



# 隐变量模型、EM算法和 变分自编码器

主讲人：林惊



# 无监督概率建模

- 在监督学习中，回归和分类都可以理解为是在学习条件概率分布

$$p(y|\mathbf{x}; \mathbf{w})$$

- 在回归任务中，条件概率密度函数的形式通常假设为

$$p(y|\mathbf{x}; \mathbf{w}) = \mathcal{N}(y; \mathbf{w}^T \mathbf{x}, \sigma^2)$$

- 在分类任务中，条件概率密度函数的形式通常假设为

$$\text{二分类: } p(y|\mathbf{x}) = (\sigma(\mathbf{xw}))^y \cdot (1 - \sigma(\mathbf{xw}))^{1-y}$$

$$\text{多分类: } p(\mathbf{y}|\mathbf{x}) = \prod_{k=1}^K [\text{softmax}_k(\mathbf{W}\mathbf{x})]^{y_k}$$

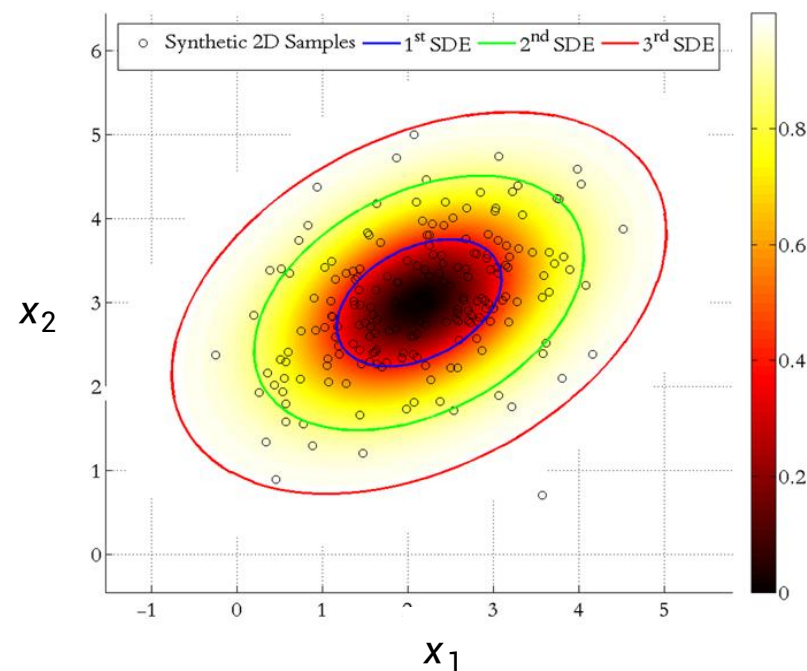
- 无监督学习也可以从学习概率分布的角度来理解。但它只关注输入数据  $\mathbf{x}$  的分布

$$p(\mathbf{x}; \mathbf{w})$$

- 建模  $\mathbf{x}$  要比标签  $y$  困难得多。一种朴素的方法是将  $p(\mathbf{x}; \mathbf{w})$  限制为高斯形式

$$p(\mathbf{x}; \mathbf{w}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- 通过优化  $\boldsymbol{\mu}$  和  $\boldsymbol{\Sigma}$  来最优地描述数据点  $\{\mathbf{x}^{(n)}\}_{n=1}^N$



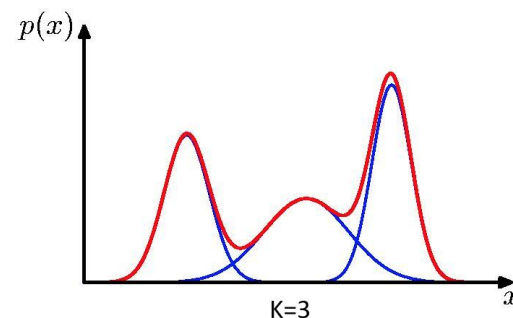
显然，该模型的表达能力非常有限

# 为什么需要隐变量？

- **理由1：** 通过组合简单模型构建更具表达能力的模型

- 假设存在一个简单的类别分布  $p(z) = \text{Cat}(K, \boldsymbol{\pi})$  和一个高斯分布  $p(x) = \mathcal{N}(x|\mu, \sigma^2)$
- 如果单独使用它们，只能对简单的统计关系进行建模
- 但如果我们将它们组合成  $p(x, z) = p(x|z)p(z)$ ，所得到的边缘分布  $p(x)$  表达能力强得多

$$p(x) = \sum_z p(x|z)p(z) = \sum_{k=1}^K \pi_z \mathcal{N}(x|\mu_z, \sigma_z^2)$$



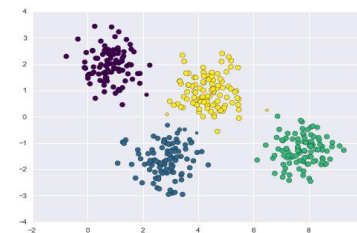
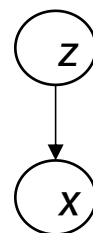
- 理论上，它能够表示任何复杂的分布

- 理由2：数据中隐藏的结构

- 1) 具有隐藏簇状结构的数据

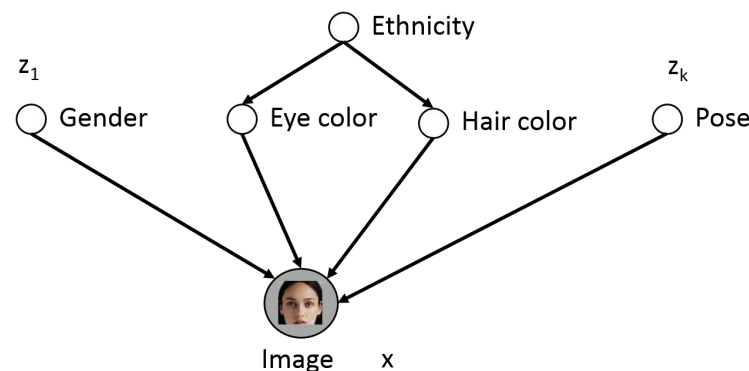
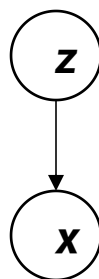
$$z_n \sim \text{Cat}(K, \boldsymbol{\pi})$$

$$x_n \sim \mathcal{N}(x | \boldsymbol{\mu}_{z_n}, \boldsymbol{\Sigma}_{z_n})$$



- 2) 面向文档的主题模型

- 3) 图像建模



- 在上面的例子中，隐变量  $z$  通常对应于高层级特征
- 如果考虑到这种隐藏结构，就可以获得更具可解释性的模型

# 隐变量模型的一般形式

- 隐变量模型 (LVMs) : 含有隐变量的概率模型

$$p(\mathbf{x}, \mathbf{z})$$

- $\mathbf{x}$  是我们感兴趣的随机变量
- $\mathbf{z}$  是隐变量 (Latent Variable)

➤ 有时可能存在多个隐变量  $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K$

$$p(\mathbf{x}, \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K)$$

- 关于我们感兴趣的变量  $\mathbf{x}$  的概率模型是

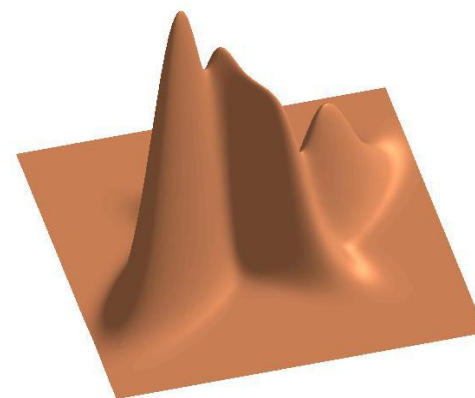
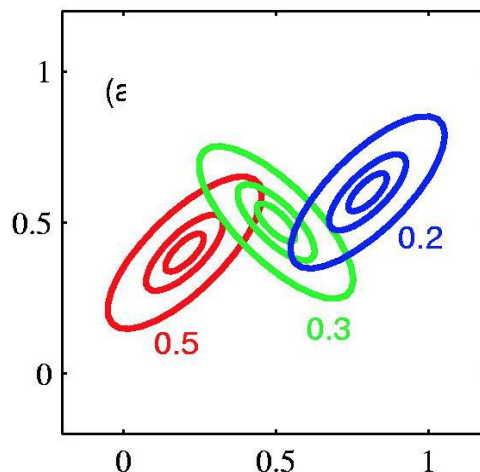
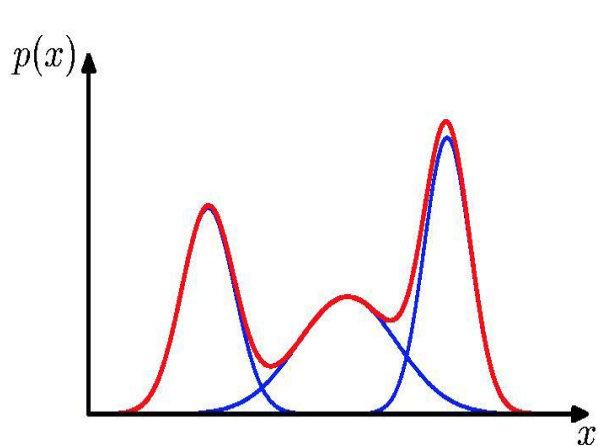
$$p(\mathbf{x}) = \int_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \quad \text{或} \quad p(\mathbf{x}) = \int_{\mathbf{z}_1 \dots \mathbf{z}_K} p(\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_K) d\mathbf{z}_1 \dots d\mathbf{z}_K$$

# 高斯混合分布

- 高斯混合分布的表达式

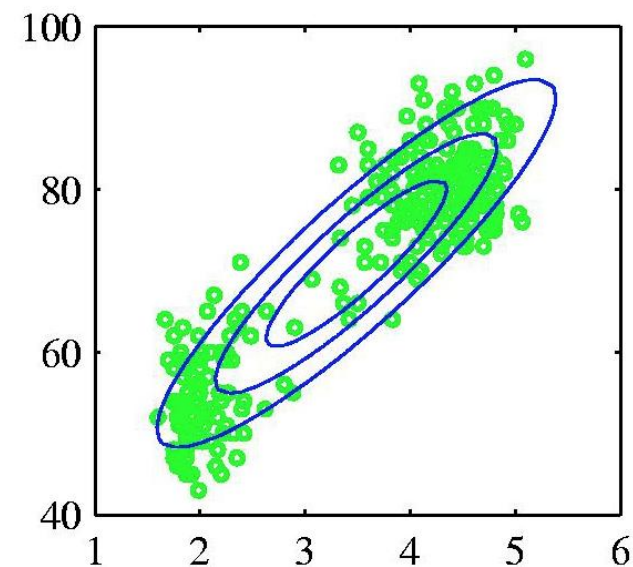
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- $K$  表示高斯分布的数量
- $\pi_k$  是第  $k$  个分布的权重，满足  $\sum_{k=1}^K \pi_k = 1$
- $\boldsymbol{\mu}_k$  和  $\boldsymbol{\Sigma}_k$  分别是第  $k$  个高斯分布的均值向量和协方差矩阵

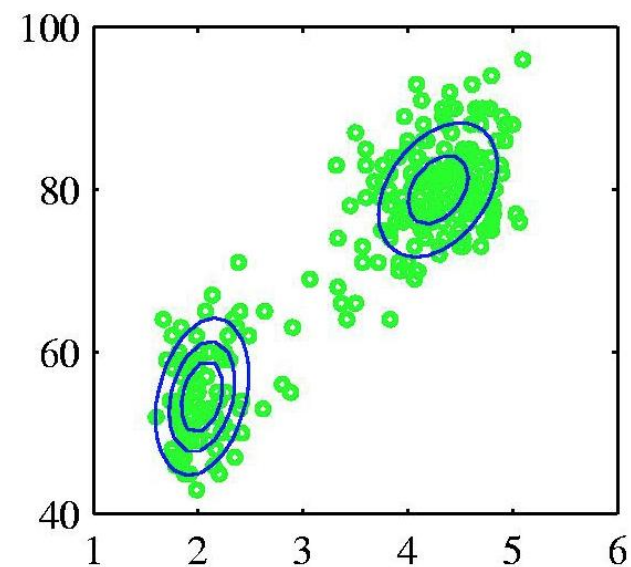




- 这些绿色的点很难用单个高斯分布建模



- 但如果用两个高斯分布的混合来建模，效果会好得多





# 将高斯混合分布表示为隐变量模型

- 对于一个隐变量模型  $p(\mathbf{x}, \mathbf{z})$ ，如果我们将其条件分布  $p(\mathbf{x}|\mathbf{z})$  和先验分布  $p(\mathbf{z})$  设为

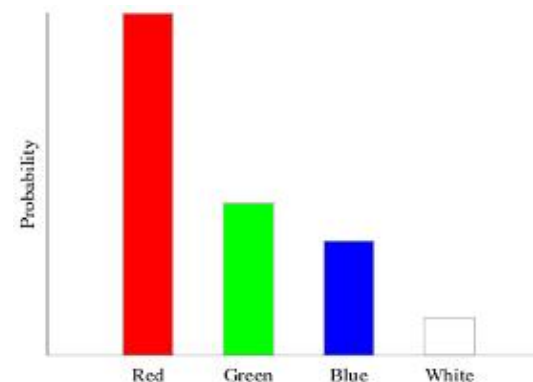
$$p(\mathbf{x}|\mathbf{z} = \mathbf{1}_k) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$p(\mathbf{z} = \mathbf{1}_k) = \pi_k$$

- $\mathbf{z}$  只能是一个 one-hot 向量，其中  $\mathbf{1}_k$  表示第  $k$  个元素为1
- 由于  $p(\mathbf{z} = \mathbf{1}_k) = \pi_k$ ， $p(\mathbf{z})$  实际上表示一个类别分布，即

$$p(\mathbf{z}) = \text{Cat}(\mathbf{z}; \boldsymbol{\pi})$$

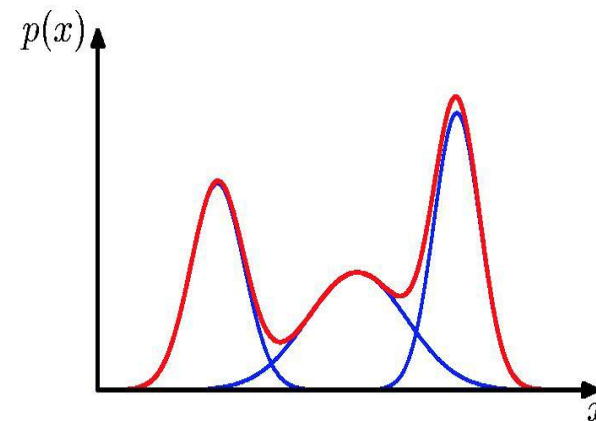
其中  $\text{Cat}(\mathbf{z} = \mathbf{1}_k; \boldsymbol{\pi}) = \pi_k$  且  $\boldsymbol{\pi} = [\pi_1, \pi_2, \dots, \pi_K]$



- 由于  $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$ ，我们可以很容易得到

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

这正是高斯混合分布



- 因此，高斯混合分布也可以等价用隐变量模型来  $p(\mathbf{x}, \mathbf{z})$  表示，其中  $p(\mathbf{x}|\mathbf{z} = \mathbf{1}_k) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ ， $p(\mathbf{z} = \mathbf{1}_k) = \pi_k$

高斯混合分布的隐变量分布  $p(\mathbf{x}, \mathbf{z})$  可以写成如下形式

$$p(\mathbf{x}, \mathbf{z}) = \prod_{k=1}^K [\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_k}$$

# 训练目标：最大化边缘概率分布

- 给定一组训练数据  $\{\mathbf{x}^{(n)}\}_{n=1}^N$ ，目标是学习分布的参数

$$\{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K \triangleq \boldsymbol{\theta}$$

- 假设数据点  $\mathbf{x}^{(n)}$  是独立同分布的，因此我们可以将联合分布写为

$$p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \prod_{n=1}^N \underbrace{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}^{(n)}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{p(\mathbf{x}^{(n)})}$$

未使用隐变量形式

- 对于概率模型，训练目标是**最大化对数似然函数**，即

$$\log p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}^{(n)}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

# 针对问题的一般形式

- 给定联合分布

$$p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}),$$

其中  $\mathbf{x}$  是观测变量,  $\mathbf{z}$  是隐变量, 我们需要最大化关于  $\boldsymbol{\theta}$  的对数似然, 即

$$\boldsymbol{\theta} = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta}),$$

其中

$$p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})$$

我们拥有的是 联合概率密度函数  $p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})$ , 但需要优化的是边缘概率密度函数

$p(\mathbf{x}; \boldsymbol{\theta})$

# 一般隐变量模型的训练

- 根据边缘分布与联合分布的关系，可知

$$p_{\theta}(\mathbf{x}_n) = \int_{\mathbf{z}_n} p_{\theta}(\mathbf{x}_n, \mathbf{z}_n) d\mathbf{z}_n$$

训练目标就是最大化对数似然函数

$$\max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(\mathbf{x}_n)$$

- 但是，由于积分的存在，很多时候，获得 $p_{\theta}(\mathbf{x}_n)$ 的表达式是很困难的
- 因此，对于一般隐变量模型，计算对数似然的梯度 $\frac{d \log p_{\theta}(\mathbf{x})}{d\theta}$ 是很困难的，使用梯度下降法来训练模型就变得非常具有挑战性

# 对数似然的重表示

- 对数似然可以重写为

$$\begin{aligned}\log p(\mathbf{x}; \boldsymbol{\theta}) &= \sum_{\mathbf{z}} \overbrace{q(\mathbf{z})}^{\forall \text{ 分布 } q(\mathbf{z})} \log p(\mathbf{x}) \\&= \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}) q(\mathbf{z})} \\&= \underbrace{\sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})}}_{\mathfrak{l}(q, \boldsymbol{\theta})} + \underbrace{\sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})}}_{KL(q||p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}))} \\&= \mathfrak{l}(q, \boldsymbol{\theta}) + KL(q||p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})), \text{ for } \forall \boldsymbol{\theta}, q(\mathbf{z})\end{aligned}$$

**注：** KL 散度用于衡量两个分布  $q$  和  $p$  之间的距离，定义为

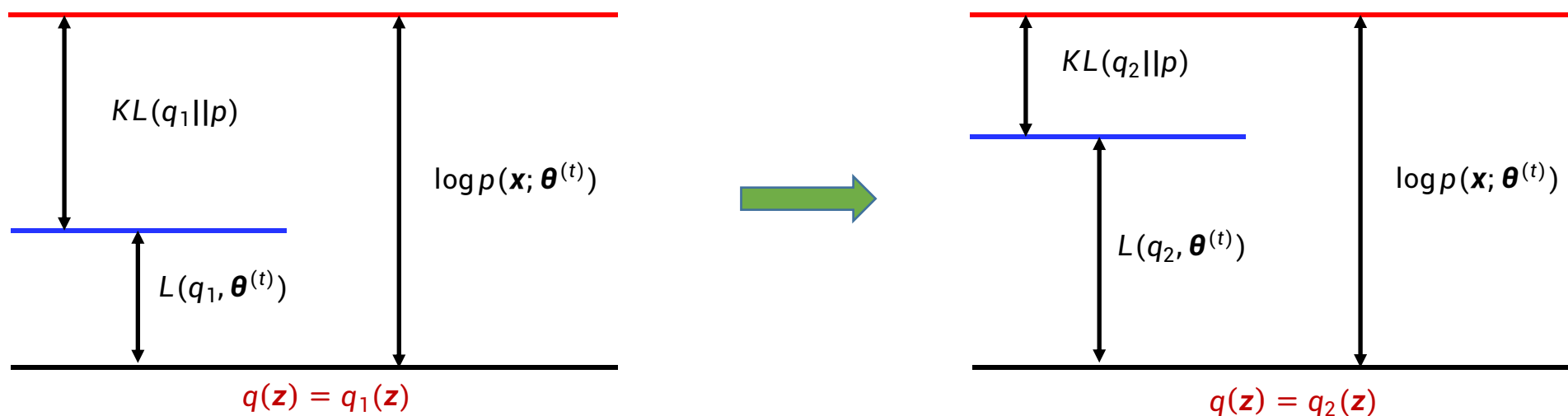
$$KL(q||p) \triangleq \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z} \geq 0$$

- 因此，我们将第  $t$  次迭代的参数记为  $\boldsymbol{\theta}^{(t)}$ ，可得

$$\log p(\mathbf{x}; \boldsymbol{\theta}^{(t)}) = \mathbb{1}(q, \boldsymbol{\theta}^{(t)}) + KL(q||p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}))$$

该等式对任何分布  $q(\mathbf{z})$  都成立

- 不同的  $q(\mathbf{z})$  会导致  $\log p(\mathbf{x}; \boldsymbol{\theta}^{(t)})$  的不同分解





# EM 算法的理论依据

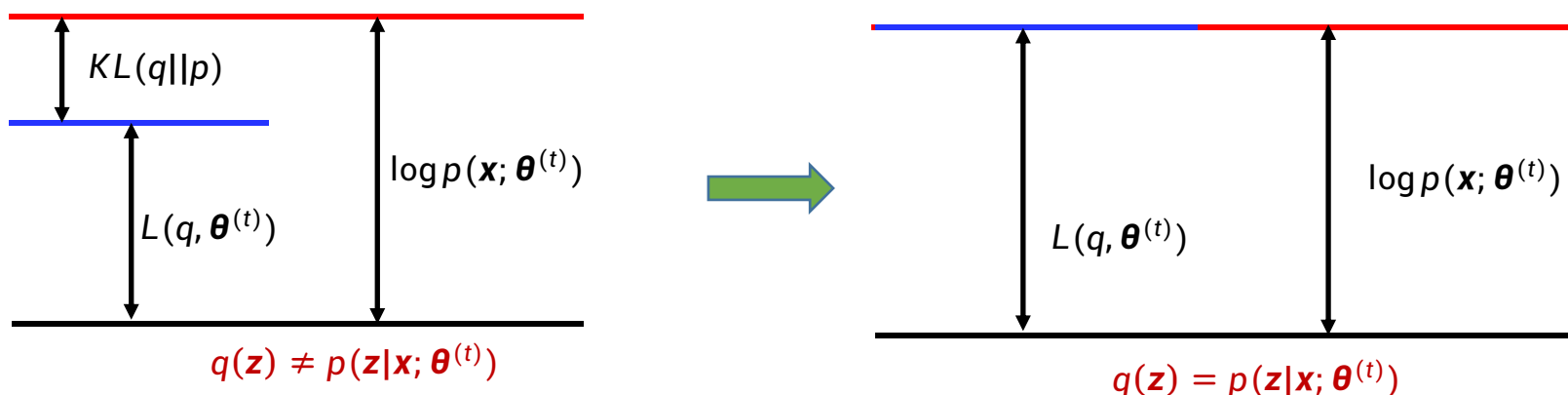
$$\log p(\mathbf{x}; \boldsymbol{\theta}^{(t)}) = \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}^{(t)})}{q(\mathbf{z})} + \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})}$$

- 如果我们设  $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$ ，那么可以得到

$$KL(q||p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})) = 0$$

因此，我们有

$$\begin{aligned} \log p(\mathbf{x}; \boldsymbol{\theta}^{(t)}) &= \mathbb{E}_{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})} \log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}^{(t)}) \\ &= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}^{(t)})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})} \end{aligned}$$



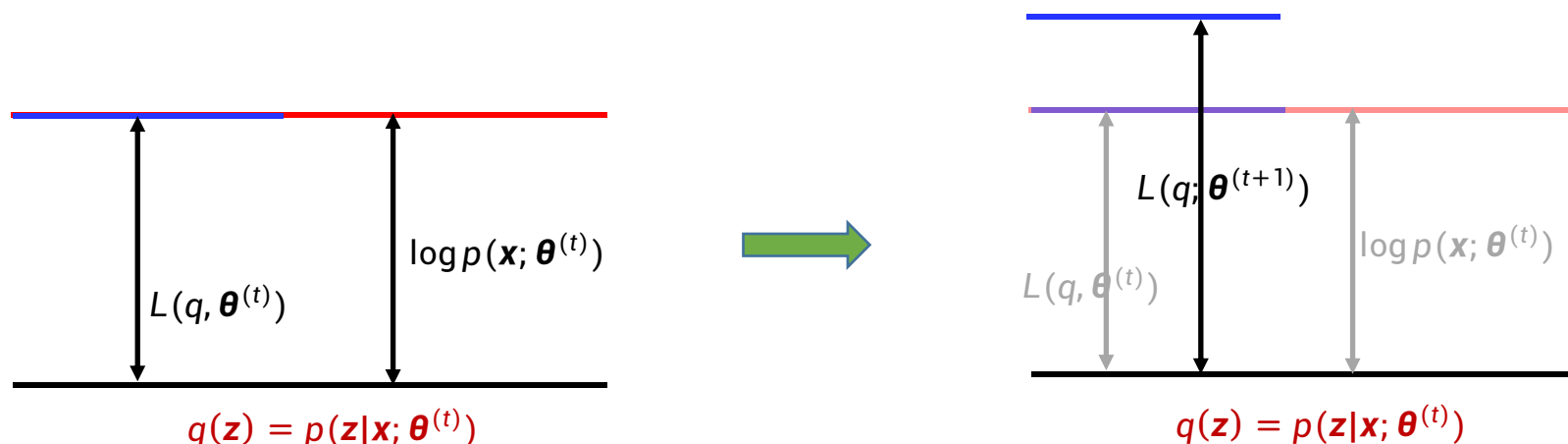
$$\log p(\mathbf{x}; \boldsymbol{\theta}^{(t)}) = \mathbb{E}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}^{(t)}) = \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}^{(t)})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})}$$

- 如果我们更新  $\boldsymbol{\theta}$  为

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} \mathbb{E}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}),$$

那么我们必然得到如下关系

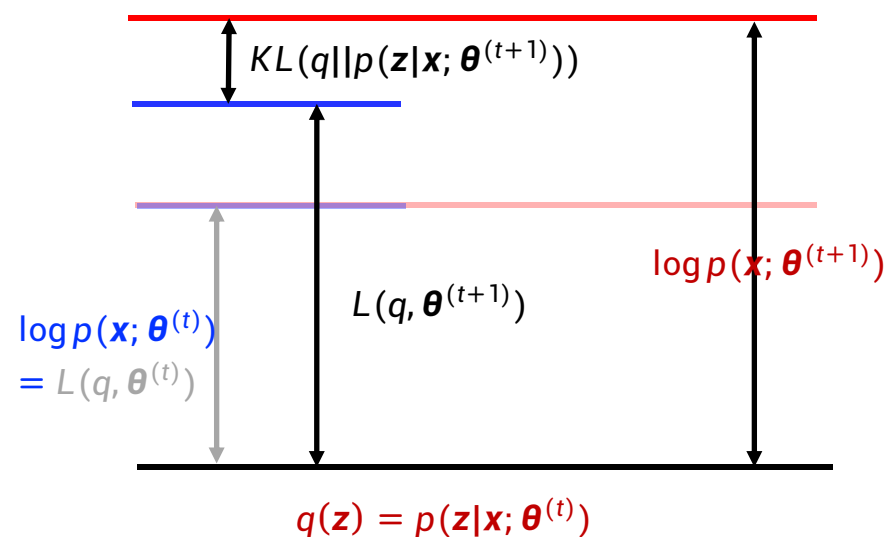
$$\mathbb{E}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}^{(t+1)}) \geq \underbrace{\mathbb{E}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}^{(t)})}_{= \log p(\mathbf{x}; \boldsymbol{\theta}^{(t)})}$$



- 通过设置  $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$ , 我们得到

$$\log p(\mathbf{x}; \boldsymbol{\theta}^{(t+1)}) = \underbrace{\mathbb{1}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}^{(t+1)})}_{\geq \log p(\mathbf{x}; \boldsymbol{\theta}^{(t)})} + \underbrace{KL(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) || p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t+1)}))}_{\geq 0}$$

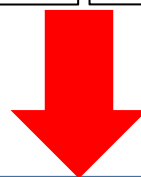
KL 散度始终为非负



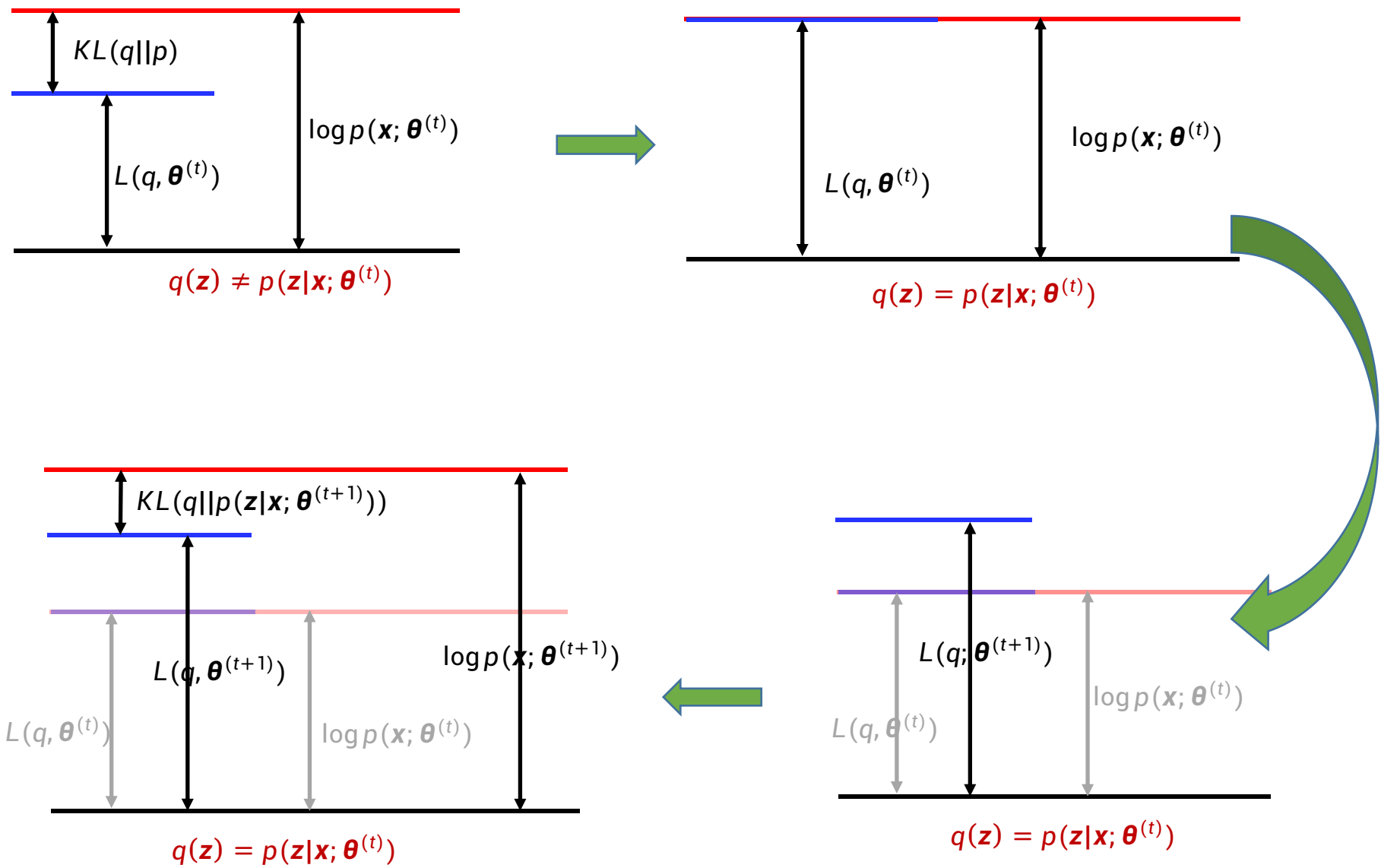
- 因此，我们可以看到

$$\log p(\mathbf{x}; \boldsymbol{\theta}^{(t+1)}) \geq \log p(\mathbf{x}; \boldsymbol{\theta}^{(t)})$$

KL 散度始终为非负



$\max_{\boldsymbol{\theta}} \mathbb{I}(p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}), \boldsymbol{\theta})$  能够保证每一步的似然函数值都会增加



# EM 算法

- 算法

*E 步*: 计算期望

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)}) = \mathbb{E}_{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})} [\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})]$$

*M 步*: 更新参数 *EM 算法可以保证每一步的似然函数值都会增加*

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)})$$

- EM 算法的关键要素

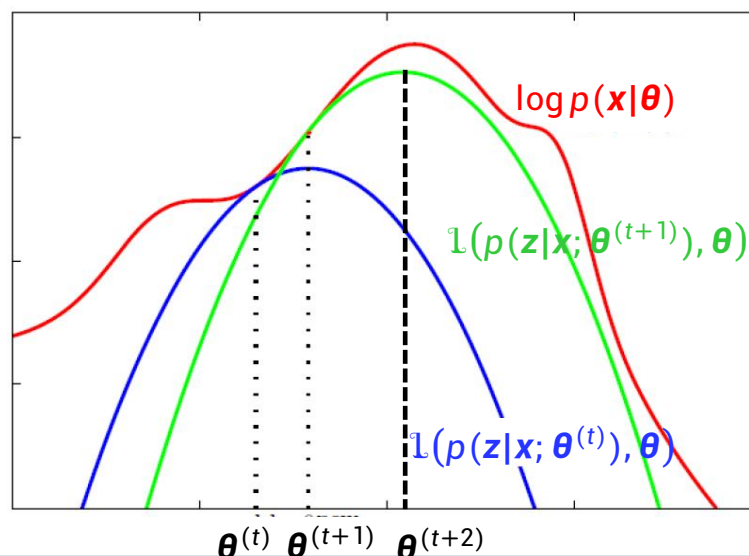
1) 后验分布  $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$

2) 联合分布对数  $\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})$  关于后验概率的期望

3) 最大化

# 参数空间视角

- 1) E 步 ( $t$ ): 给定模型参数  $\theta^{(t)}$ , 推导出  $\mathbb{L}(p(z|x; \theta^{(t)}), \theta)$  的表达式
- 2) M 步 ( $t$ ): 计算最优值  $\theta^{(t+1)} = \arg \max_{\theta} \mathbb{L}(p(z|x; \theta^{(t)}), \theta)$
- 3) E 步 ( $t+1$ ): 给定模型参数  $\theta^{(t+1)}$ , 推导出  $\mathbb{L}(p(z|x; \theta^{(t+1)}), \theta)$  的表达式
- 4) 重复以上过程直至收敛

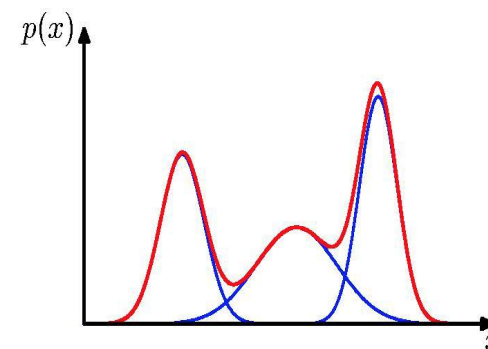




# 高斯混合模型回顾

- 对于一个高斯混合分布，即

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



它可以表示为联合分布的边缘分布

$$\begin{aligned} p(\mathbf{x}, \mathbf{z}) &= p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \\ &= \prod_{k=1}^K [\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_k} \end{aligned}$$

- $\mathbf{z} = [z_1, z_2, \dots, z_K]$  服从参数为  $\pi$  的类别分布

# EM 算法的两个步骤

- 这是一个隐变量模型，因此我们可以使用 EM 算法来优化它

注:  $\max_{\theta} \mathbb{I}(p(\mathbf{z}|\mathbf{x}; \theta^{(t)}), \theta)$  等价于  $\max_{\theta} Q(\theta; \theta^{(t)})$

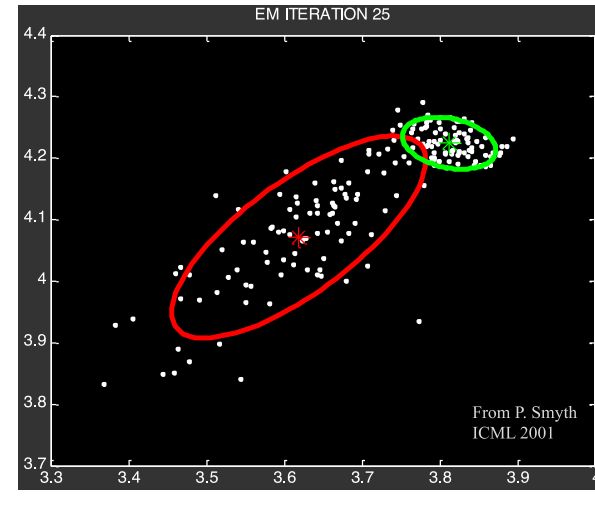
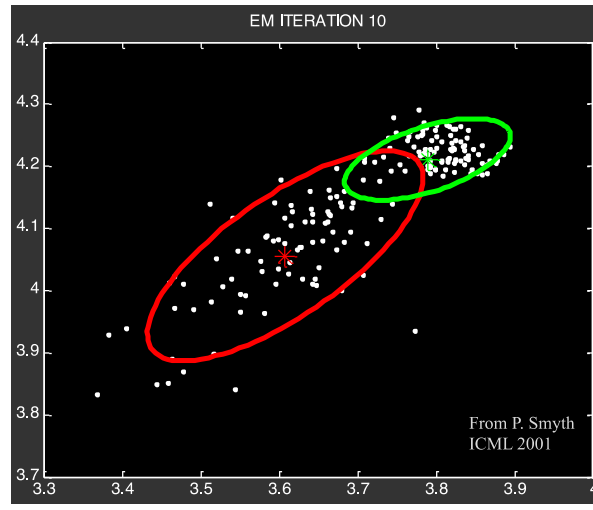
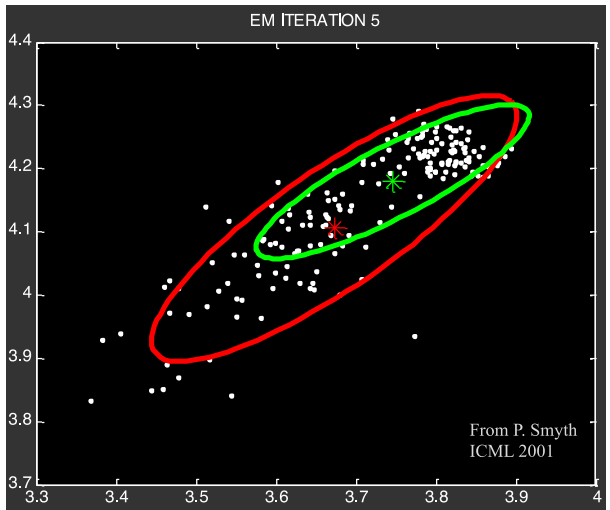
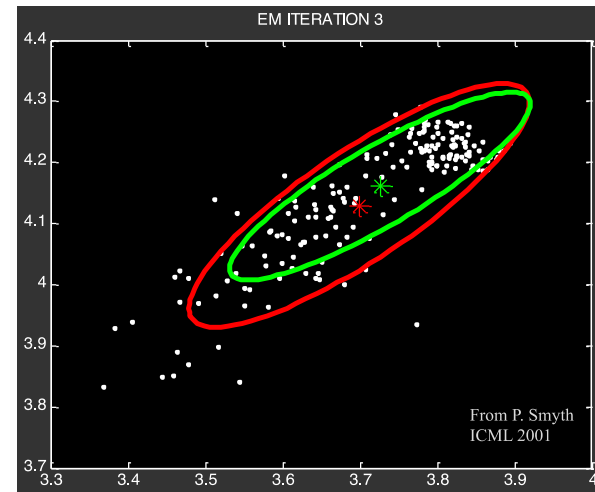
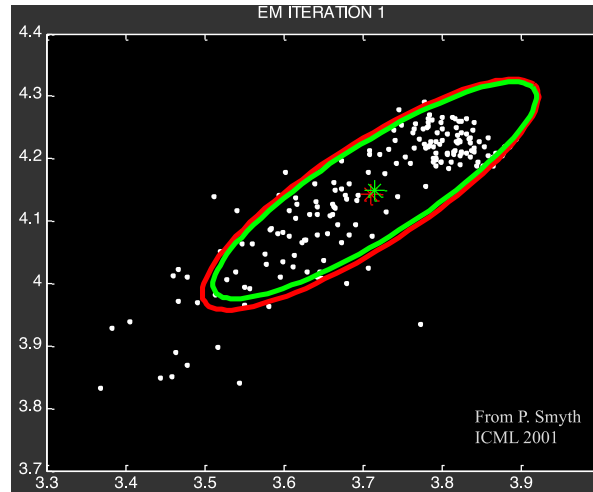
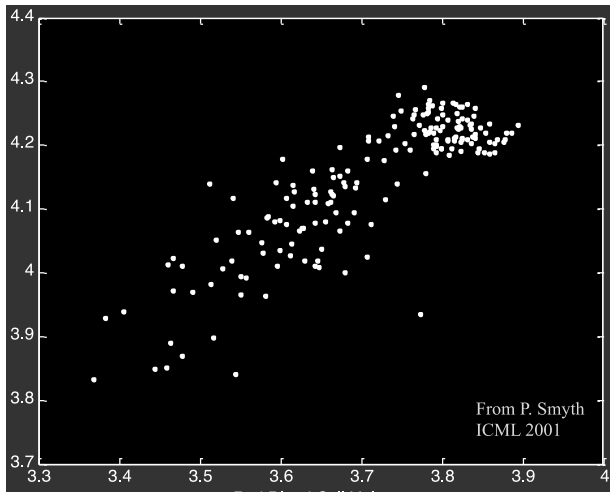
- 回顾: EM 算法的关键要素

➤ E 步: 计算关于后验概率  $p(\mathbf{z}|\mathbf{x}; \theta^{(t)})$  的期望

$$Q(\theta; \theta^{(t)}) = \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}^{(n)}|\mathbf{x}^{(n)}; \theta^{(t)})} [\log p(\mathbf{x}^{(n)}, \mathbf{z}^{(n)}; \theta)]$$

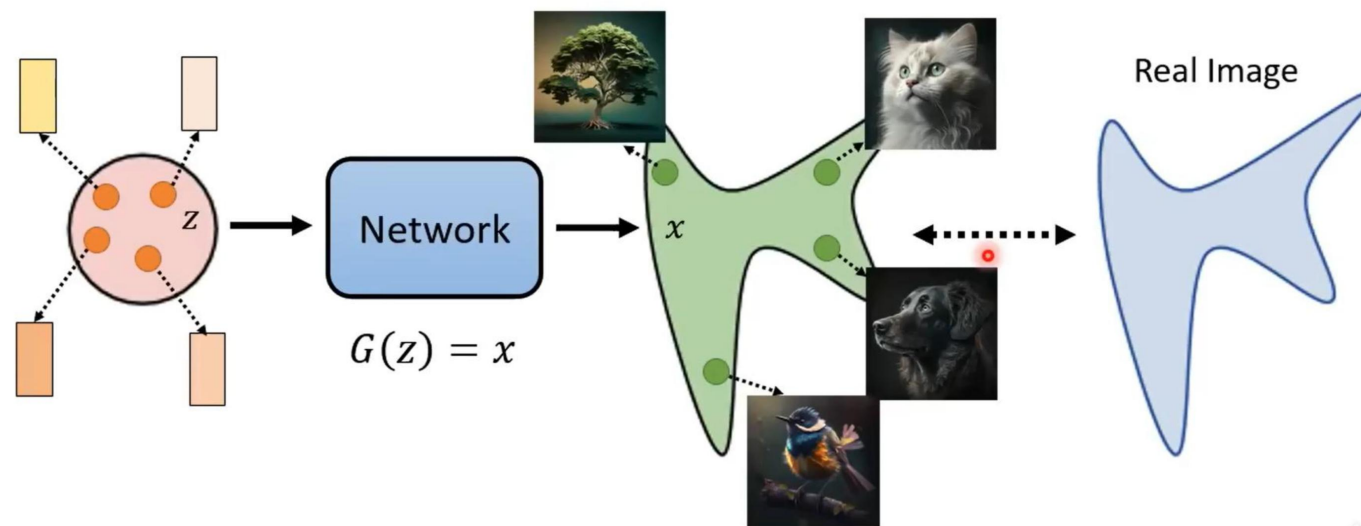
➤ M 步: 最大化

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta; \theta^{(t)})$$



# 从 EM 算法到变分自编码器

- 变分自编码器（Variational AutoEncoder, VAE）的直接数学目标就是：  
输入一个向量 $z$ ，经过模型得到 $G(z) = x$ ，希望这个得到的 $x$ 的分布 $p_{\theta}(x)$   
与实际已有数据的分布 $p_{data}(x)$ 尽可能接近



- 优化目标：

$$\theta^* = \arg \min_{\theta} KL(p_{\theta}(x), p_{data}(x)) = \arg \max_{\theta} \mathbb{E}_{x \sim p(data)} [p_{\theta}(x)]$$

# 从 EM 算法到变分自编码器

- 优化目标:

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim p(\text{data})} [p_{\theta}(x)]$$

- EM算法最核心的地方在于，为隐变量 $z$ 引入了一个分布 $q(z)$ ：

$$L(\theta) = \log p_{\theta}(x) = \int_z q(z) \log p_{\theta}(x|z) dz$$

$$= \int_z q(z) \log \left( \frac{p(x, z|\theta)}{p(z|x, \theta)} \cdot \frac{q(z)}{q(z)} \right) dz$$

$$= \int_z q(z) \log \frac{p(x|z, \theta) \cdot p(z)}{q(z)} dz$$

$$= \underbrace{\mathbb{E}_{z \sim q(z)} \left[ \log p_{\theta}(x|z) \right]}_{\text{ELBO}} + \underbrace{KL(q(z) \| p_{\theta}(z|x))}_{\text{KL Divergence}}$$

# 从 EM 算法到变分自编码器

- 优化目标:

$$L(\theta) = \underbrace{\mathbb{E}_{z \sim q(z)} [\log p_{\theta}(x|z)] - KL(q(z) \| p(z))}_{ELBO} + \underbrace{KL(q(z) \| p_{\theta}(z|x))}$$

- EM算法的基本思路: 先让KL散度取0, 这样只剩一个ELBO项, 最后让ELBO项最大。不断迭代直到收敛

**E 步:** 令  $q^{(t+1)}(z) = p_{\theta}(z|x)$ , 计算期望:

$$Q(\theta; \theta^{(t)}) = \mathbb{E}_{z \sim p_{\theta}(z|x)} [\log p_{\theta}(x|z)]$$

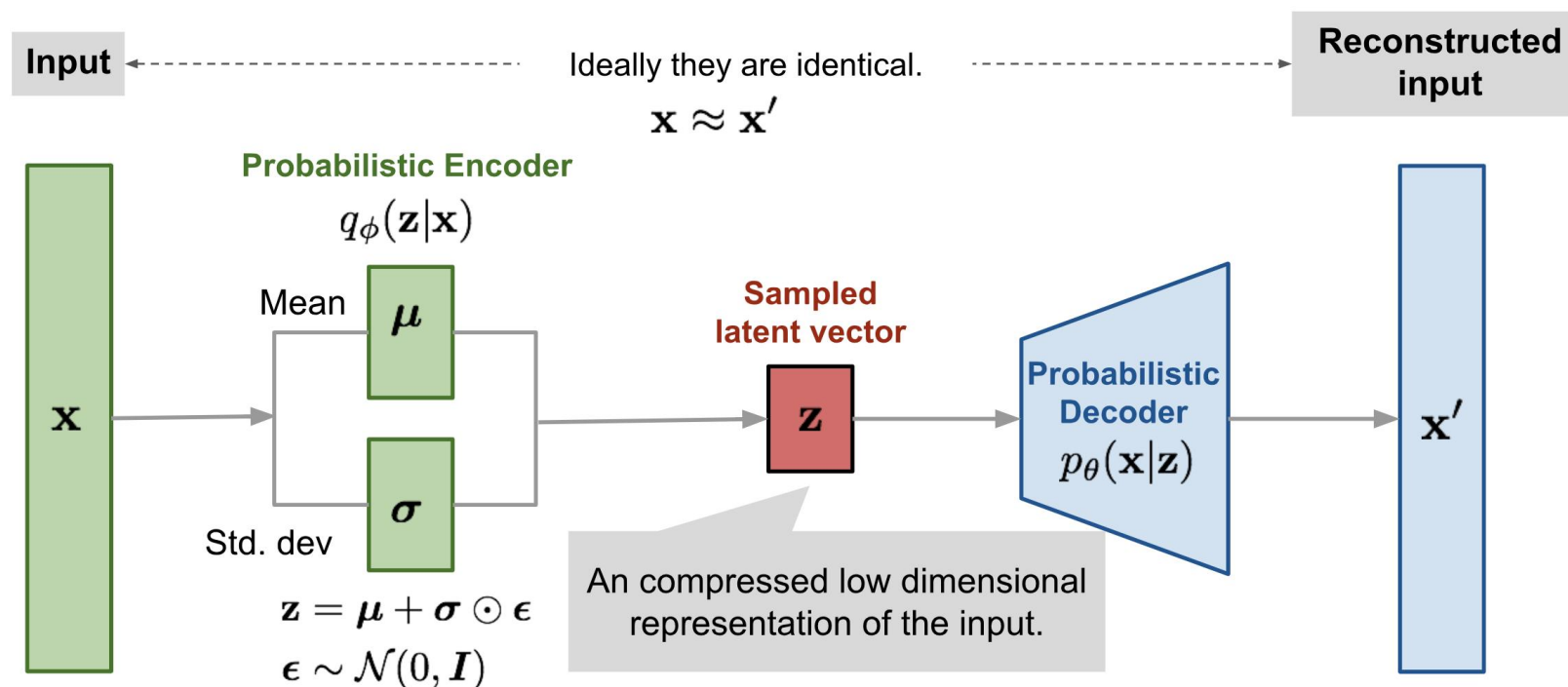
**M步:** 更新参数:

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta; \theta^{(t)})$$

# 变分自编码器

$$L(\theta) = \underbrace{\mathbb{E}_{z \sim q(z)} [\log p_{\theta}(x|z)] - KL(q(z) \| p(z))}_{ELBO} + \underbrace{KL(q(z) \| p_{\theta}(z|x))}$$

- 第一个改进：合并E-step和M-step，同时优化 $q(z)$ ,  $\theta$ 来达到一步到位 $\max_{q, \theta} ELBO$
- 第二个改进：对于 $q(z)$ 这个很难预测的分布，用深度学习模型 $q_{\phi}(z|x)$ 来进行控制





# 变分自编码器

- 变分自编码器的优化目标：

$$L(\theta) = \underbrace{\mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)]}_{\text{重建项}} + \underbrace{KL(q_{\phi}(z|x) \| p(z))}_{\text{正则项}}$$

