Generative Models for Clustering, GMM, and Intro to EM

Piyush Rai

Machine Learning (CS771A)

Sept 26, 2016

Machine Learning (CS771A)

Generative Models for Clustering, GMM, and Intro to EM

3

• A probabilistic way to think about the data generation process

$$\overset{p(z|\Phi)}{\longleftarrow} \underbrace{z} \overset{p(x|z,\theta)}{\longleftarrow} \underbrace{x} \underbrace{4}$$

3

• A probabilistic way to think about the data generation process

$$p(z|\Phi)$$
 z $p(x|z,\theta)$ x 4

Idea: First generate a random latent variable z from a prior distr. p(z|φ) and then generate x conditioned on z from the data distr. p(x|z, θ)

- -

• A probabilistic way to think about the data generation process

$$\overbrace{z}^{p(z|\Phi)} (z) \xrightarrow{(x|z,\theta)} (x)$$

- Idea: First generate a random latent variable z from a prior distr. p(z|φ) and then generate x conditioned on z from the data distr. p(x|z, θ)
- Some exceptions to this general view/definition of a generative model:

・ロト ・同ト ・ヨト ・ヨト ・ヨー

• A probabilistic way to think about the data generation process

$$\overbrace{z}^{p(z|\Phi)} (z) \xrightarrow{(x|z,\theta)} (x)$$

- Idea: First generate a random latent variable z from a prior distr. p(z|φ) and then generate x conditioned on z from the data distr. p(x|z, θ)
- Some exceptions to this general view/definition of a generative model:
 - Not all generative models have latent variables (e.g., for each observation, the other observations may play the role of latent variables)



3

イロト イポト イヨト イヨト

• A probabilistic way to think about the data generation process

$$\overbrace{z}^{p(z|\Phi)} (z) \xrightarrow{(x|z,\theta)} (x)$$

- Idea: First generate a random latent variable z from a prior distr. p(z|φ) and then generate x conditioned on z from the data distr. p(x|z, θ)
- Some exceptions to this general view/definition of a generative model:
 - Not all generative models have latent variables (e.g., for each observation, the other observations may play the role of latent variables)



• The z to x map may be a deterministic fn. (e.g., a neural or deep neural net)

イロン 不同と 不同と 不同と

• A probabilistic way to think about the data generation process

$$\overbrace{z}^{p(z|\Phi)} (z) \xrightarrow{(x|z,\theta)} (x)$$

- Idea: First generate a random latent variable z from a prior distr. p(z|φ) and then generate x conditioned on z from the data distr. p(x|z, θ)
- Some exceptions to this general view/definition of a generative model:
 - Not all generative models have latent variables (e.g., for each observation, the other observations may play the role of latent variables)



- The z to x map may be a deterministic fn. (e.g., a neural or deep neural net)
- We will focus on probabilistic generative models with latent variables

イロン 不同と 不同と 不同と

Generative Models for Clustering

-

- A generative model for data clustering
- Data assumed generated from a mixture of K Gaussians



3

イロン 不同と イヨン イヨン

- A generative model for data clustering
- Data assumed generated from a mixture of K Gaussians



- Let $0 \le \pi_k \le 1$ denote the "mixing weight" of the *k*-th Gaussian. It means:
 - π_k is the fraction of points generated from the k-th Gaussian

A D F A R F A B F A B F

- A generative model for data clustering
- Data assumed generated from a mixture of K Gaussians



- Let $0 \le \pi_k \le 1$ denote the "mixing weight" of the *k*-th Gaussian. It means:
 - π_k is the fraction of points generated from the k-th Gaussian
 - $\pi_k = p(\mathbf{z}_n = k)$ is the prior prob. of \mathbf{x}_n belonging to the k-th Gaussian

イロン 不同と 不同と 不同と

- A generative model for data clustering
- Data assumed generated from a mixture of K Gaussians



- Let $0 \le \pi_k \le 1$ denote the "mixing weight" of the *k*-th Gaussian. It means:
 - π_k is the fraction of points generated from the k-th Gaussian
 - $\pi_k = p(\mathbf{z}_n = k)$ is the prior prob. of \mathbf{x}_n belonging to the k-th Gaussian
- Let $\pi = (\pi_1, \pi_2, \dots, \pi_K)$ denote the vector of mixing wts of K Gaussians. This is a probability vector and sums to 1, i.e., $\sum_{k=1}^{K} \pi_k = 1$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ●□

- A generative model for data clustering
- Data assumed generated from a mixture of K Gaussians



- Let $0 \le \pi_k \le 1$ denote the "mixing weight" of the *k*-th Gaussian. It means:
 - π_k is the fraction of points generated from the k-th Gaussian
 - $\pi_k = p(\mathbf{z}_n = k)$ is the prior prob. of \mathbf{x}_n belonging to the k-th Gaussian
- Let $\pi = (\pi_1, \pi_2, \dots, \pi_K)$ denote the vector of mixing wts of K Gaussians. This is a probability vector and sums to 1, i.e., $\sum_{k=1}^{K} \pi_k = 1$
- Notation $z_n = k$ is equivalent to a size K one-hot vector z_n

$$\mathbf{z}_n = \underbrace{[0 \ 0 \ \dots \ 1 \ 0 \ 0]}_{\mathbf{z}_n}$$

all zeros except the k-th bit, i.e., $z_{nk} = 1$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ●□

- A generative model for data clustering
- Data assumed generated from a mixture of K Gaussians



- Let $0 \le \pi_k \le 1$ denote the "mixing weight" of the *k*-th Gaussian. It means:
 - π_k is the fraction of points generated from the k-th Gaussian
 - $\pi_k = p(\mathbf{z}_n = k)$ is the prior prob. of \mathbf{x}_n belonging to the k-th Gaussian
- Let $\pi = (\pi_1, \pi_2, \dots, \pi_K)$ denote the vector of mixing wts of K Gaussians. This is a probability vector and sums to 1, i.e., $\sum_{k=1}^{K} \pi_k = 1$
- Notation $z_n = k$ is equivalent to a size K one-hot vector z_n

$$\mathbf{z}_n = [\underbrace{0 \ 0 \ \dots \ 1 \ 0 \ 0}]$$

all zeros except the k-th bit, i.e., $z_{nk}\,=\,1$

• Note: The prior $p(\boldsymbol{z}|\pi)$ on each \boldsymbol{z}_n is a multinomial, i.e., $p(\boldsymbol{z}_n|\pi) = \prod_{k=1}^{K} \pi_k^{\boldsymbol{z}_{nk}}$

• The generative story for each \boldsymbol{x}_n , $n = 1, 2, \dots, N$

・ロト ・同ト ・ヨト ・ヨト ・ヨー

- The generative story for each \boldsymbol{x}_n , $n = 1, 2, \dots, N$
 - First choose one of the K mixture components as

 $z_n \sim \text{Multinomial}(z_n | \pi)$ (from the prior p(z) over z)

・ロト ・同ト ・ヨト ・ヨト ・ヨー

- The generative story for each \boldsymbol{x}_n , $n = 1, 2, \dots, N$
 - First choose one of the K mixture components as

 $z_n \sim \text{Multinomial}(z_n | \pi)$ (from the prior p(z) over z)

• Suppose $z_n = k$. Now generate x_n from the k-th Gaussian as

 $\mathbf{x}_n \sim \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ (from the data distr. $p(\mathbf{x} | \mathbf{z})$)

- The generative story for each \boldsymbol{x}_n , $n = 1, 2, \dots, N$
 - First choose one of the K mixture components as

 $z_n \sim \text{Multinomial}(z_n | \pi)$ (from the prior p(z) over z)

• Suppose $z_n = k$. Now generate x_n from the k-th Gaussian as

 $\mathbf{x}_n \sim \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ (from

(from the data distr. p(x|z))



イロト 不得 トイヨト イヨト 二日

• Multivariate Gaussian in D dimensions

$$p(\boldsymbol{x}|\mu,\boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{x}-\mu)^{\top}\boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\mu)\right)$$

・ロト ・同ト ・ヨト ・ヨト ・ヨー

• Multivariate Gaussian in D dimensions

$$p(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^{\top}\boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right)$$

• Given N i.i.d. observations $\{x_n\}_{n=1}^N$ from this Gaussian

◆□> ◆□> ◆目> ◆目> ◆目> □目 − のへで

• Multivariate Gaussian in D dimensions

$$p(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right)$$

- Given N i.i.d. observations $\{x_n\}_{n=1}^N$ from this Gaussian
 - MLE for the D imes 1 mean $\mu\in\mathbb{R}^D$ $\hat{\mu}=rac{1}{N}\sum_{n=1}^N oldsymbol{x}_n$

◆□ → ◆□ → ◆三 → ◆三 → ● ● ●

• Multivariate Gaussian in D dimensions

$$p(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right)$$

- Given N i.i.d. observations $\{\mathbf{x}_n\}_{n=1}^N$ from this Gaussian
 - MLE for the D imes 1 mean $\mu\in\mathbb{R}^D$ $\hat{\mu}=rac{1}{N}\sum_{n=1}^N oldsymbol{x}_n$
 - MLE for the $D \times D$ p.s.d. covariance matrix Σ

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

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

• Multivariate Gaussian in D dimensions

$$p(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right)$$

- Given N i.i.d. observations $\{x_n\}_{n=1}^N$ from this Gaussian
 - MLE for the D imes 1 mean $\mu\in\mathbb{R}^D$ $\hat{\mu}=rac{1}{N}\sum_{n=1}^N x_n$
 - MLE for the $D \times D$ p.s.d. covariance matrix Σ

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

• Note: The "trace trick" simplifies the derivative calculations

$$\underbrace{\mathbf{x}^{\top}\boldsymbol{\Sigma}^{-1}\mathbf{x}}_{\text{a scalar}} = \operatorname{trace}(\mathbf{x}^{\top}\boldsymbol{\Sigma}^{-1}\mathbf{x}) = \operatorname{trace}(\boldsymbol{\Sigma}^{-1}\mathbf{x}\mathbf{x}^{\top})$$

◆□ > ◆□ > ◆三 > ◆三 > ・三 ● のへの

• Let's do MLE for estimating GMM parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

- Let's do MLE for estimating GMM parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$
- We basically have to estimate K multivariate Gaussians with wts π_1, \ldots, π_K

- Let's do MLE for estimating GMM parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$
- We basically have to estimate K multivariate Gaussians with wts π_1, \ldots, π_K
- The conditional p.d.f. of a data point x if it comes from Gaussian k

 $p(\mathbf{x}|\mathbf{z}=k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

- Let's do MLE for estimating GMM parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$
- We basically have to estimate K multivariate Gaussians with wts π_1, \ldots, π_K
- The conditional p.d.f. of a data point \boldsymbol{x} if it comes from Gaussian k

 $p(\mathbf{x}|\mathbf{z}=k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

• Since z is NOT known, we need to look at the marginal p.d.f. of x

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) = \sum_{k=1}^{K} p(\mathbf{x}, \mathbf{z} = k) = \sum_{k=1}^{K} p(\mathbf{z} = k) p(\mathbf{x} | \mathbf{z} = k) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● ○ ○ ○ ○ ○

- Let's do MLE for estimating GMM parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$
- We basically have to estimate K multivariate Gaussians with wts π_1, \ldots, π_K
- The conditional p.d.f. of a data point \boldsymbol{x} if it comes from Gaussian k

 $p(\mathbf{x}|\mathbf{z}=k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

• Since z is NOT known, we need to look at the marginal p.d.f. of x

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) = \sum_{k=1}^{K} p(\mathbf{x}, \mathbf{z} = k) = \sum_{k=1}^{K} p(\mathbf{z} = k) p(\mathbf{x} | \mathbf{z} = k) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

• Note: Here $p(\mathbf{x})$ means $p(\mathbf{x}|\Theta)$ where $\Theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$

- Let's do MLE for estimating GMM parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$
- We basically have to estimate K multivariate Gaussians with wts π_1, \ldots, π_K
- The conditional p.d.f. of a data point \boldsymbol{x} if it comes from Gaussian k

 $p(\mathbf{x}|\mathbf{z}=k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

• Since z is NOT known, we need to look at the marginal p.d.f. of x

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) = \sum_{k=1}^{K} p(\mathbf{x}, \mathbf{z} = k) = \sum_{k=1}^{K} p(\mathbf{z} = k) p(\mathbf{x} | \mathbf{z} = k) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Note: Here $p(\mathbf{x})$ means $p(\mathbf{x}|\Theta)$ where $\Theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$
- To learn the GMM parameters Θ , we have to do MLE on p(x)

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- Let's do MLE for estimating GMM parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$
- We basically have to estimate K multivariate Gaussians with wts π_1, \ldots, π_K
- The conditional p.d.f. of a data point x if it comes from Gaussian k

 $p(\mathbf{x}|\mathbf{z}=k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

• Since z is NOT known, we need to look at the marginal p.d.f. of x

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) = \sum_{k=1}^{K} p(\mathbf{x}, \mathbf{z} = k) = \sum_{k=1}^{K} p(\mathbf{z} = k) p(\mathbf{x} | \mathbf{z} = k) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Note: Here $p(\mathbf{x})$ means $p(\mathbf{x}|\Theta)$ where $\Theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$
- To learn the GMM parameters Θ , we have to do MLE on $p(\mathbf{x})$
- In general, it is not an easy problem due to the difficult form of p(x) (for "why", see the next slide)

Machine Learning (CS771A)

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

• Given N observations x_1, \ldots, x_N , the log-likelihood will be

$$\mathcal{L} = \log \prod_{n=1}^{N} p(\mathbf{x}_n) = \sum_{n=1}^{N} \log p(\mathbf{x}_n) = \sum_{n=1}^{N} \log \underbrace{\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{\text{params get coupled}}$$

Machine Learning (CS771A)

・ロト ・同ト ・ヨト ・ヨト ・ヨー

• Given N observations x_1, \ldots, x_N , the log-likelihood will be

$$\mathcal{L} = \log \prod_{n=1}^{N} p(\mathbf{x}_n) = \sum_{n=1}^{N} \log p(\mathbf{x}_n) = \sum_{n=1}^{N} \log \underbrace{\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{\text{params get coupled!}}$$

Due to the coupling of parameters, MLE by simply taking derivatives of *L* and setting to zero won't give a closed form solution of Θ = {π_k, μ_k, Σ_k}^K_{k=1}

◆□ → ◆□ → ◆三 → ◆三 → ● ● ●

• Given N observations x_1, \ldots, x_N , the log-likelihood will be

$$\mathcal{L} = \log \prod_{n=1}^{N} p(\boldsymbol{x}_n) = \sum_{n=1}^{N} \log p(\boldsymbol{x}_n) = \sum_{n=1}^{N} \log \underbrace{\sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{\text{params get coupled}}$$

- Due to the coupling of parameters, MLE by simply taking derivatives of *L* and setting to zero won't give a closed form solution of Θ = {π_k, μ_k, Σ_k}^K_{k=1}
 - Gradient based iterative methods can be used

◆□ → ◆□ → ◆三 → ◆三 → ● ● ●

• Given N observations x_1, \ldots, x_N , the log-likelihood will be

$$\mathcal{L} = \log \prod_{n=1}^{N} p(\mathbf{x}_n) = \sum_{n=1}^{N} \log p(\mathbf{x}_n) = \sum_{n=1}^{N} \log \underbrace{\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{\text{params get coupled}}$$

- Due to the coupling of parameters, MLE by simply taking derivatives of *L* and setting to zero won't give a closed form solution of Θ = {π_k, μ_k, Σ_k}^K_{k=1}
 - Gradient based iterative methods can be used
 - However, we will use Expectation Maximization (EM) a more general way of solving such problems (i.e., parameter estimation with latent variables)

(ロ) (同) (三) (三) (三) (000)

• Given N observations x_1, \ldots, x_N , the log-likelihood will be

$$\mathcal{L} = \log \prod_{n=1}^{N} p(\mathbf{x}_n) = \sum_{n=1}^{N} \log p(\mathbf{x}_n) = \sum_{n=1}^{N} \log \underbrace{\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{\text{params get coupled}}$$

- Due to the coupling of parameters, MLE by simply taking derivatives of *L* and setting to zero won't give a closed form solution of Θ = {π_k, μ_k, Σ_k}^K_{k=1}
 - Gradient based iterative methods can be used
 - However, we will use Expectation Maximization (EM) a more general way of solving such problems (i.e., parameter estimation with latent variables)
- Note: For problems where z_n is continuous or comes from a combinatorially large space (e.g., it's a binary vector), doing MLE will be even harder!

$$\mathcal{L} = \sum_{n=1}^{N} \log \underbrace{\int_{z_n} p(x_n | z_n) p(z_n) dz_n}_{\sum z_n}$$

Ouch! Intractable integral!!!

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のへの

• Given N observations x_1, \ldots, x_N , the log-likelihood will be

$$\mathcal{L} = \log \prod_{n=1}^{N} p(\mathbf{x}_n) = \sum_{n=1}^{N} \log p(\mathbf{x}_n) = \sum_{n=1}^{N} \log \underbrace{\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{\text{params get coupled}}$$

- Due to the coupling of parameters, MLE by simply taking derivatives of \mathcal{L} and setting to zero won't give a closed form solution of $\Theta = \{\pi_k, \mu_k, \boldsymbol{\Sigma}_k\}_{k=1}^{K}$
 - Gradient based iterative methods can be used
 - However, we will use Expectation Maximization (EM) a more general way of solving such problems (i.e., parameter estimation with latent variables)
- Note: For problems where z_n is continuous or comes from a combinatorially large space (e.g., it's a binary vector), doing MLE will be even harder!

$$\mathcal{L} = \sum_{n=1}^{N} \log \underbrace{\int_{z_n} p(x_n | z_n) p(z_n) dz_n}_{\sum z_n}$$

Ouch! Intractable integral!!!

In such cases, something like EM becomes even more important

・ロト ・ 同 ト ・ 三 ト ・ 三 ・ つへの
• MLE for GMM will be simplified if we "knew" the z_n for each x_n

<ロ> (四) (四) (三) (三) (三) (三)

- MLE for GMM will be simplified if we "knew" the z_n for each x_n
- Reason: If z_n is known, the summation over z_n isn't required

$$\sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

-

イロト イポト イヨト イヨト

- MLE for GMM will be simplified if we "knew" the z_n for each x_n
- Reason: If z_n is known, the summation over z_n isn't required

$$\sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

• With z_n "known", we can try doing MLE on $p(x_n, z_n)$, instead of on $p(x_n)$

- MLE for GMM will be simplified if we "knew" the z_n for each x_n
- Reason: If z_n is known, the summation over z_n isn't required

$$\sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- With z_n "known", we can try doing MLE on $p(x_n, z_n)$, instead of on $p(x_n)$
 - $p(x_n, z_n)$ is known as "complete data likelihood" (z_n makes x_n "complete")

- MLE for GMM will be simplified if we "knew" the z_n for each x_n
- Reason: If z_n is known, the summation over z_n isn't required

$$\sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- With z_n "known", we can try doing MLE on $p(x_n, z_n)$, instead of on $p(x_n)$
 - $p(x_n, z_n)$ is known as "complete data likelihood" (z_n makes x_n "complete")
 - $p(x_n)$ is known as "incomplete data likelihood"

<ロ> (四) (四) (三) (三) (三) (三)

- MLE for GMM will be simplified if we "knew" the z_n for each x_n
- Reason: If z_n is known, the summation over z_n isn't required

$$\sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- With z_n "known", we can try doing MLE on $p(x_n, z_n)$, instead of on $p(x_n)$
 - $p(x_n, z_n)$ is known as "complete data likelihood" (z_n makes x_n "complete")
 - $p(x_n)$ is known as "incomplete data likelihood"
- Since we don't know the "true" z_n , we will have to rely on a "guess" for z_n

- MLE for GMM will be simplified if we "knew" the z_n for each x_n
- Reason: If z_n is known, the summation over z_n isn't required

$$\sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- With z_n "known", we can try doing MLE on $p(x_n, z_n)$, instead of on $p(x_n)$
 - $p(x_n, z_n)$ is known as "complete data likelihood" (z_n makes x_n "complete")
 - $p(x_n)$ is known as "incomplete data likelihood"
- Since we don't know the "true" z_n , we will have to rely on a "guess" for z_n
 - This guess for z_n will be based on the current values of params Θ

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ ○○

- MLE for GMM will be simplified if we "knew" the z_n for each x_n
- Reason: If z_n is known, the summation over z_n isn't required

$$\sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- With z_n "known", we can try doing MLE on $p(x_n, z_n)$, instead of on $p(x_n)$
 - $p(x_n, z_n)$ is known as "complete data likelihood" (z_n makes x_n "complete")
 - $p(x_n)$ is known as "incomplete data likelihood"
- Since we don't know the "true" z_n , we will have to rely on a "guess" for z_n
 - This guess for z_n will be based on the current values of params Θ
 - Can do MLE on $p(x_n, z_n)$ to re-estimate Θ using these guesses, and repeat

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ ○○

- MLE for GMM will be simplified if we "knew" the z_n for each x_n
- Reason: If z_n is known, the summation over z_n isn't required

$$\sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- With z_n "known", we can try doing MLE on $p(x_n, z_n)$, instead of on $p(x_n)$
 - $p(x_n, z_n)$ is known as "complete data likelihood" (z_n makes x_n "complete")
 - $p(x_n)$ is known as "incomplete data likelihood"
- Since we don't know the "true" z_n , we will have to rely on a "guess" for z_n
 - This guess for z_n will be based on the current values of params Θ
 - Can do MLE on $p(x_n, z_n)$ to re-estimate Θ using these guesses, and repeat
 - A more formal view of this iterative procedure is given by the Expectation Maximization (EM) algorithm (next lecture)

• The complete data log-likelihood over the *N* obs.

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

(note that, for each n, only one z_{nk} will be 1)

イロン 不得 とくほ とくほ とうほう

• The complete data log-likelihood over the *N* obs.

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

(note that, for each n, only one z_{nk} will be 1)

• The above gets further simplified to

$$\sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \log p(\boldsymbol{z}_{n} = k) p(\boldsymbol{x}_{n} | \boldsymbol{z}_{n} = k) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} [\log \pi_{k} + \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})]$$

• The complete data log-likelihood over the *N* obs.

$$\sum_{n=1}^{N} \log p(\mathbf{x}_n, \mathbf{z}_n) = \sum_{n=1}^{N} \log p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n) = \sum_{n=1}^{N} \log \prod_{k=1}^{K} [p(\mathbf{x}_n | \mathbf{z}_n = k) p(\mathbf{z}_n = k)]^{z_{nk}}$$

• The above gets further simplified to

$$\sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \log p(\boldsymbol{z}_{n} = k) p(\boldsymbol{x}_{n} | \boldsymbol{z}_{n} = k) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} [\log \pi_{k} + \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})]$$

• If we know value of each z_{nk} deterministically, we can plug these in and do MLE on the above objective (which has a simple and separable structure)

⁽note that, for each n, only one z_{nk} will be 1)

• The complete data log-likelihood over the *N* obs.

$$\sum_{n=1}^{N} \log p(\mathbf{x}_n, \mathbf{z}_n) = \sum_{n=1}^{N} \log p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n) = \sum_{n=1}^{N} \log \prod_{k=1}^{K} [p(\mathbf{x}_n | \mathbf{z}_n = k) p(\mathbf{z}_n = k)]^{z_{nk}}$$

• The above gets further simplified to

$$\sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \log p(\boldsymbol{z}_{n} = k) p(\boldsymbol{x}_{n} | \boldsymbol{z}_{n} = k) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} [\log \pi_{k} + \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})]$$

- If we know value of each z_{nk} deterministically, we can plug these in and do MLE on the above objective (which has a simple and separable structure)
- What if we don't have a *deterministic* guess for z_{nk} ?

イロン 不得 とくほど 不良 とうほう

⁽note that, for each n, only one z_{nk} will be 1)

• The complete data log-likelihood over the *N* obs.

$$\sum_{n=1}^{N} \log p(\mathbf{x}_n, \mathbf{z}_n) = \sum_{n=1}^{N} \log p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n) = \sum_{n=1}^{N} \log \prod_{k=1}^{K} [p(\mathbf{x}_n | \mathbf{z}_n = k) p(\mathbf{z}_n = k)]^{z_{nk}}$$

• The above gets further simplified to

$$\sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \log p(\boldsymbol{z}_{n} = k) p(\boldsymbol{x}_{n} | \boldsymbol{z}_{n} = k) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} [\log \pi_{k} + \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})]$$

- If we know value of each z_{nk} deterministically, we can plug these in and do MLE on the above objective (which has a simple and separable structure)
- What if we don't have a *deterministic* guess for z_{nk} ?
 - In such cases, we can use the posterior expectations of the z_{nk}'s (which are basically posterior probabilities of cluster assignments of points to clusters)

Machine Learning (CS771A)

⁽note that, for each n, only one z_{nk} will be 1)

• Posterior expectations $\mathbb{E}[z_{nk}]$ can be computed using current estimates of Θ

3

• Posterior expectations $\mathbb{E}[z_{nk}]$ can be computed using current estimates of Θ

 $\mathbb{E}[z_{nk}] = 0 \times p(z_{nk} = 0 | \mathbf{x}_n) + 1 \times p(z_{nk} = 1 | \mathbf{x}_n)$

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● ○ ○ ○ ○ ○

• Posterior expectations $\mathbb{E}[z_{nk}]$ can be computed using current estimates of Θ

$$\mathbb{E}[z_{nk}] = 0 \times p(z_{nk} = 0 | \mathbf{x}_n) + 1 \times p(z_{nk} = 1 | \mathbf{x}_n)$$
$$= p(z_{nk} = 1 | \mathbf{x}_n)$$

3

• Posterior expectations $\mathbb{E}[z_{nk}]$ can be computed using current estimates of Θ

$$\mathbb{E}[z_{nk}] = 0 \times p(z_{nk} = 0|\mathbf{x}_n) + 1 \times p(z_{nk} = 1|\mathbf{x}_n)$$
$$= p(z_{nk} = 1|\mathbf{x}_n)$$
$$\propto p(z_{nk} = 1)p(\mathbf{x}_n|z_{nk} = 1)$$
(Bayes Rule)

- E

• Posterior expectations $\mathbb{E}[z_{nk}]$ can be computed using current estimates of Θ

$$\begin{split} \mathbb{E}[z_{nk}] &= 0 \times p(z_{nk} = 0 | x_n) + 1 \times p(z_{nk} = 1 | x_n) \\ &= p(z_{nk} = 1 | x_n) \\ &\propto p(z_{nk} = 1) p(x_n | z_{nk} = 1) \quad (\text{Bayes Rule}) \end{split}$$
Thus $\mathbb{E}[z_{nk}] \propto \pi_k \mathcal{N}(x_n | \mu_k, \mathbf{\Sigma}_k) \quad (\text{Posterior prob. of } x_n \text{ belonging to cluster } k)$

3

• Posterior expectations $\mathbb{E}[z_{nk}]$ can be computed using current estimates of Θ

$$\begin{split} \mathbb{E}[z_{nk}] &= 0 \times p(z_{nk} = 0|\mathbf{x}_n) + 1 \times p(z_{nk} = 1|\mathbf{x}_n) \\ &= p(z_{nk} = 1|\mathbf{x}_n) \\ &\propto p(z_{nk} = 1)p(\mathbf{x}_n|z_{nk} = 1) \quad (\text{Bayes Rule}) \end{split}$$
Thus $\mathbb{E}[z_{nk}] \propto \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (\text{Posterior prob. of } \boldsymbol{x}_n \text{ belonging to cluster } k)$

• The final expression of $\mathbb{E}[z_{nk}]$ makes intuitive sense

3

• Posterior expectations $\mathbb{E}[z_{nk}]$ can be computed using current estimates of Θ

$$\mathbb{E}[z_{nk}] = 0 \times p(z_{nk} = 0|\mathbf{x}_n) + 1 \times p(z_{nk} = 1|\mathbf{x}_n)$$

$$= p(z_{nk} = 1|\mathbf{x}_n)$$

$$\propto p(z_{nk} = 1)p(\mathbf{x}_n|z_{nk} = 1) \quad (Bayes Rule)$$
Thus $\mathbb{E}[z_{nk}] \propto \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (Posterior prob. of \boldsymbol{x}_n \text{ belonging to cluster } k)$

- The final expression of $\mathbb{E}[z_{nk}]$ makes intuitive sense
- Note: We can finally normalize $\mathbb{E}[z_{nk}]$ as $\mathbb{E}[z_{nk}] = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)}$ since $\sum_{k=1}^K \mathbb{E}[z_{nk}] = 1$

・ロト ・ 日 ・ モ ト ・ ヨ ・ うへの

• Posterior expectations $\mathbb{E}[z_{nk}]$ can be computed using current estimates of Θ

$$\mathbb{E}[z_{nk}] = 0 \times p(z_{nk} = 0|\mathbf{x}_n) + 1 \times p(z_{nk} = 1|\mathbf{x}_n)$$

$$= p(z_{nk} = 1|\mathbf{x}_n)$$

$$\propto p(z_{nk} = 1)p(\mathbf{x}_n|z_{nk} = 1) \quad (Bayes Rule)$$
Thus $\mathbb{E}[z_{nk}] \propto \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (Posterior prob. of \boldsymbol{x}_n \text{ belonging to cluster } k)$

- The final expression of $\mathbb{E}[z_{nk}]$ makes intuitive sense
- Note: We can finally normalize $\mathbb{E}[z_{nk}]$ as $\mathbb{E}[z_{nk}] = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)}$ since $\sum_{k=1}^K \mathbb{E}[z_{nk}] = 1$
- Given $\mathbb{E}[z_{nk}]$, we can now define expected complete data log-lik.

• Posterior expectations $\mathbb{E}[z_{nk}]$ can be computed using current estimates of Θ

$$\mathbb{E}[z_{nk}] = 0 \times p(z_{nk} = 0|\mathbf{x}_n) + 1 \times p(z_{nk} = 1|\mathbf{x}_n)$$

$$= p(z_{nk} = 1|\mathbf{x}_n)$$

$$\propto p(z_{nk} = 1)p(\mathbf{x}_n|z_{nk} = 1) \quad (Bayes Rule)$$
Thus $\mathbb{E}[z_{nk}] \propto \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (Posterior prob. of \boldsymbol{x}_n \text{ belonging to cluster } k)$

- The final expression of $\mathbb{E}[z_{nk}]$ makes intuitive sense
- Note: We can finally normalize $\mathbb{E}[z_{nk}]$ as $\mathbb{E}[z_{nk}] = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)}$ since $\sum_{k=1}^K \mathbb{E}[z_{nk}] = 1$
- Given $\mathbb{E}[z_{nk}]$, we can now define expected complete data log-lik.

$$\mathcal{L} = \sum_{n=1}^{N} \sum_{k=1}^{K} \mathbb{E}[\boldsymbol{z}_{nk}] [\log \pi_k + \log \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

.. and do MLE for the parameters $\boldsymbol{\Theta}$ using this as the objective function

Machine Learning (CS771A)

・ロト ・ 同 ト ・ 三 ト ・ 三 ・ つへの

• Given
$$\mathbb{E}[z_{nk}] = \gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell=1}^{K} \pi_\ell \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)}$$
, the expected complete data log-lik.

3

• Given
$$\mathbb{E}[z_{nk}] = \gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)}$$
, the expected complete data log-lik.
$$\mathcal{L} = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

3

• Given
$$\mathbb{E}[z_{nk}] = \gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)}$$
, the expected complete data log-lik.
$$\mathcal{L} = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

• Taking derivatives w.r.t. $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$, $orall k=1,\ldots,K$

<ロ> (四) (四) (三) (三) (三) (三)

• Given
$$\mathbb{E}[z_{nk}] = \gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)}$$
, the expected complete data log-lik.
$$\mathcal{L} = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

• Taking derivatives w.r.t. $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$, $orall k=1,\ldots,K$

 $\frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}_{k}} = \frac{\partial}{\partial \boldsymbol{\mu}_{k}} \sum_{n=1}^{N} \gamma_{nk} \log \mathcal{N}(\mathbf{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) = 0 \qquad \text{(note: constants w.r.t. } \boldsymbol{\mu}_{k} \text{ can be ignored)}$

• Given
$$\mathbb{E}[z_{nk}] = \gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)}$$
, the expected complete data log-lik.
$$\mathcal{L} = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

• Taking derivatives w.r.t. $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$, $orall k=1,\ldots,K$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}_{k}} = \frac{\partial}{\partial \boldsymbol{\mu}_{k}} \sum_{n=1}^{N} \gamma_{nk} \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) = 0 \qquad \text{(note: constants w.r.t. } \boldsymbol{\mu}_{k} \text{ can be ignored)}$$
$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\Sigma}_{k}} = \frac{\partial}{\partial \boldsymbol{\Sigma}_{k}} \sum_{n=1}^{N} \gamma_{nk} \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) = 0 \qquad \text{(note: constants w.r.t. } \boldsymbol{\Sigma}_{k} \text{ can be ignored)}$$

・ロト ・御 ト ・ヨト ・ヨト 三日 -

• Given
$$\mathbb{E}[z_{nk}] = \gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)}$$
, the expected complete data log-lik.
$$\mathcal{L} = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

• Taking derivatives w.r.t. $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$, $orall k=1,\ldots,K$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}_{k}} = \frac{\partial}{\partial \boldsymbol{\mu}_{k}} \sum_{n=1}^{N} \gamma_{nk} \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) = 0 \quad \text{(note: constants w.r.t. } \boldsymbol{\mu}_{k} \text{ can be ignored)}$$
$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\Sigma}_{k}} = \frac{\partial}{\partial \boldsymbol{\Sigma}_{k}} \sum_{n=1}^{N} \gamma_{nk} \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) = 0 \quad \text{(note: constants w.r.t. } \boldsymbol{\Sigma}_{k} \text{ can be ignored)}$$

• For each k, it's a "weighted" version of the MLE problem for the multivariate Gaussian $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, given observations $\{\mathbf{x}_n\}_{n=1}^N$ with weights $\{\gamma_{nk}\}_{n=1}^N$

・ロト ・御 ト ・ヨト ・ヨト 三日 -

• Given
$$\mathbb{E}[z_{nk}] = \gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)}$$
, the expected complete data log-lik.
$$\mathcal{L} = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

• Taking derivatives w.r.t. $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$, $orall k=1,\ldots,K$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}_{k}} = \frac{\partial}{\partial \boldsymbol{\mu}_{k}} \sum_{n=1}^{N} \gamma_{nk} \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) = 0 \quad \text{(note: constants w.r.t. } \boldsymbol{\mu}_{k} \text{ can be ignored)}$$
$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\Sigma}_{k}} = \frac{\partial}{\partial \boldsymbol{\Sigma}_{k}} \sum_{n=1}^{N} \gamma_{nk} \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) = 0 \quad \text{(note: constants w.r.t. } \boldsymbol{\Sigma}_{k} \text{ can be ignored)}$$

- For each k, it's a "weighted" version of the MLE problem for the multivariate Gaussian $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, given observations $\{\mathbf{x}_n\}_{n=1}^N$ with weights $\{\gamma_{nk}\}_{n=1}^N$
- Can also solve for π_k likewise (subject to contraint $\sum_{k=1}^{K} \pi_k = 1$)

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ ○○

• Given
$$\mathbb{E}[z_{nk}] = \gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)}$$
, the expected complete data log-lik.
$$\mathcal{L} = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

• Taking derivatives w.r.t. $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$, $orall k=1,\ldots,K$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}_{k}} = \frac{\partial}{\partial \boldsymbol{\mu}_{k}} \sum_{n=1}^{N} \gamma_{nk} \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) = 0 \qquad \text{(note: constants w.r.t. } \boldsymbol{\mu}_{k} \text{ can be ignored)}$$
$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\Sigma}_{k}} = \frac{\partial}{\partial \boldsymbol{\Sigma}_{k}} \sum_{n=1}^{N} \gamma_{nk} \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) = 0 \qquad \text{(note: constants w.r.t. } \boldsymbol{\Sigma}_{k} \text{ can be ignored)}$$

- For each k, it's a "weighted" version of the MLE problem for the multivariate Gaussian $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, given observations $\{\mathbf{x}_n\}_{n=1}^N$ with weights $\{\gamma_{nk}\}_{n=1}^N$
- Can also solve for π_k likewise (subject to contraint $\sum_{k=1}^{K} \pi_k = 1$)
- Derivations are a bit tedious (but straightforward). I will provide a note.

Machine Learning (CS771A)

GMM Parameter Update Equations

• The final expressions for updates of $\{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$

$$\boldsymbol{\mu}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma_{nk} \boldsymbol{x}_{n}$$
$$\boldsymbol{\Sigma}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma_{nk} (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k}) (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k})^{\top}$$
$$\boldsymbol{\pi}_{k} = \frac{N_{k}}{N}$$

3

GMM Parameter Update Equations

• The final expressions for updates of $\{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$

$$\boldsymbol{\mu}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma_{nk} \boldsymbol{x}_{n}$$
$$\boldsymbol{\Sigma}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma_{nk} (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k}) (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k})^{\top}$$
$$\pi_{k} = \frac{N_{k}}{N}$$

• Note: $N_k = \sum_{n=1}^N \gamma_{nk}$ is the effective num. of pts. assigned to Gaussian k

• Update equations make intuitive sense

GMM Parameter Update Equations

• The final expressions for updates of $\{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$

$$\boldsymbol{\mu}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma_{nk} \boldsymbol{x}_{n}$$
$$\boldsymbol{\Sigma}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma_{nk} (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k}) (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k})^{\top}$$
$$\boldsymbol{\pi}_{k} = \frac{N_{k}}{N}$$

• Note: $N_k = \sum_{n=1}^N \gamma_{nk}$ is the effective num. of pts. assigned to Gaussian k

- Update equations make intuitive sense
- Also note that each point x_n contributes to each $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ but fractionally (based on the values of γ_{nk})

Machine Learning (CS771A)

The Full Algorithm for Learning GMM

• Initialize the parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$ randomly, or using K-means

The Full Algorithm for Learning GMM

- Initialize the parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$ randomly, or using K-means
- Iterate until convergence (e.g., when $\log p(\mathbf{x}|\Theta)$ ceases to increase)
- Initialize the parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$ randomly, or using K-means
- Iterate until convergence (e.g., when $\log p(\mathbf{x}|\Theta)$ ceases to increase)
 - Given Θ , compute each expectation z_{nk} (post. prob. of $z_{nk} = 1$), $\forall n, k$

イロン 不通 とうほう イロン しゅう

- Initialize the parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$ randomly, or using K-means
- Iterate until convergence (e.g., when $\log p(\mathbf{x}|\Theta)$ ceases to increase)
 - Given Θ , compute each expectation z_{nk} (post. prob. of $z_{nk} = 1$), $\forall n, k$

$$\gamma_{nk} = \mathbb{E}[z_{nk}] \propto \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad \text{(and re-normalize s.t. } \sum_{k=1}^{K} \gamma_{nk} = 1)$$

- Initialize the parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$ randomly, or using K-means
- Iterate until convergence (e.g., when $\log p(\mathbf{x}|\Theta)$ ceases to increase)
 - Given Θ , compute each expectation z_{nk} (post. prob. of $z_{nk} = 1$), $\forall n, k$

$$\begin{split} \gamma_{nk} &= \mathbb{E}[z_{nk}] \propto \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \qquad (\text{and re-normalize s.t. } \sum_{k=1}^K \gamma_{nk} = 1) \\ \bullet \text{ Given } \gamma_{nk} &= \mathbb{E}[z_{nk}] \text{ and } N_k = \sum_{n=1}^N \gamma_{nk}, \text{ update } \Theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K \text{ as } \end{split}$$

イロン 不通 とうほう イロン しゅう

- Initialize the parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$ randomly, or using K-means
- Iterate until convergence (e.g., when $\log p(\mathbf{x}|\Theta)$ ceases to increase)
 - Given Θ , compute each expectation z_{nk} (post. prob. of $z_{nk} = 1$), $\forall n, k$

$$\gamma_{nk} = \mathbb{E}[z_{nk}] \propto \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (\text{and re-normalize s.t. } \sum_{k=1}^n \gamma_{nk} = 1)$$

• Given $\gamma_{nk} = \mathbb{E}[z_{nk}]$ and $N_k = \sum_{n=1}^N \gamma_{nk}$, update $\Theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$ as
 $\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n$
 $\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$
 $\pi_k = \frac{N_k}{N}$

イロン 不通 とうほう イロン しゅう

- Initialize the parameters $\Theta = \{\pi_k, \mu_k, \mathbf{\Sigma}_k\}_{k=1}^K$ randomly, or using K-means
- Iterate until convergence (e.g., when $\log p(\mathbf{x}|\Theta)$ ceases to increase)
 - Given Θ , compute each expectation z_{nk} (post. prob. of $z_{nk} = 1$), $\forall n, k$

$$\gamma_{nk} = \mathbb{E}[z_{nk}] \propto \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (\text{and re-normalize s.t. } \sum_{k=1}^n \gamma_{nk} = 1)$$

• Given $\gamma_{nk} = \mathbb{E}[z_{nk}]$ and $N_k = \sum_{n=1}^N \gamma_{nk}$, update $\Theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$ as
 $\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n$
 $\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$
 $\pi_k = \frac{N_k}{N}$

(It's basically an Expectation Maximization (EM) algorithm for learning GMM. We will look at EM more formally in the next class.)

Machine Learning (CS771A)

v

Illustration of GMM Clustering



Notice the "mixed" colored points in the overlapping regions in the final clustering

Illustration of GMM Clustering



Notice the "mixed" colored points in the overlapping regions in the final clustering Also check out this demo of GMM: https://www.youtube.com/watch?v=B36fzChfyGU

GMM vs K-means

For the GMM clustering (rightmost figure), the most probable cluster for each point has been labeled



Pic courtesy: https://en.wikipedia.org/wiki/Expectation-maximization_algorithm/

GMM vs K-means

For the GMM clustering (rightmost figure), the most probable cluster for each point has been labeled



Note that K-means, unlike GMM, tends to learn equi-sized clusters.

Pic courtesy: https://en.wikipedia.org/wiki/Expectation-maximization_algorithm/

GMM vs K-means

For the GMM clustering (rightmost figure), the most probable cluster for each point has been labeled



Note that K-means, unlike GMM, tends to learn equi-sized clusters.

GMM with $\Sigma_k = I$ and $\pi_k = 1/K$, and soft assignments converted to hard assign. (setting the largest prob. to 1, rest to 0), is equivalent to K-means.

Pic courtesy: https://en.wikipedia.org/wiki/Expectation-maximization_algorithm/

• Other types of component distributions can be used

3

イロト イポト イヨト イヨト

- Other types of component distributions can be used
- Sequence data models, such as HMM, can also be seen as mixture models

 $p(\mathbf{z}_n|\phi) = p(\mathbf{z}_n|\mathbf{z}_{n-1})$ (Current cluster/state depends on previous)

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● ○ ○ ○ ○ ○

- Other types of component distributions can be used
- Sequence data models, such as HMM, can also be seen as mixture models

 $p(\mathbf{z}_n|\phi) = p(\mathbf{z}_n|\mathbf{z}_{n-1})$ (Current cluster/state depends on previous)

• Also used in supervised learning problems (mixture of experts)

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ ○○

- Other types of component distributions can be used
- Sequence data models, such as HMM, can also be seen as mixture models

 $p(\mathbf{z}_n|\phi) = p(\mathbf{z}_n|\mathbf{z}_{n-1})$ (Current cluster/state depends on previous)

- Also used in supervised learning problems (mixture of experts)
 - For data x_n , first choose one of K experts, and then use that expert's predictive model for $p(y_n|x_n)$

- Other types of component distributions can be used
- Sequence data models, such as HMM, can also be seen as mixture models

 $p(\mathbf{z}_n|\phi) = p(\mathbf{z}_n|\mathbf{z}_{n-1})$ (Current cluster/state depends on previous)

- Also used in supervised learning problems (mixture of experts)
 - For data x_n , first choose one of K experts, and then use that expert's predictive model for $p(y_n|x_n)$
 - Experts and points-to-experts assignments can be learned iteratively as in GMM

・ロト ・ 日 ・ モ ト ・ ヨ ・ うへの

- Other types of component distributions can be used
- Sequence data models, such as HMM, can also be seen as mixture models

 $p(\mathbf{z}_n|\phi) = p(\mathbf{z}_n|\mathbf{z}_{n-1})$ (Current cluster/state depends on previous)

- Also used in supervised learning problems (mixture of experts)
 - For data x_n , first choose one of K experts, and then use that expert's predictive model for $p(y_n|x_n)$
 - Experts and points-to-experts assignments can be learned iteratively as in GMM
- Can be used for performing generative classification (e.g., naïve Bayes)

・ロト ・ 日 ・ モ ト ・ ヨ ・ うへの

- Other types of component distributions can be used
- Sequence data models, such as HMM, can also be seen as mixture models

 $p(\mathbf{z}_n|\phi) = p(\mathbf{z}_n|\mathbf{z}_{n-1})$ (Current cluster/state depends on previous)

- Also used in supervised learning problems (mixture of experts)
 - For data x_n , first choose one of K experts, and then use that expert's predictive model for $p(y_n|x_n)$
 - Experts and points-to-experts assignments can be learned iteratively as in GMM
- Can be used for performing generative classification (e.g., naïve Bayes)

 $p(y_n = k | \boldsymbol{x}_n) \propto p(y_n = k) p(\boldsymbol{x}_n | y_n = k)$ (cluster ids are the known classes)

- Other types of component distributions can be used
- Sequence data models, such as HMM, can also be seen as mixture models

 $p(\mathbf{z}_n|\phi) = p(\mathbf{z}_n|\mathbf{z}_{n-1})$ (Current cluster/state depends on previous)

- Also used in supervised learning problems (mixture of experts)
 - For data x_n , first choose one of K experts, and then use that expert's predictive model for $p(y_n|x_n)$
 - Experts and points-to-experts assignments can be learned iteratively as in GMM
- Can be used for performing generative classification (e.g., naïve Bayes)

 $p(y_n = k | \boldsymbol{x}_n) \propto p(y_n = k) p(\boldsymbol{x}_n | y_n = k)$ (cluster ids are the known classes)

p(y = k) and p(x|z = k) can be efficiently estimated using training data

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ □臣 = のへで

- Other types of component distributions can be used
- Sequence data models, such as HMM, can also be seen as mixture models

 $p(\mathbf{z}_n|\phi) = p(\mathbf{z}_n|\mathbf{z}_{n-1})$ (Current cluster/state depends on previous)

- Also used in supervised learning problems (mixture of experts)
 - For data x_n , first choose one of K experts, and then use that expert's predictive model for $p(y_n|x_n)$
 - Experts and points-to-experts assignments can be learned iteratively as in GMM
- Can be used for performing generative classification (e.g., naïve Bayes)

 $p(y_n = k | \boldsymbol{x}_n) \propto p(y_n = k) p(\boldsymbol{x}_n | y_n = k)$ (cluster ids are the known classes)

p(y = k) and p(x|z = k) can be efficiently estimated using training data

Number of clusters (K) in a mixture model can be learned from data using nonparametric Bayesian methods (e.g., "infinite" mixture models)

Machine Learning (CS771A)

- The general Expectation Maximization (EM) algorithm
- Generative models for dimensionality reduction
 - Factor Analysis and Probabilistic PCA (and extensions)
 - EM based parameter estimation for these models

-

イロト イポト イヨト イヨト