k) = \pi_k$, the model's params $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$
- The MLE objective $\log p(\boldsymbol{X}, \boldsymbol{Z}|\Theta)$ is (will provide a note for the proof)

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

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

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

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

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

# 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 $\hat{z}_n = \text{argmax}_{z_n} \, p(z_n|x_n, \Theta)$, the MLE problem for GMM

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

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

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

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

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

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

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

- 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 $\widehat{Z}$
  - Case 2 solved using Expectation Maximization (EM) alternating b/w estimating $\Theta_{MLE}$ and $\mathbb{E}[Z]$

# ALT-OPT for GMM

- We will assume we have a "hard" (most probable) guess of $z_n$, say $\hat{z}_n$

- ALT-OPT which maximizes $\log p(X, \hat{Z}|\Theta)$ 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 $x_n$ belonging to cluster $k$, given current $\Theta$

    - For each $n$, compute most probable value (our best guess) of $z_n$ as

$$\hat{z}_n = \text{argmax}_{k=1,2,\ldots,K} \; p(z_n = k|\hat{\Theta}, x_n)$$

    - Solve MLE problem for $\Theta$ using most probable $z_n$'s

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

$N_k$ : Effective number of points in cluster k

$$\hat{\pi}_k = \frac{1}{N} \sum_{n=1}^N \hat{z}_{nk} \qquad \hat{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \hat{z}_{nk} x_n$$

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \hat{z}_{nk} (x_n - \hat{\mu}_k)(x_n - \hat{\mu}_k)^\top$$

# Expectation-Maximization (EM) for GMM

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

- Note: Expectation will be w.r.t. the CP of $Z$, i.e., $p(Z|X, \Theta)$

- The EM algorithm for GMM operates as follows

  Note that EM for GMM also gives a soft clustering $z_n = [\gamma_{n1}, \gamma_{n2}, \ldots, \gamma_{nK}]$ for each input $x_n$

  - Initialize $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$ as $\widehat{\Theta}$

  - Repeat until convergence

    - Compute CP $p(Z|X, \widehat{\Theta})$ using current estimate of $\Theta$. Since obs are i.i.d, compute for each $n$ (and for $k = 1,2,\ldots K$)

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

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

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

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

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

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

$$\hat{\pi}_k = \frac{1}{N} \sum_{n=1}^{N} \mathbb{E}[z_{nk}] \qquad \hat{\mu}_k = \frac{1}{N_k} \sum_{n=1}^{N} \mathbb{E}[z_{nk}] x_n$$

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

$$\mathbb{E}[z_{nk}] = \gamma_{nk} = 0 \times p(z_{nk} = 0 | x_n, \widehat{\Theta}) + 1 \times p(z_{nk} = 1 | x_n, \widehat{\Theta})$$

$$= p(z_{nk} = 1 | x_n, \widehat{\Theta})$$

Posterior probability of $x_n$ belonging to $k^{th}$ cluster

$$\propto \hat{\pi}_k \mathcal{N}(x_n | \hat{\mu}_k, \hat{\Sigma}_k)$$

# 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)

Accounts for fraction of points in each cluster

Accounts for cluster shapes (since each cluster is a Gaussian

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

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

3. 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

M-step:

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

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

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

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