# CS182: Introduction to Machine Learning –Expectation-Maximization (EM) algorithm and Gaussian Mixture Models (GMM)

Yujiao Shi

SIST, ShanghaiTech

Spring, 2025

上海科技大学
ShanghaiTech University

- Consider a generative model with joint distr. $p(\mathbf{X}, \mathbf{Z} | \Theta) = \prod_{n=1}^{N} p(\mathbf{x}_n, \mathbf{z}_n)$
  - Observed data: $\mathbf{X} = \{\mathbf{x}_n\}_{n=1}^{N}$
  - Latent variables: $\mathbf{Z} = \{\mathbf{z}_n\}_{n=1}^{N}$. All the model parameters: $\Theta$

上海科技大学

ShanghaiTech University

- Consider a generative model with joint distr. $p(\mathbf{X}, \mathbf{Z}|\Theta) = \prod_{n=1}^{N} p(\mathbf{x}_n, \mathbf{z}_n)$

  - Observed data: $\mathbf{X} = \{\mathbf{x}_n\}_{n=1}^{N}$

  - Latent variables: $\mathbf{Z} = \{\mathbf{z}_n\}_{n=1}^{N}$. All the model parameters: $\Theta$

- Goal: Estimate the model parameters $\Theta$ via MLE (or MAP)

$$\hat{\Theta} = \arg\max_{\Theta} \log p(\mathbf{X}|\Theta) \quad = \quad \arg\max_{\Theta} \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) \qquad \text{(when } \mathbf{Z} \text{ is discrete)}$$

$$= \quad \arg\max_{\Theta} \log \int_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) d\mathbf{Z} \quad \text{(when } \mathbf{Z} \text{ is continuous)}$$

**3**

- Consider the 'incomplete" data log likelihood

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

- Consider the 'incomplete'' data log likelihood

$$\log p(\mathbf{X}|\Theta) \quad = \quad \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

# Parameter Estimation with Latent Variables

- Consider the 'incomplete' data log likelihood

$$\log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})}$$

上海科技大学

- Consider the 'incomplete" data log likelihood

$$\log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta}$$

# Parameter Estimation with Latent Variables

- Consider the 'incomplete' data log likelihood

$$\log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

**8**

- Consider the 'incomplete" data log likelihood

$$\log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X},\mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X},\mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X},\mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X},\mathbf{Z}|\Theta) - \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X},\mathbf{Z}|\Theta) + \text{const.}$$

- If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X},\Theta)$, the above inequality becomes equality

**9**

上海科技大学

- Consider the 'incomplete' data log likelihood

$$\log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \geq \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

- If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, the above inequality becomes equality

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})}$$

**10**

- Consider the 'incomplete' data log likelihood

$$\log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \geq \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

- If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, the above inequality becomes equality

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log \frac{p(\mathbf{Z}|\mathbf{X}, \Theta) p(\mathbf{X}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)}$$

**11**

上海科技大学

- Consider the 'incomplete' data log likelihood

$$\log p(\mathbf{X}|\Theta) \quad = \quad \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \quad \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \quad \geq \quad \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

- If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, the above inequality becomes equality

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad = \quad \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log \frac{p(\mathbf{Z}|\mathbf{X}, \Theta) p(\mathbf{X}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}|\Theta)$$

上海科技大学

- Consider the 'incomplete' data log likelihood

$$\log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \geq \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

- If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, the above inequality becomes equality

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log \frac{p(\mathbf{Z}|\mathbf{X}, \Theta) p(\mathbf{X}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}|\Theta)$$

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

- Consider the 'incomplete' data log likelihood

$$\log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \geq \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

- If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, the above inequality becomes equality

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log \frac{p(\mathbf{Z}|\mathbf{X}, \Theta) p(\mathbf{X}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}|\Theta)$$

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

# Parameter Estimation with Latent Variables

- Consider the 'incomplete" data log likelihood

$$\log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \geq \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

- If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, the above inequality becomes equality

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log \frac{p(\mathbf{Z}|\mathbf{X}, \Theta) p(\mathbf{X}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}|\Theta)$$

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

- Thus for $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, we have

上海科技大学

- Consider the 'incomplete" data log likelihood

$$\log p(\mathbf{X}|\Theta) \quad = \quad \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \quad \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \quad \geq \quad \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

- If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, the above inequality becomes equality

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad = \quad \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log \frac{p(\mathbf{Z}|\mathbf{X}, \Theta) p(\mathbf{X}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}|\Theta)$$

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

- Thus for $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, we have

$$\log p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

# Parameter Estimation with Latent Variables

- Consider the 'incomplete' data log likelihood

$$\log p(\mathbf{X}|\Theta) \quad = \quad \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \quad \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \quad \geq \quad \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

- If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, the above inequality becomes equality

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad = \quad \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log \frac{p(\mathbf{Z}|\mathbf{X}, \Theta) p(\mathbf{X}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}|\Theta)$$

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

- Thus for $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, we have

$$\log p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.} = \mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] + \text{const.}$$

# Parameter Estimation with Latent Variables

- Consider the 'incomplete" data log likelihood

$$\log p(\mathbf{X}|\Theta) \quad = \quad \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \quad \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \quad \geq \quad \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

- If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, the above inequality becomes equality

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad = \quad \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log \frac{p(\mathbf{Z}|\mathbf{X}, \Theta) p(\mathbf{X}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}|\Theta)$$

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

- Thus for $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, we have

$$\log p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.} = \mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] + \text{const.}$$

- Thus $\log p(\mathbf{X}|\Theta)$ is tightly lower-bounded by $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$ which EM maximizes

**18**

Initialize the parameters: $\Theta^{old}$. Then alternate between these steps:

# The Expectation Maximization (EM) Algorithm

Initialize the parameters: $\Theta^{old}$. Then alternate between these steps:

- **E (Expectation) step:**

# The Expectation Maximization (EM) Algorithm

Initialize the parameters: $\Theta^{old}$. Then alternate between these steps:

- **E (Expectation) step:**
  - Compute the posterior $p(\mathbf{Z}|\mathbf{X}, \Theta^{old})$ over latent variables $\mathbf{Z}$ using $\Theta^{old}$

# The Expectation Maximization (EM) Algorithm

Initialize the parameters: $\Theta^{old}$. Then alternate between these steps:

- **E (Expectation) step:**
    - Compute the posterior $p(\mathbf{Z}|\mathbf{X}, \Theta^{old})$ over latent variables $\mathbf{Z}$ using $\Theta^{old}$
    - Compute the expected complete data log-likelihood w.r.t. *this* posterior

$$\mathcal{Q}(\Theta, \Theta^{old}) = \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \Theta^{old})}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

Initialize the parameters: $\Theta^{old}$. Then alternate between these steps:

- **E (Expectation) step:**
  - Compute the posterior $p(\mathbf{Z}|\mathbf{X}, \Theta^{old})$ over latent variables $\mathbf{Z}$ using $\Theta^{old}$
  - Compute the expected complete data log-likelihood w.r.t. *this* posterior

$$\mathcal{Q}(\Theta, \Theta^{old}) = \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \Theta^{old})}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

- **M (Maximization) step:**

Initialize the parameters: $\Theta^{old}$. Then alternate between these steps:

- **E (Expectation) step:**

  - Compute the posterior $p(\mathbf{Z}|\mathbf{X}, \Theta^{old})$ over latent variables $\mathbf{Z}$ using $\Theta^{old}$
  - Compute the expected complete data log-likelihood w.r.t. *this* posterior

$$\mathcal{Q}(\Theta, \Theta^{old}) = \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \Theta^{old})}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

- **M (Maximization) step:**
  - Maximize the expected complete data log-likelihood w.r.t. $\Theta$

$$\Theta^{new} = \arg\max_{\Theta} \mathcal{Q}(\Theta, \Theta^{old}) \quad \text{(if doing MLE)}$$

$$\Theta^{new} = \arg\max_{\Theta} \{\mathcal{Q}(\Theta, \Theta^{old}) + \log p(\Theta)\} \quad \text{(if doing MAP)}$$

Initialize the parameters: $\Theta^{old}$. Then alternate between these steps:

- **E (Expectation) step:**
  - Compute the posterior $p(\mathbf{Z}|\mathbf{X}, \Theta^{old})$ over latent variables $\mathbf{Z}$ using $\Theta^{old}$
  - Compute the expected complete data log-likelihood w.r.t. *this* posterior

$$\mathcal{Q}(\Theta, \Theta^{old}) = \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \Theta^{old})}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

- **M (Maximization) step:**
  - Maximize the expected complete data log-likelihood w.r.t. $\Theta$

$$\Theta^{new} = \arg\max_{\Theta} \mathcal{Q}(\Theta, \Theta^{old}) \quad \text{(if doing MLE)}$$

$$\Theta^{new} = \arg\max_{\Theta}\{\mathcal{Q}(\Theta, \Theta^{old}) + \log p(\Theta)\} \quad \text{(if doing MAP)}$$

- If the log-likelihood or the parameter values not converged then set $\Theta^{old} = \Theta^{new}$ and go to the E step.

Initialize the parameters: $\Theta^{old}$. Then alternate between these steps:

- ## E (Expectation) step:

  - Compute the posterior $p(\mathbf{Z}|\mathbf{X}, \Theta^{old})$ over latent variables $\mathbf{Z}$ using $\Theta^{old}$
  - Compute the expected complete data log-likelihood w.r.t. *this* posterior

$$\mathcal{Q}(\Theta, \Theta^{old}) = \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \Theta^{old})}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

- ## M (Maximization) step:

  - Maximize the expected complete data log-likelihood w.r.t. $\Theta$

$$\Theta^{new} = \arg\max_{\Theta} \mathcal{Q}(\Theta, \Theta^{old}) \quad \text{(if doing MLE)}$$

$$\Theta^{new} = \arg\max_{\Theta}\{\mathcal{Q}(\Theta, \Theta^{old}) + \log p(\Theta)\} \quad \text{(if doing MAP)}$$
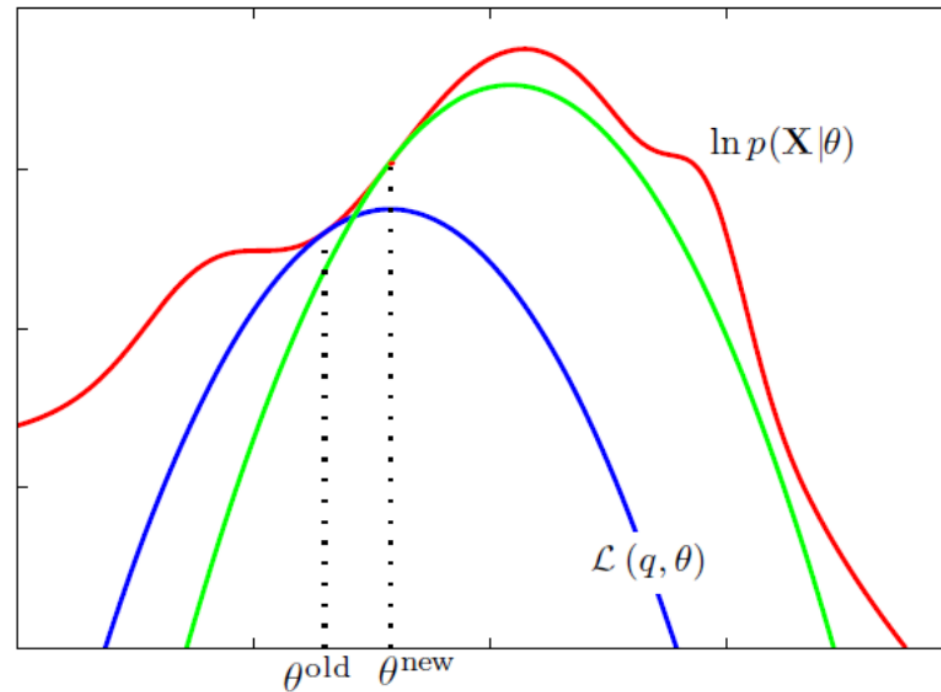
- If the log-likelihood or the parameter values not converged then set $\Theta^{old} = \Theta^{new}$ and go to the E step.

The algorithm converges to a local maxima of $p(\mathbf{X}|\Theta)$ (as we saw)

- E-step: Update of $q$ makes the $\mathcal{L}(q, \Theta)$ curve touch the $\log p(\mathbf{X}|\Theta)$ curve
- M-step gives the maxima $\Theta^{new}$ of $\mathcal{L}(q, \Theta)$
- Next E-step readjusts $\mathcal{L}(q, \Theta)$ curve (green) to meet $\log p(\mathbf{X}|\Theta)$ curve again
- This continues until a local maxima of $\log p(\mathbf{X}|\Theta)$ is reached

# EM: Some Comments

- A general framework for parameter estimation in latent variable models

- Very widely used in problems with "missing data", e.g., missing features, or missing labels (semi-supervised learning)

  - "Missing" parts can be treated as latent variables $z$ and estimated using EM

- More advanced probabilistic inference algorithms are based on similar ideas

  - E.g., variational Bayesian inference

- Very easy to extend to online learning setting and gives SGD like algorithms (will post a reading on "Online EM" on the class webpage)

- Note: The E and M steps may not always be possible to perform exactly (approximate inference methods may be needed in such cases)
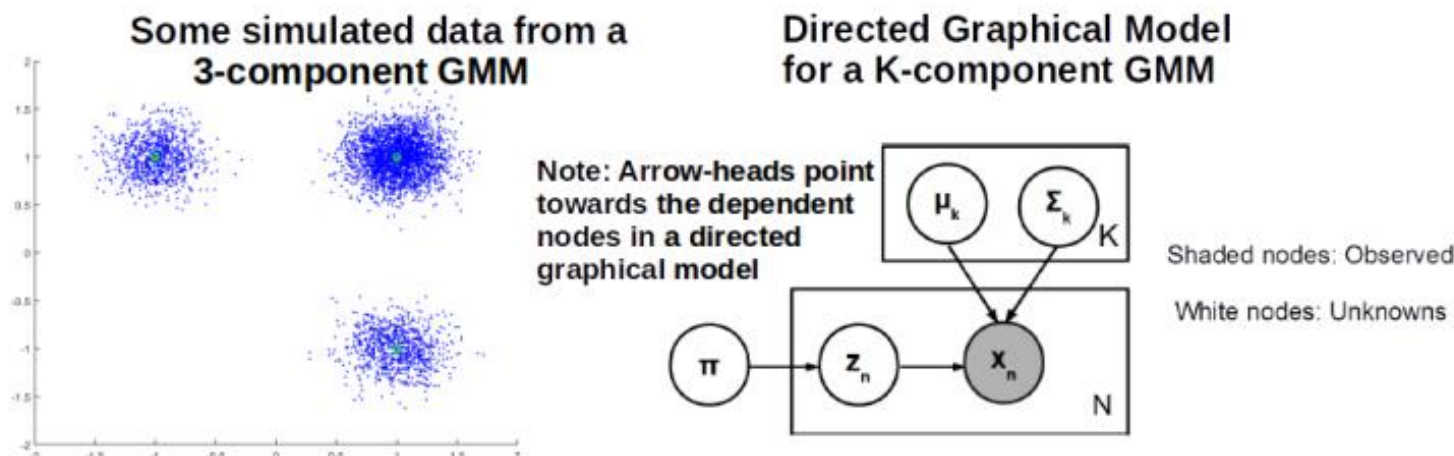
- The generative story for each $x_n$, $n = 1, 2, \ldots, N$

  - First choose one of the $K$ mixture components as

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

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

  $$x_n \sim \mathcal{N}(x_n|\mu_k, \Sigma_k) \qquad \text{(from the data distr. } p(x|z))$$

**Some simulated data from a 3-component GMM**

**Directed Graphical Model for a K-component GMM**

Note: Arrow-heads point towards the dependent nodes in a directed graphical model

Shaded nodes: Observed

White nodes: Unknowns

# Recap: Learning GMM

We derive the Expectation-Maximization (EM) algorithm for GMM with $K$ components:
- Observed data $X = \{x_1, \dots, x_N\}$
- Latent variables $Z = \{z_1, \dots, z_N\}$ , where $z_i \in \{1, \dots, K\}$.
- The parameters are $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$.

## 1. Complete-Data Likelihood

The joint distribution of observed data and latent data is:

$$p(X, Z|\Theta) = \prod_{i=1}^{N} p(x_i, z_i)\, p(z_i|\Theta) = \prod_{i=1}^{N} \pi_{z_i} \mathcal{N}\left(x_i \middle| \mu_{z_i}, \Sigma_{z_i}\right)$$

The complete-data log-likelihood is:

$$\log p(X, Z|\Theta) = \sum_{i=1}^{N} \log \pi_{z_i} + \log \mathcal{N}\left(x_i \middle| \mu_{z_i}, \Sigma_{z_i}\right)$$

## 2. E-Step

☐ Compute the posterior responsibility $\gamma_{ik} = p(z_i = k | x_i, \Theta^{old})$:

$$\gamma_{ik} = p(z_i = k | x_i, \Theta^{old})$$

$$= \frac{p(x_i | z_i = k, \Theta^{old}) p(z_i = k | \Theta^{old})}{p(x_i | \Theta^{old})}$$

$$= \frac{\pi_k^{old} \mathcal{N}(x_i | \mu_k^{old}, \Sigma_k^{old})}{\sum_{j=1}^{K} \pi_j^{old} \mathcal{N}(x_i | \mu_j^{old}, \Sigma_j^{old})}$$

☐ This is a soft assignment of $x_i$ to cluster $k$.

## 3. M-step: Maximize $Q(\Theta, \Theta^{old})$

☐ The $Q$-function is the expected complete-data log-likeligood

$$Q(\Theta, \Theta^{old}) = \mathrm{E}_{Z|X,\Theta^{old}}[\log p(X, Z|\Theta)] = \sum_{i=1}^{N} \sum_{k=1}^{K} \gamma_{ik}[\log \pi_k + \log \mathcal{N}(x_i|\mu_k, \Sigma_k)]$$

☐ Update Mixing coefficients $\pi_k$:

- Maximize $Q$ w.r.t. $\pi_k$ under the constraint $\sum_k \pi_k = 1$:

$$\mathcal{L} = Q(\Theta, \Theta^{old}) + \lambda\left(1 - \sum_{k=1}^{K} \pi_k\right)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \sum_{i=1}^{N} \frac{\gamma_{ik}}{\pi_k} - \lambda = 0 \quad \Rightarrow \quad \pi_k \propto \sum_{i=1}^{N} \gamma_{ik}$$

## 3. M-step: Maximize $Q(\Theta, \Theta^{old})$

$$Q(\Theta, \Theta^{old}) = \mathrm{E}_{Z|X,\Theta^{old}}[\log p(X, Z|\Theta)] = \sum_{i=1}^{N}\sum_{k=1}^{K} \gamma_{ik}[\log \pi_k + \log \mathcal{N}(x_i|\mu_k, \Sigma_k)]$$

□ Update Mixing coefficients $\pi_k$:

$$\mathcal{L} = Q(\Theta, \Theta^{old}) + \lambda \left(1 - \sum_{k=1}^{K} \pi_k\right)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \sum_{i=1}^{N} \frac{\gamma_{ik}}{\pi_k} - \lambda = 0 \quad \Rightarrow \quad \pi_k \propto \sum_{i=1}^{N} \gamma_{ik}$$

- Enforce constraint $\sum_k \pi_k = 1$:

$$\lambda = \sum_{k=1}^{K}\sum_{i=1}^{N} \gamma_{ik} = N \quad \Rightarrow \quad \pi_k^{new} = \frac{1}{N}\sum_{i=1}^{N} \gamma_{ik}$$

## 3. M-step: Maximize $Q(\Theta, \Theta^{old})$

$$Q(\Theta, \Theta^{old}) = E_{Z|X,\Theta^{old}}[\log p(X, Z|\Theta)] = \sum_{i=1}^{N}\sum_{k=1}^{K} \gamma_{ik}[\log \pi_k + \log \mathcal{N}(x_i|\mu_k, \Sigma_k)]$$

☐ The probability density function for a $D$-dimensional Gaussian:

$$\mathcal{N}(x \mid \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{D/2}|\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1}(x - \mu_k)\right)$$

☐ Update Means $\mu_k$:

$$\frac{\partial Q}{\partial \mu_k} = \sum_{i=1}^{N}\gamma_{ik}\Sigma_k^{-1}(x_i - \mu_k) = 0 \implies \mu_k^{new} = \frac{\sum_{i=1}^{N}\gamma_{ik}x_i}{\sum_{i=1}^{N}\gamma_{ik}}$$

## 3. M-step: Maximize $Q(\boldsymbol{\Theta}, \boldsymbol{\Theta}^{old})$

$$Q(\Theta, \Theta^{old}) = \mathrm{E}_{Z|X,\Theta^{old}}[\log p(X, Z|\Theta)] = \sum_{i=1}^{N} \sum_{k=1}^{K} \gamma_{ik}[\log \pi_k + \log \mathcal{N}(x_i|\mu_k, \Sigma_k)]$$

☐ The probability density function for a $D$ -dimensional Gaussian:

$$\mathcal{N}(x \mid \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{D/2}|\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1}(x - \mu_k)\right)$$

☐ Update Covariances $\Sigma_k$ :  $\boxed{\dfrac{\partial \log|\Sigma_k|}{\partial \Sigma_k^{-1}} = -\Sigma_k \qquad \dfrac{\partial \mathrm{tr}(A\Sigma_k^{-1})}{\partial \Sigma_k^{-1}} = A^T}$

$$\frac{\partial Q}{\partial \Sigma_k^{-1}} = \frac{1}{2} \sum_{i=1}^{N} \gamma_{ik}[\Sigma_k - (x_i - \mu_k)(x_i - \mu_k)^{\mathsf{T}}] = 0$$

$$\Sigma_k^{new} = \frac{\sum_{i=1}^{N} \gamma_{ik}(x_i - \mu_k^{new})(x_i - \mu_k^{new})^{\mathsf{T}}}{\sum_{i=1}^{N} \gamma_{ik}}$$

# Learning GMM: Summary

1. Initialize the $\text{Parameters } \Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$ randomly, or using K-means

2. Iterate until convergence (e.g., when $\log p(X|\Theta)$ ceases to increase

   a) E-step:

   $$\gamma_{ik} = p(z_i = k | x_i, \Theta^{old}) = \frac{\pi_k^{old}\mathcal{N}(x_i|\mu_k^{old}, \Sigma_k^{old})}{\sum_{j=1}^{K} \pi_j^{old}\mathcal{N}(x_i|\mu_j^{old}, \Sigma_j^{old})}$$

   b) M-step:

   $$\pi_k^{new} = \frac{1}{N}\sum_{i=1}^{N} \gamma_{ik}$$

   $$\mu_k^{new} = \frac{\sum_{i=1}^{N} \gamma_{ik} x_i}{\sum_{i=1}^{N} \gamma_{ik}}$$

   $$\Sigma_k^{new} = \frac{\sum_{i=1}^{N} \gamma_{ik}(x_i - \mu_k^{new})(x_i - \mu_k^{new})^\top}{\sum_{i=1}^{N} \gamma_{ik}}$$