
Generative Models: Fundamentals and Applications

Lecture 4:

Generative models for continuous data

Shuigeng Zhou, Yuxi Mi
College of CSAI

October 13, 2025



目录



- 生成分类器
- 单高斯模型
- 高斯混合模型
- 隐马尔可夫模型



生成分类器

- 关键：类条件密度函数 $p(\mathbf{x}|y = c, \theta)$

$$\begin{aligned} p(y = c|\mathbf{x}, \theta) &= \frac{p(\mathbf{x}, y=c | \theta)}{p(\mathbf{x}|\theta)} \\ &= \frac{p(\mathbf{x}, y=c | \theta)}{\sum_{c' \in \mathcal{C}} p(\mathbf{x}, y=c' | \theta)} \\ &= \frac{p(y=c | \theta)p(\mathbf{x}|y=c, \theta)}{\sum_{c' \in \mathcal{C}} p(y=c' | \theta)p(\mathbf{x}|y=c', \theta)} \\ &\propto p(y = c | \theta)p(\mathbf{x}|y = c, \theta) \end{aligned}$$

生成分类器



■ 分类器

$$\hat{y}(\mathbf{x}) = \operatorname{argmax}_c [\log p(y = c | \boldsymbol{\pi}) + \log p(\mathbf{x} | \boldsymbol{\theta}_c)]$$

□ $\boldsymbol{\pi}$: 类先验的参数

■ $\pi_c = \frac{N_c}{N}$

□ $\boldsymbol{\theta}_c$: 第c个类的模型参数

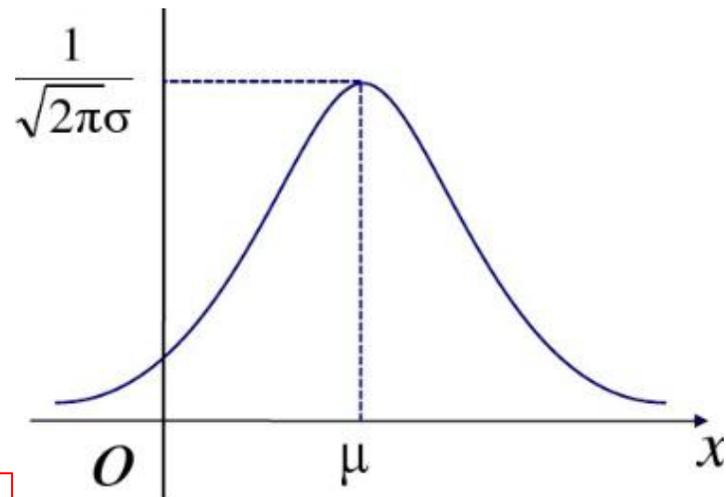
目录



- 生成分类器
- 单高斯模型
- 高斯混合模型
- 隐马尔可夫模型

单高斯模型

- 一元高斯分布
 - 一维随机变量 x
 - 均值 μ
 - 方差 σ^2
 - 概率密度函数



$$\mathcal{N}(x|\mu, \sigma^2) \triangleq \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$



单高斯模型

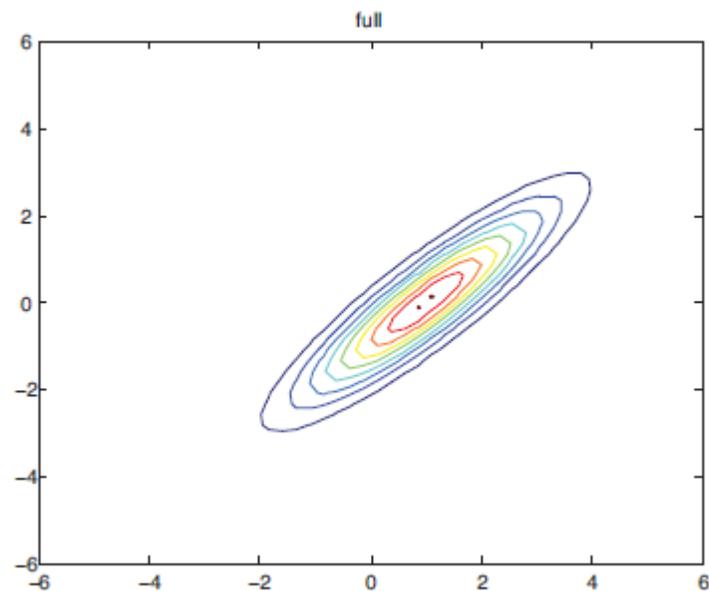
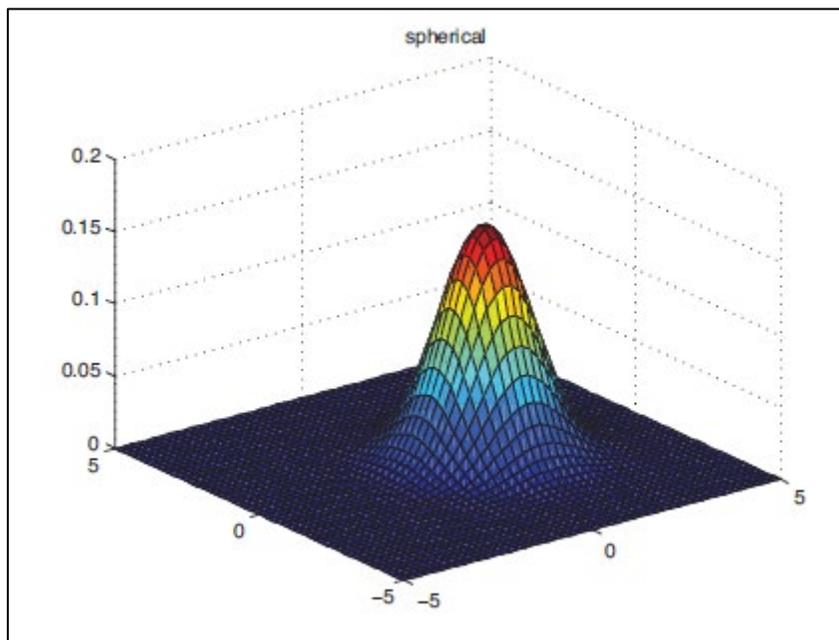
- 多元高斯分布 (multivariate Gaussian)
 - D维随机变量 \mathbf{x}
 - 均值 $\boldsymbol{\mu}$
 - 协方差矩阵 $\boldsymbol{\Sigma}$
 - 联合概率密度函数

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

单高斯模型



■ 二元高斯分布示意图





单高斯模型

- 马氏距离 (Mahalanobis distance)

$$D_M = (x - \mu)^T \Sigma^{-1} (x - \mu)$$

- 等价于数据x**变换**之后的欧氏距离

- 变换：平移 μ ，旋转 U

- 对协方差矩阵 Σ 进行特征分解： $\Sigma = U\Lambda U^T$

- U ：特征向量组成的正交矩阵，且 $U^T U = I$

- Λ ：特征值构成的对角矩阵



单高斯模型

- 马氏距离 (Mahalanobis distance)

$$D_M = (x - \mu)^T \Sigma^{-1} (x - \mu)$$

- 变换：平移 μ ，旋转 U

- 注意到 $\Sigma = U\Lambda U^T$ ，并记 $z_i = u_i^T (x - \mu)$

$$\begin{aligned} D_M &= (x - \mu)^T \mathbf{U}\Lambda^{-1}\mathbf{U}^T (x - \mu) \\ &= (x - \mu)^T \left(\sum_{i=1}^D \frac{1}{\lambda_i} u_i u_i^T \right) (x - \mu) \\ &= \sum_{i=1}^D \frac{1}{\lambda_i} (x - \mu)^T u_i u_i^T (x - \mu) \\ &= \sum_{i=1}^D \frac{1}{\lambda_i} z_i^T z_i \end{aligned}$$

单高斯模型



- 给定N个独立同分布的样本 $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ，估计模型参数 $\boldsymbol{\mu}, \boldsymbol{\Sigma}$?
 - 极大似然估计 MLE
 - 最大后验概率估计 MAP



极大似然估计

- 对数似然函数

- 令 $\Sigma = U\Lambda^{-1}U^T$

$$\begin{aligned} L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \log p(\mathcal{D}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^N \log p(\mathbf{x}_i|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \sum_{i=1}^N \log \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right] \\ &= -\frac{DN}{2} \log 2\pi + \frac{N}{2} \log |\Lambda| - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \end{aligned}$$

极大似然估计



Theorem 4.1.1 (MLE for a Gaussian). *If we have N iid samples $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then the MLE for the parameters is given by*

$$\hat{\boldsymbol{\mu}}_{mle} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \triangleq \bar{\mathbf{x}} \quad (4.6)$$

$$\hat{\boldsymbol{\Sigma}}_{mle} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T = \frac{1}{N} \left(\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right) - \bar{\mathbf{x}} \bar{\mathbf{x}}^T \quad (4.7)$$

That is, the MLE is just the empirical mean and empirical covariance. In the univariate case, we get the following familiar results:

$$\hat{\mu} = \frac{1}{N} \sum_i x_i = \bar{x} \quad (4.8)$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_i (x_i - \bar{x})^2 = \left(\frac{1}{N} \sum_i x_i^2 \right) - (\bar{x})^2 \quad (4.9)$$



最大后验概率估计

■ 目标

- 求后验概率 $p(\mu, \Sigma | \mathcal{D}) \propto p(\mathcal{D} | \mu, \Sigma) p(\mu, \Sigma)$ 最大化参数 μ, Σ

■ 似然函数

$$p(\mathcal{D} | \mu, \Sigma) = (2\pi)^{-ND/2} |\Sigma|^{-\frac{N}{2}} \exp \left(-\frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu) \right)$$

其中

$$\sum_{i=1}^N (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu) = \text{tr}(\Sigma^{-1} \mathbf{S}_{\bar{\mathbf{x}}}) + N(\bar{\mathbf{x}} - \mu)^T \Sigma^{-1} (\bar{\mathbf{x}} - \mu)$$

$$\mathbf{S}_{\bar{\mathbf{x}}} = \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$



最大后验概率估计

■ 似然函数

$$p(\mathcal{D}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-ND/2} |\boldsymbol{\Sigma}|^{-\frac{N}{2}} \exp\left(-\frac{N}{2}(\boldsymbol{\mu} - \bar{\mathbf{x}})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu} - \bar{\mathbf{x}})\right) \exp\left(-\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}_{\bar{\mathbf{x}}})\right)$$



最大后验概率估计

- 共轭先验 $p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = p(\boldsymbol{\mu}|\boldsymbol{\Sigma})p(\boldsymbol{\Sigma})$

$$\begin{aligned} \text{NIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathbf{m}_0, \kappa_0, \nu_0, \mathbf{S}_0) &\triangleq \\ &\mathcal{N}(\boldsymbol{\mu} | \mathbf{m}_0, \frac{1}{\kappa_0} \boldsymbol{\Sigma}) \times \text{IW}(\boldsymbol{\Sigma} | \mathbf{S}_0, \nu_0) \\ &= \frac{1}{Z_{\text{NIW}}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{\kappa_0}{2} (\boldsymbol{\mu} - \mathbf{m}_0)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \mathbf{m}_0)\right) \\ &\quad \times |\boldsymbol{\Sigma}|^{-\frac{\nu_0 + D + 1}{2}} \exp\left(-\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}_0)\right) \\ &= \frac{1}{Z_{\text{NIW}}} |\boldsymbol{\Sigma}|^{-\frac{\nu_0 + D + 2}{2}} \\ &\quad \times \exp\left(-\frac{\kappa_0}{2} (\boldsymbol{\mu} - \mathbf{m}_0)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \mathbf{m}_0) - \frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}_0)\right) \\ Z_{\text{NIW}} &= 2^{\nu_0 D / 2} \Gamma_D(\nu_0 / 2) (2\pi / \kappa_0)^{D/2} |\mathbf{S}_0|^{-\nu_0 / 2} \end{aligned}$$



最大后验概率估计

- 共轭先验 $p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = p(\boldsymbol{\mu}|\boldsymbol{\Sigma})p(\boldsymbol{\Sigma})$
 - \mathbf{m}_0 : $\boldsymbol{\mu}$ 的先验均值
 - 设 $\mathbf{m}_0 = \bar{\mathbf{x}}$
 - κ_0 : 对 \mathbf{m}_0 的信任程度
 - 通常设 κ_0 为一个极小的数, 如 $\kappa_0 = 0.01$
 - \mathbf{S}_0 : $\boldsymbol{\Sigma}$ 的先验均值
 - 设 $\mathbf{S}_0 = \text{diag}(\mathbf{S}_{\bar{\mathbf{x}}})/N$
 - ν_0 : 对 \mathbf{S}_0 的信任程度
 - 通常设 $\nu_0 = D + 2$



最大后验概率估计

■ 后验分布 $p(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathcal{D})$

$$\begin{aligned} p(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathcal{D}) &= \text{NIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathbf{m}_N, \kappa_N, \nu_N, \mathbf{S}_N) \\ \mathbf{m}_N &= \frac{\kappa_0 \mathbf{m}_0 + N \bar{\mathbf{x}}}{\kappa_N} = \frac{\kappa_0}{\kappa_0 + N} \mathbf{m}_0 + \frac{N}{\kappa_0 + N} \bar{\mathbf{x}} \\ \kappa_N &= \kappa_0 + N \\ \nu_N &= \nu_0 + N \\ \mathbf{S}_N &= \mathbf{S}_0 + \mathbf{S}_{\bar{\mathbf{x}}} + \frac{\kappa_0 N}{\kappa_0 + N} (\bar{\mathbf{x}} - \mathbf{m}_0)(\bar{\mathbf{x}} - \mathbf{m}_0)^T \\ &= \mathbf{S}_0 + \mathbf{S} + \kappa_0 \mathbf{m}_0 \mathbf{m}_0^T - \kappa_N \mathbf{m}_N \mathbf{m}_N^T \end{aligned}$$

- 后验均值 \mathbf{m}_N ：先验均值 \mathbf{m}_0 和 MLE 的凸组合
- 后验散度矩阵 \mathbf{S}_N ：先验散度矩阵 \mathbf{S}_0 + 经验散度矩阵 $\mathbf{S}_{\bar{\mathbf{x}}}$ + 偏移项



最大后验概率估计

- 最大化后验分布 $p(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathcal{D})$
 - 后验众数 (mode)

$$\operatorname{argmax} p(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathcal{D}) = \left(\mathbf{m}_N, \frac{\mathbf{S}_N}{\nu_N + D + 2} \right)$$

- 如果设置 $\kappa_0 = 0$

$$\operatorname{argmax} p(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathcal{D}) = \left(\bar{\mathbf{x}}, \frac{\mathbf{S}_0 + \mathbf{S}_{\bar{\mathbf{x}}}}{\nu_0 + N + D + 2} \right)$$



高斯判别分析

■ 高斯判别分析 (GDA)

- 假设每个类 c 的样本服从多元高斯分布

$$p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$$

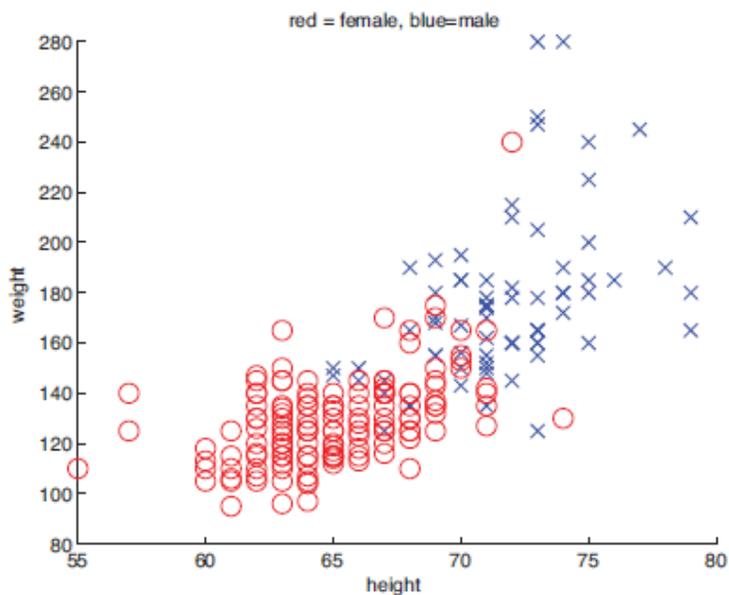
- 给定样本 \mathbf{x} , 生成分类器预测最可能的类别

$$\begin{aligned}\hat{y}(\mathbf{x}) &= \operatorname{argmax}_c [\log p(y = c|\boldsymbol{\pi}) + \log p(\mathbf{x}|\boldsymbol{\theta}_c)] \\ &= \operatorname{argmax}_c [\log p(y = c|\boldsymbol{\pi}) + \log \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)]\end{aligned}$$

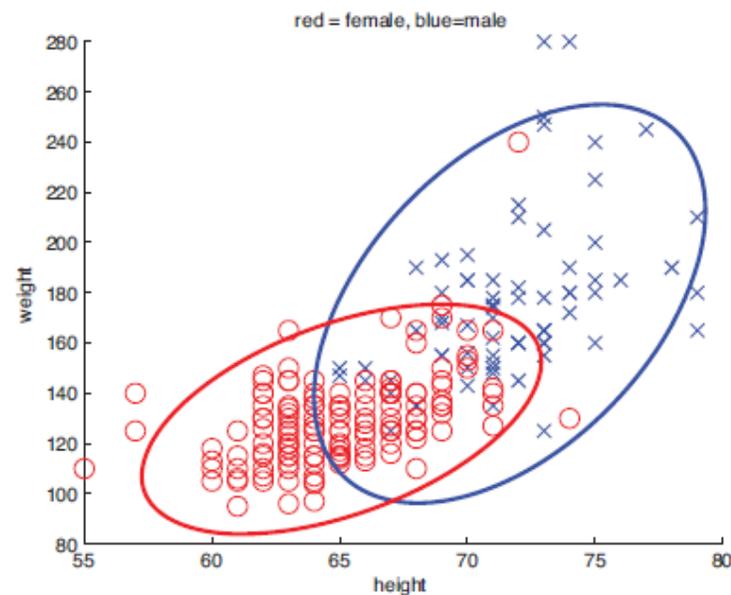
高斯判别分析



- 左图：高度/体重数据
- 右图：每个类的可视化



(a)



(b)



高斯判别分析

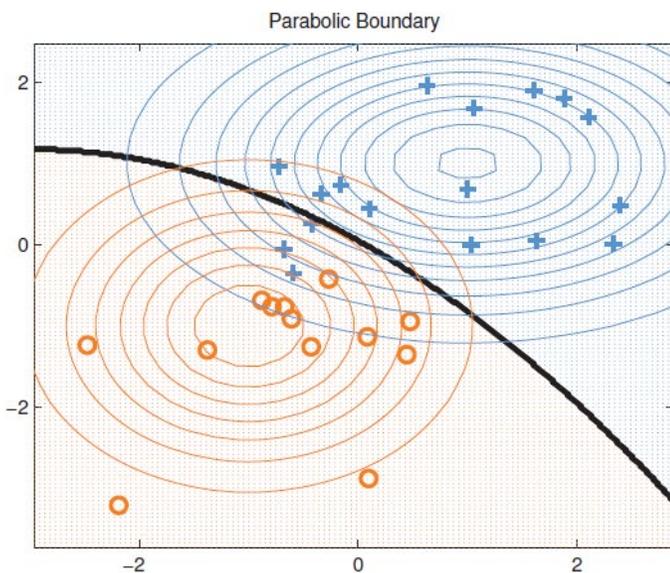
- 生成分类器
 - 如果类先验 π_c 是均匀的

$$\begin{aligned}\hat{y}(\mathbf{x}) &= \operatorname{argmax}_c [\log p(y = c | \boldsymbol{\pi}_c) + \log \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)] \\ &= \operatorname{argmax}_c \log \left\{ \frac{1}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}_c|^{\frac{1}{2}}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \right] \right\} \\ &= \operatorname{argmax}_c -\frac{1}{2} \log |\boldsymbol{\Sigma}_c| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \\ &= \operatorname{argmin}_c (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) + \log |\boldsymbol{\Sigma}_c|\end{aligned}$$

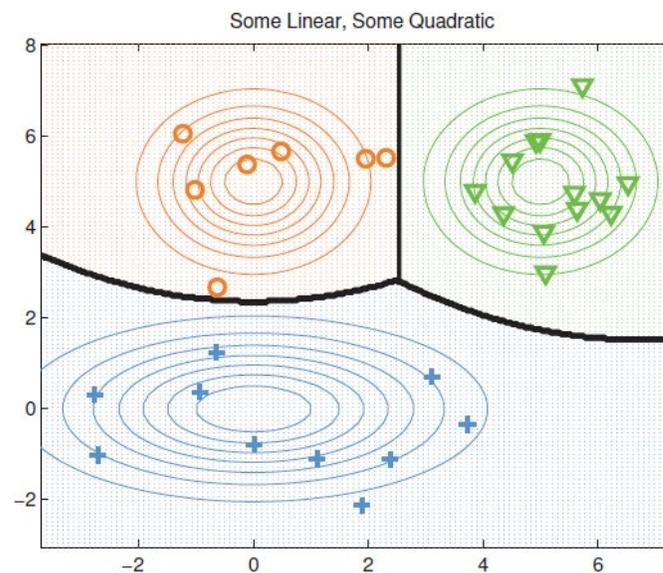
高斯判别分析



- 高斯判别分析又称为**二次判别分析 (QDA)**
 - 判别边界 $(\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) + \log |\boldsymbol{\Sigma}_c|$ 是二次形式



(a)

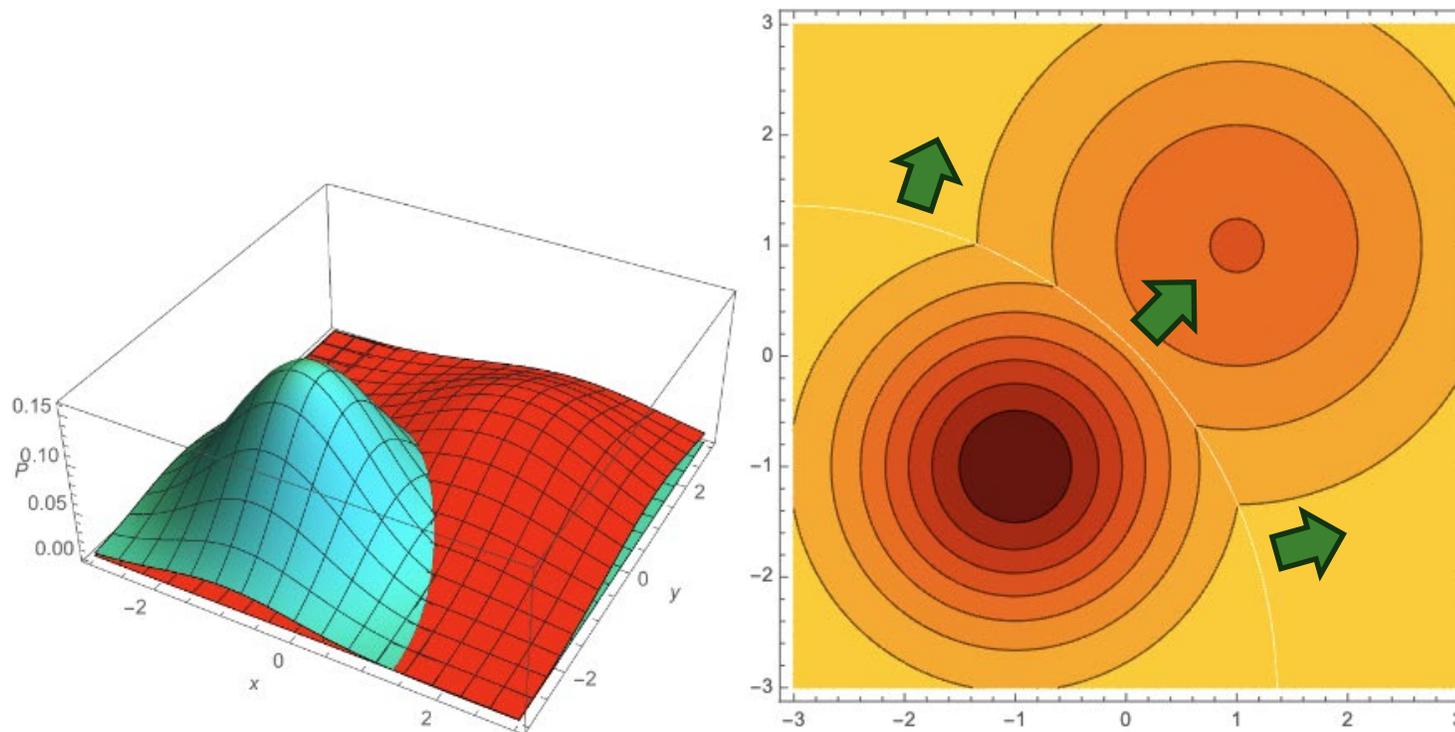


(b)

高斯判别分析



- 判别边界：两个类密度分布的“等势面”





高斯判别分析

■ 特例： $\Sigma_c = \Sigma$

$$\begin{aligned} p(y = c | \mathbf{x}, \boldsymbol{\theta}) &\propto \pi_c \exp \left[\boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c \right] \\ &= \exp \left[\boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c \right] \exp \left[-\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} \right] \end{aligned}$$

□ 定义

$$\begin{aligned} \gamma_c &= -\frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c \\ \boldsymbol{\beta}_c &= \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c \end{aligned}$$

与类别c无关的二次项，
判别时可以忽略



高斯判别分析

- 特例： $\Sigma_c = \Sigma$
 - 生成分类器

$$p(y = c | \mathbf{x}, \boldsymbol{\theta}) \propto \exp[\beta_c^T \mathbf{x} + \gamma_c]$$

- 任意两个类 c 和 c' 的决策边界恰好是一条直线

$$\begin{aligned} p(y = c | x, \theta) = p(y = c' | x, \theta) &\implies \beta_c^T x + \gamma_c = \beta_{c'}^T x + \gamma_{c'} \\ &\implies (\beta_c^T - \beta_{c'}^T)x = \gamma_{c'} - \gamma_c \end{aligned}$$

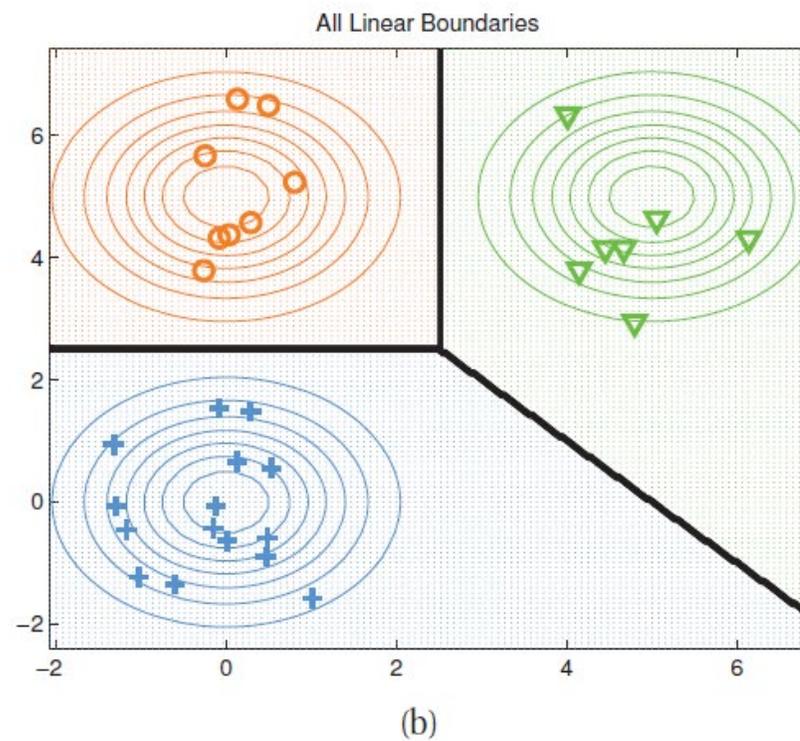
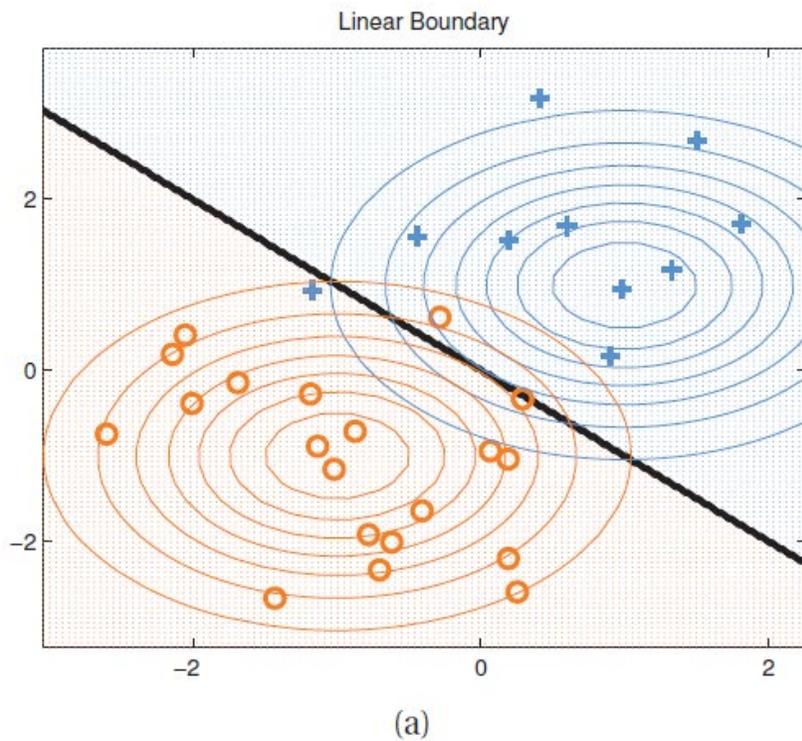
- 亦称为线性判别分析 (linear discriminant analysis, LDA)

$$\log p(y = c | \mathbf{x}, \boldsymbol{\theta}) \propto \beta_c^T \mathbf{x} + \gamma_c$$

高斯判别分析



- 特例: $\Sigma_c = \Sigma$



高斯判别分析



■ 对数似然函数

$$\log p(\mathcal{D}|\boldsymbol{\theta}) = \left[\sum_{i=1}^N \sum_{c=1}^C \mathbb{I}(y_i = c) \log \pi_c \right] + \sum_{c=1}^C \left[\sum_{i:y_i=c} \log \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \right]$$

□ π_c : 第c个类的先验

$$0 \leq \pi_c \leq 1$$

$$\sum_{c=1}^C \pi_c = 1$$

□ $\theta_c = (\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$: 第c个类的高斯分布参数



高斯判别分析

■ 极大似然估计

- π_c : 第c个类的先验

$$\hat{\pi}_c = \frac{N_c}{N}$$

- $\theta_c = (\mu_c, \Sigma_c)$: 第c个类的高斯分布参数

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{i:y_i=c} \mathbf{x}_i$$

$$\hat{\Sigma}_c = \frac{1}{N_c} \sum_{i:y_i=c} (\mathbf{x}_i - \hat{\mu}_c)(\mathbf{x}_i - \hat{\mu}_c)^T$$

目录

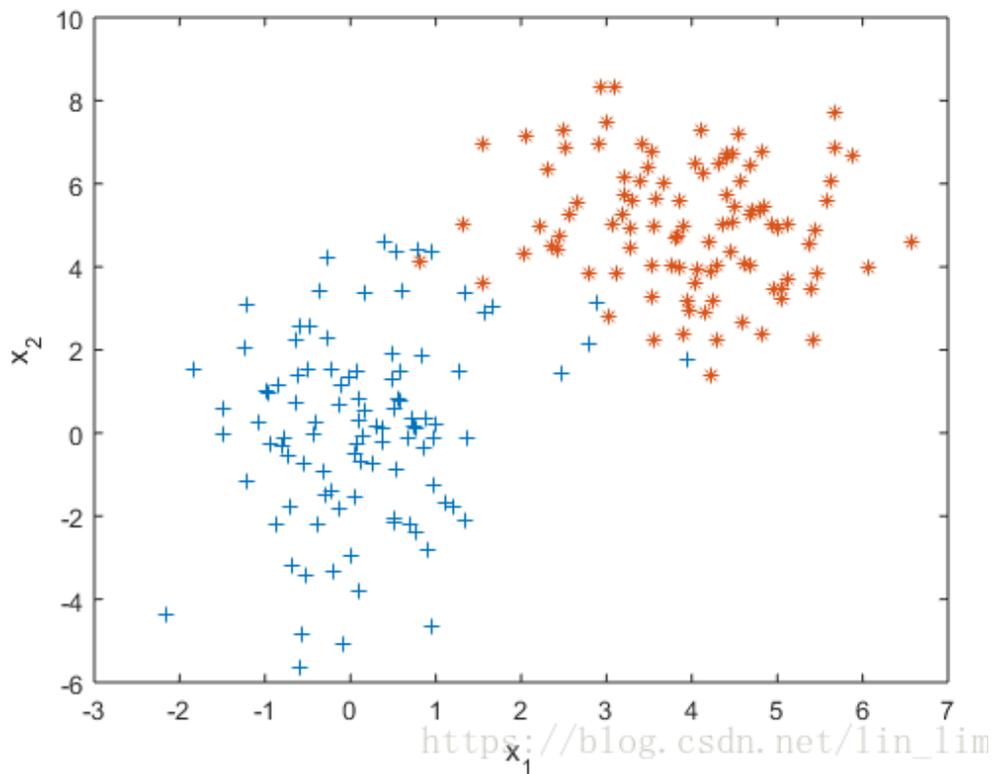


- 生成分类器
- 单高斯模型
- 高斯混合模型
- 隐马尔可夫模型

为什么有高斯混合模型

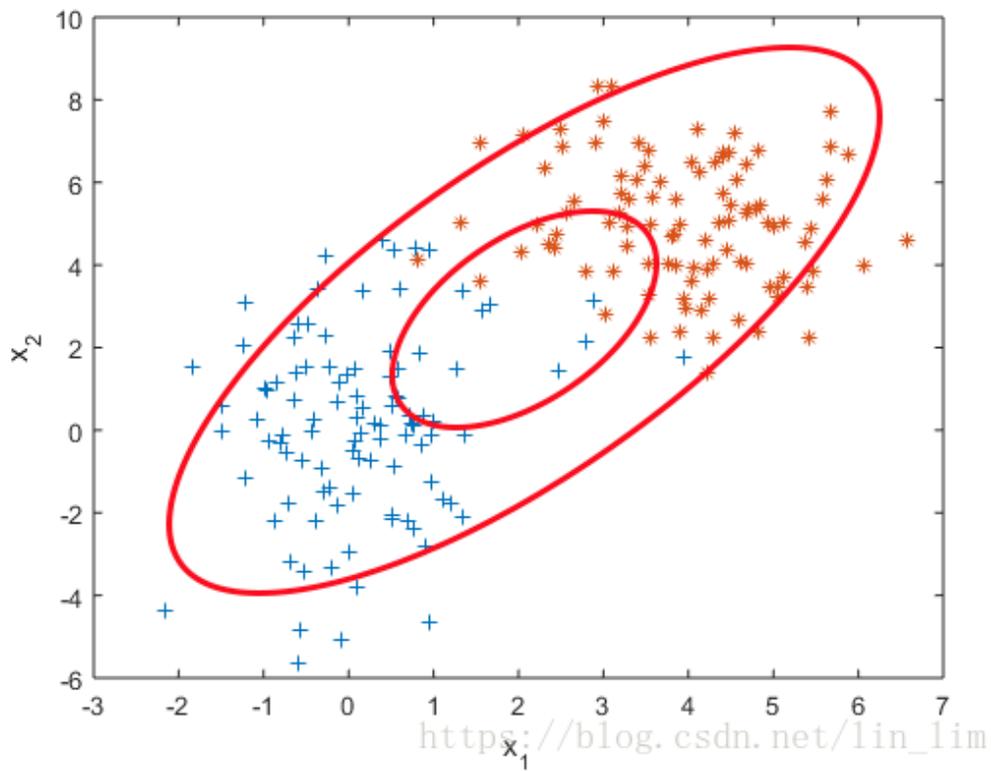


- 假设这组数据是由某个高斯分布产生



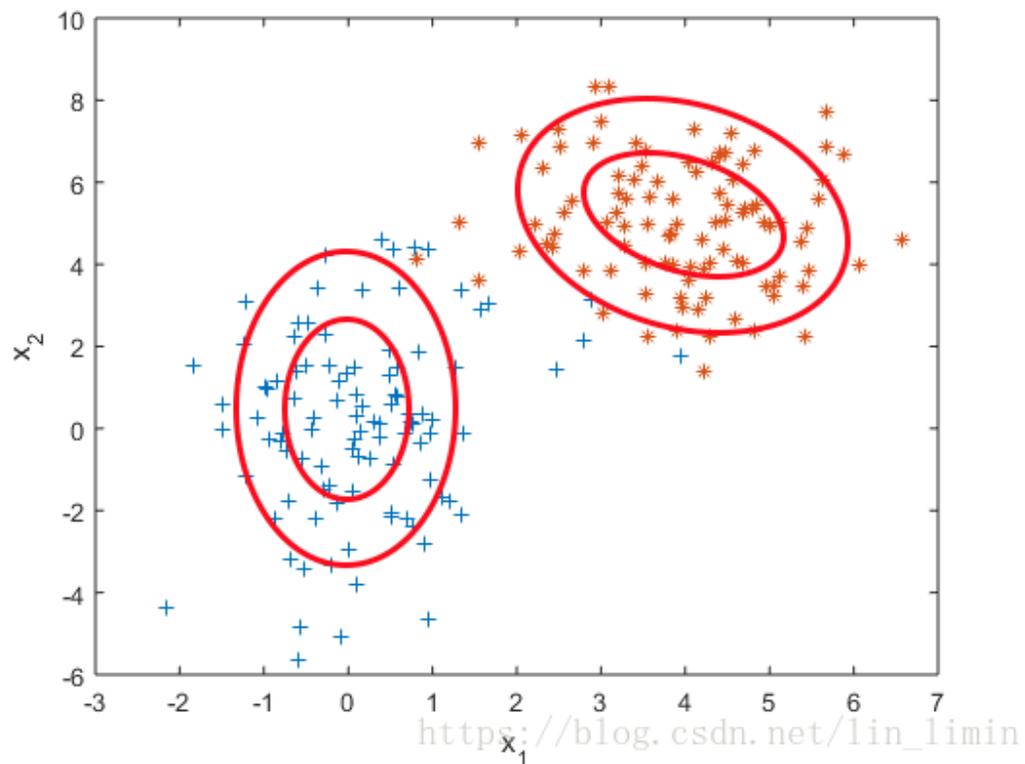
为什么有高斯混合模型

- 样本服从单高斯分布的**假设并不合理**



为什么有高斯混合模型

- 这是用两个不同的高斯分布模型产生的数据





高斯混合模型

- 假设混合高斯模型由 K 个高斯模型组成，其概率密度为

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\theta}) &= \sum_z p(z|\boldsymbol{\theta})p(\mathbf{x}|z, \boldsymbol{\theta}) \\ &= \sum_{k=1}^K p(z = k|\boldsymbol{\theta})p(\mathbf{x}|z = k, \boldsymbol{\theta}) \\ &= \sum_{k=1}^K \pi_k p_k(\mathbf{x}|\boldsymbol{\theta}_k) \end{aligned}$$

- 隐含类别 $z \in \{1, \dots, K\}$: 确定数据 \mathbf{x} 是由哪个高斯模型产生的

高斯混合模型



- 高斯混合模型参数为 θ

$$\theta = (\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K)$$

- π_k : 第k个类的混合权重

$$0 \leq \pi_k \leq 1$$

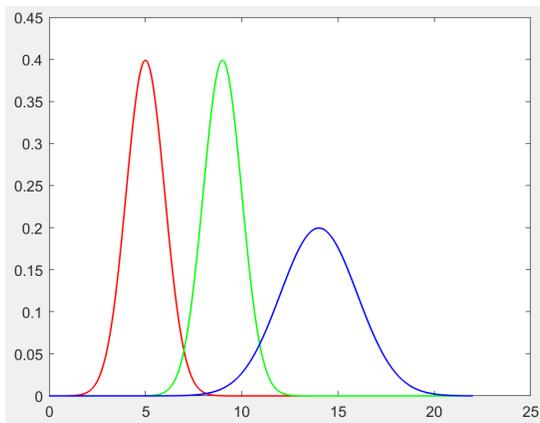
$$\sum_{k=1}^K \pi_k = 1$$

- θ_k : 第k个单高斯模型参数（均值、协方差），概率密度函数为 p_k

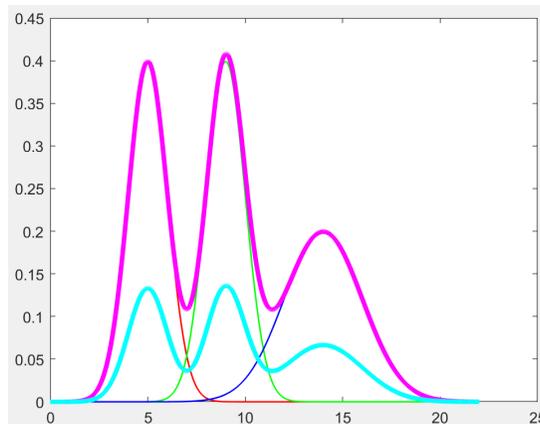
$$p_k(x|\theta_k) \sim \mathcal{N}(x|\mu_k, \Sigma_k)$$

高斯混合模型

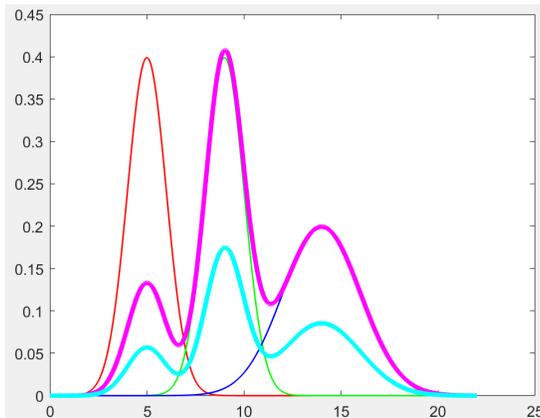
■ 一维的例子



红线: $N(5,1)$
 绿线: $N(9,1)$
 蓝线: $N(14,2)$



品红: $N(5,1)+N(9,1)+N(14,2)$
 淡蓝: $(1/3)*N(5,1)+ (1/3)*N(9,1)+ (1/3)*N(14,2)$

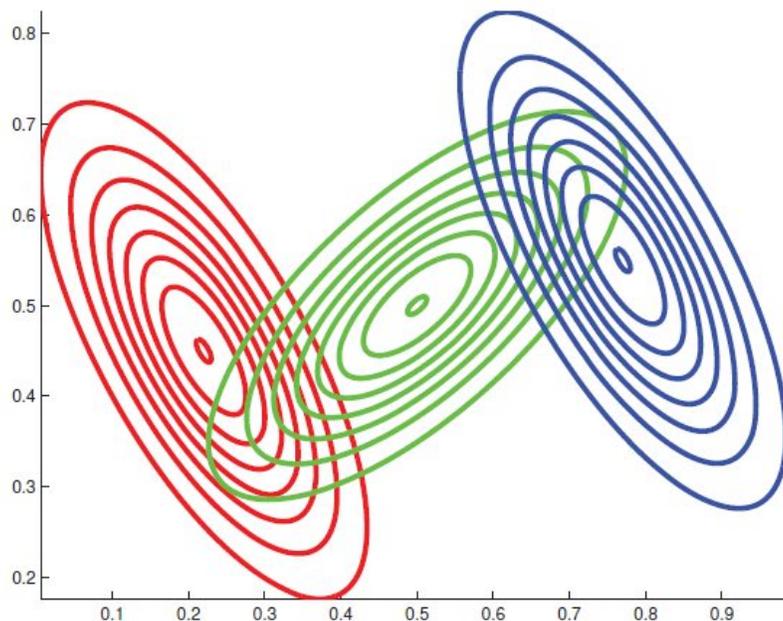


品红: $(1/3)*N(5,1)+ N(9,1)+ N(14,2)$
 淡蓝: $(1/7)*N(5,1)+ (3/7)*N(9,1) + (3/7)*N(14,2)$

二维高斯混合模型



- 直观上，高斯混合模型可以模拟任意分布函数





高斯混合模型

- 高斯混合模型可用于聚类

- 响应度 (responsibility) r_{ik}

- $p(z_i = k | \mathbf{x}_i, \theta)$: 第 i 个样本属于第 k 个类的后验概率

$$r_{ik} \triangleq p(z_i = k | \mathbf{x}_i, \theta) = \frac{p(z_i = k | \theta) p(\mathbf{x}_i | z_i = k, \theta)}{\sum_{k'=1}^K p(z_i = k' | \theta) p(\mathbf{x}_i | z_i = k', \theta)}$$

- 最大后验概率估计 MAP

$$z_i^* = \arg \max_k r_{ik} = \arg \max_k \log p(\mathbf{x}_i | z_i = k, \theta) + \log p(z_i = k | \theta)$$

- 等价于使用生成分类器的计算公式

高斯混合模型

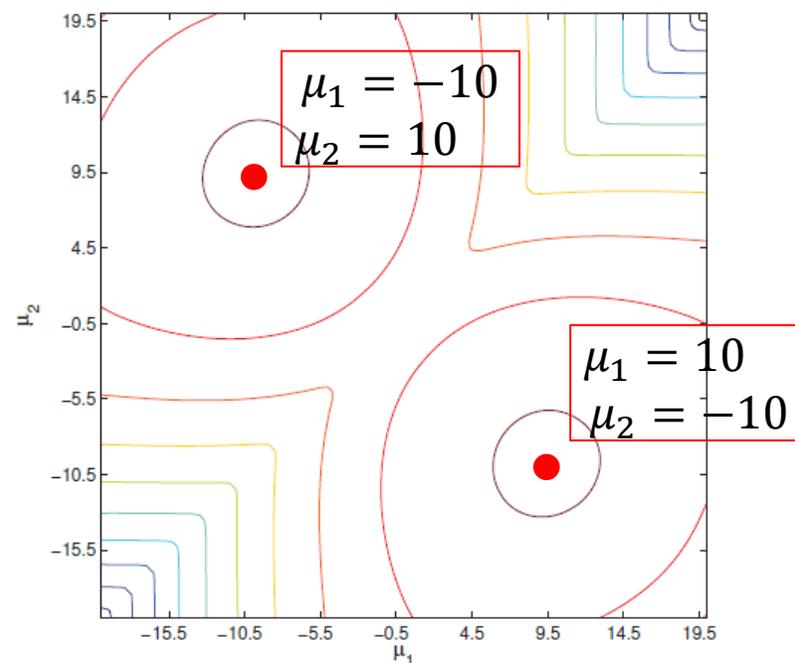
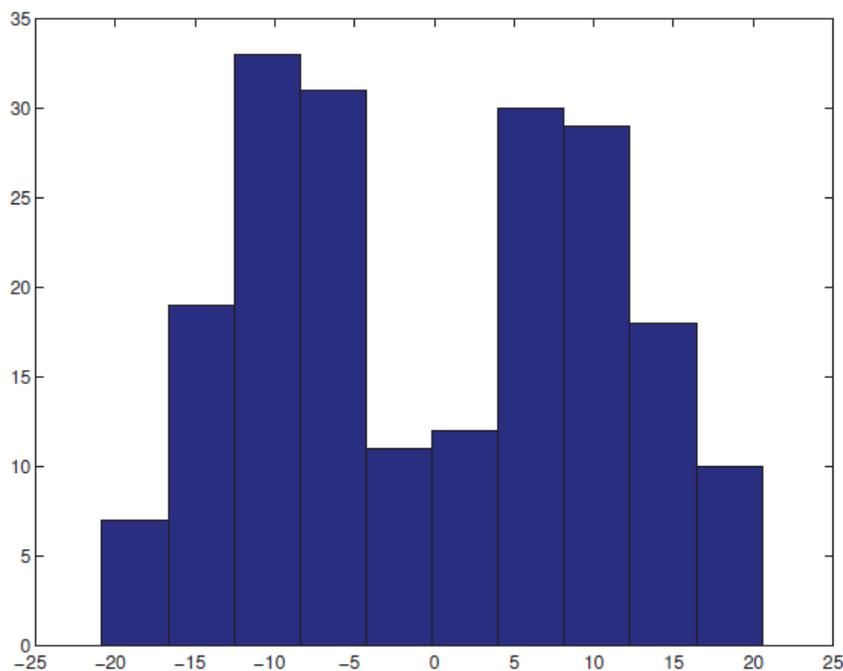


- 与生成分类器的差别
 - **训练**阶段是否有类标签
 - 分类：已知类别信息 y_i
 - 聚类：没有类别信息，使用隐变量 z_i

参数估计



- 左图：2个一维高斯模型混合，采样200个点
- 右图：似然函数表面 $\log p(\mathcal{D}|\mu_1, \mu_2)$



参数估计

- 极大似然估计 $\max_{\theta} \log p(\mathcal{D}|\theta)$

$$p(\mathcal{D}|\theta) = \prod_{i=1}^N \sum_{k=1}^K \pi_k \mathcal{N}(x_i|\mu_k, \Sigma_k)$$

$$p(x_i, z_i|\theta) = p(z_i|\theta)p(x_i|z_i, \theta)$$

$$p(z_i = k|\theta) = \pi_k, \quad p(x_i|z_i = k, \theta) = \mathcal{N}(x_i|\mu_k, \Sigma_k)$$



$$\log p(\mathcal{D}|\theta) = \sum_i \log \left[\sum_{z_i} p(\mathbf{x}_i, \mathbf{z}_i|\theta) \right]$$

- 最大后验概率估计 $\max_{\theta} p(\theta|\mathcal{D})$

$$\log p(\theta|\mathcal{D}) \propto \log p(\mathcal{D}|\theta) + \log p(\theta)$$

- 注意：这两个参数估计的解不唯一
可能有 $K!$ 种不同的解

EM算法



■ E步: Q函数

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t-1)}) &\triangleq \mathbb{E} \left[\sum_i \log p(\mathbf{x}_i, z_i | \boldsymbol{\theta}) \right] \\ &= \sum_i \mathbb{E} \left[\log \left[\prod_{k=1}^K (\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k))^{\mathbb{I}(z_i=k)} \right] \right] \\ &= \sum_i \sum_k \mathbb{E} [\mathbb{I}(z_i = k)] \log [\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k)] \\ &= \sum_i \sum_k p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}^{t-1}) \log [\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k)] \\ &= \sum_i \sum_k r_{ik} \log \pi_k + \sum_i \sum_k r_{ik} \log p(\mathbf{x}_i | \boldsymbol{\theta}_k) \end{aligned}$$



EM算法

■ E步：Q函数

□ 计算

$$r_{ik} \triangleq p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}^{(t-1)})$$

$$r_{ik} = \frac{\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k^{(t-1)})}{\sum_{k'} \pi_{k'} p(\mathbf{x}_i | \boldsymbol{\theta}_{k'}^{(t-1)})}$$

■ M步：最大化Q函数

□ 类先验

$$\pi_k = \frac{1}{N} \sum_i r_{ik} = \frac{r_k}{N}$$

- r_k ：第k个类的加权数据量

EM算法

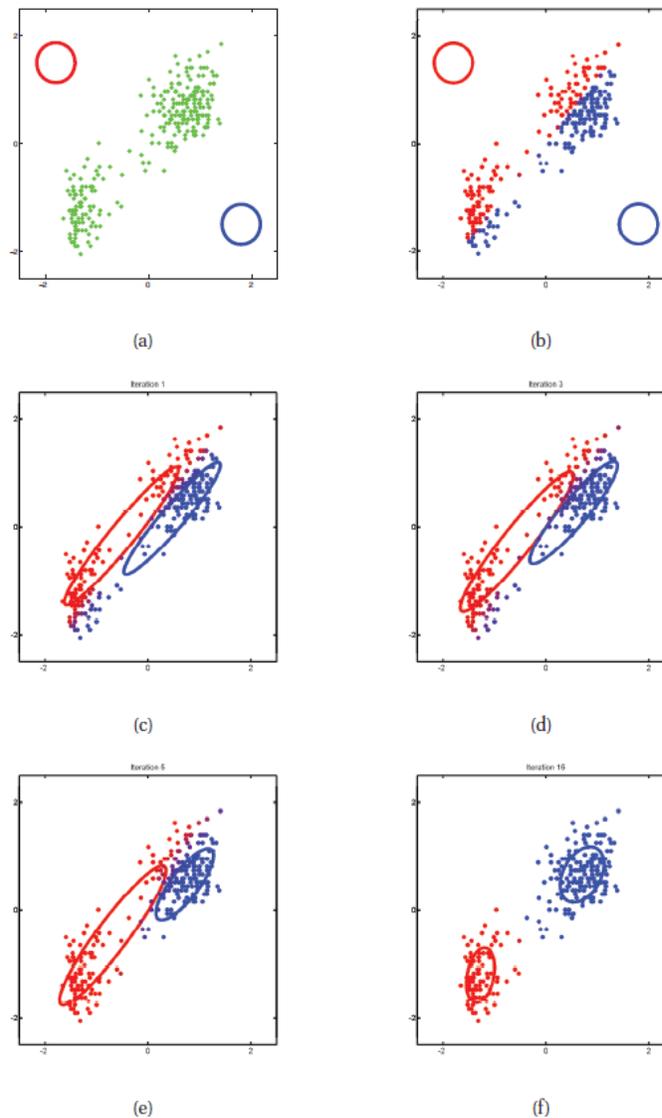
- M步：最大化Q函数
 - 第k个类的高斯模型参数

$$\begin{aligned}\mu_k &= \frac{\sum_i r_{ik} \mathbf{x}_i}{r_k} \\ \Sigma_k &= \frac{\sum_i r_{ik} (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T}{r_k} = \frac{\sum_i r_{ik} \mathbf{x}_i \mathbf{x}_i^T}{r_k} - \mu_k \mu_k^T\end{aligned}$$

- 均值 μ_k ：第k个类的数据加权平均
- 协方差 Σ_k ：第k个类的数据的加权散度矩阵

实例：EM算法

- (a) 原始数据
- (b) 模型参数赋初值
- (c) 第1次迭代以后
- (d) 第3次迭代以后
- (e) 第5次迭代以后
- (f) 第16次迭代以后





特例： K-means算法

- K-means算法是高斯混合模型的EM算法变种
 - 假设
 - 协方差矩阵 $\Sigma_k = \sigma^2 I_D$ 已知
 - 类先验 $\pi_k = 1/K$ 已知
 - 均值 μ_k 未知
 - EM算法
 - E步：delta-函数近似

$$z_i^* = \operatorname{argmax}_k p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta})$$

$$p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}) \approx \mathbb{I}(k = z_i^*)$$



特例： K-means算法

- E步：计算 z_i^*

$$\begin{aligned} & \max_k p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}) \\ &= \max_k \frac{p(z_i = k | \boldsymbol{\theta}) p(\mathbf{x}_i | z_i = k, \boldsymbol{\theta})}{\sum_{k'=1}^K p(z_i = k' | \boldsymbol{\theta}) p(\mathbf{x}_i | z_i = k', \boldsymbol{\theta})} \\ &= \max_k p(z_i = k | \boldsymbol{\theta}) p(\mathbf{x}_i | z_i = k, \boldsymbol{\theta}) \\ &= \max_k \frac{1}{K} \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x}_i - \mu_k)^T \Sigma^{-1} (\mathbf{x}_i - \mu_k) \right] \\ &= \min_k (\mathbf{x}_i - \mu_k)^T (\mathbf{x}_i - \mu_k) \end{aligned}$$



特例： K-means算法

- E步：计算 z_i^*

$$\begin{aligned} z_i^* &= \operatorname{argmin}_k (\mathbf{x}_i - \mu_k)^T (\mathbf{x}_i - \mu_k) \\ &= \operatorname{argmin}_k \|\mathbf{x}_i - \mu_k\|_2^2 \end{aligned}$$

- M步：计算均值 μ_k

$$\mu_k = \frac{\sum_i r_{ik} \mathbf{x}_i}{r_k} = \frac{1}{N_k} \sum_{i:z_i=k} \mathbf{x}_i$$



特例：K-means算法

■ 算法伪代码

Algorithm 11.1: K-means algorithm

```
1 initialize  $\mu_k$  ;
2 repeat
3   Assign each data point to its closest cluster center:  $z_i = \arg \min_k \|\mathbf{x}_i - \mu_k\|_2^2$ ;
4   Update each cluster center by computing the mean of all points assigned to it:
   
$$\mu_k = \frac{1}{N_k} \sum_{i:z_i=k} \mathbf{x}_i$$

5 until converged;
```



特例：K-means算法

- 高斯混合模型相比于K-means算法的**优点**:
 - 计算一个样本属于某类的概率值；
 - 不仅仅可以用于聚类，还可以用于概率密度估计；
 - 可以用于生成新的样本点

极大似然估计：过拟合现象

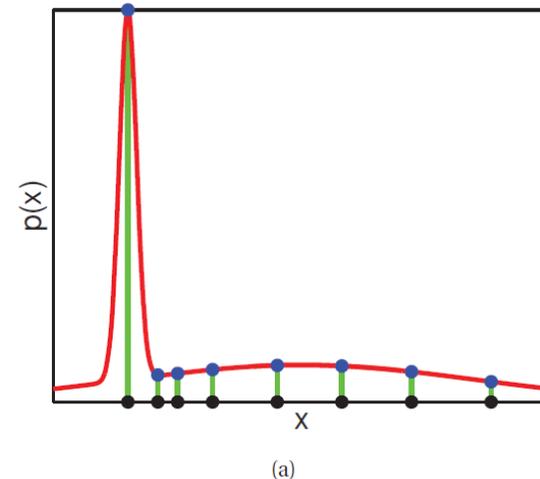
■ 实例

□ $K=2$

□ 协方差矩阵 $\Sigma_k = \sigma_k^2 I_D$

□ EM算法初始化: $\mu_2 = \mathbf{x}_1$

$$\log p(\mathcal{D}|\theta) = \sum_i \log \left[\sum_{z_i} p(\mathbf{x}_i, z_i|\theta) \right]$$



$$\begin{aligned} \mathcal{N}(\mathbf{x}_1|\mu_2, \sigma_2^2 I_D) &= \frac{1}{(2\pi)^{D/2} |\sigma_2^2 I_D|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x}_1 - \mu_2)^T (\sigma_2^2 I_D)^{-1} (\mathbf{x}_1 - \mu_2) \right] \\ &= \frac{1}{(2\pi\sigma_2^2)^{D/2}} \end{aligned}$$

■ 当 $\sigma_2 \rightarrow 0$ 时, $\mathcal{N}(\mathbf{x}_1|\mu_2, \sigma_2^2 I_D) \rightarrow \infty$



最大后验概率估计

- 后验概率

$$\log p(\boldsymbol{\theta}|\mathcal{D}) \propto \log p(\mathcal{D}|\boldsymbol{\theta}) + \log p(\boldsymbol{\theta})$$

- EM算法

- E步:

$$Q'(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t-1)}) = \left[\sum_i \sum_k r_{ik} \log \pi_{:k} + \sum_i \sum_k r_{ik} \log p(\mathbf{x}_i | \boldsymbol{\theta}_k) \right] + \log p(\boldsymbol{\pi}) + \sum_k \log p(\boldsymbol{\theta}_k)$$

- 计算

$$r_{ik} \triangleq p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}^{(t-1)})$$

$$r_{ik} = \frac{\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k^{(t-1)})}{\sum_{k'} \pi_{k'} p(\mathbf{x}_i | \boldsymbol{\theta}_{k'}^{(t-1)})}$$



最大后验概率估计

$$Q'(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t-1)}) = \left[\sum_i \sum_k r_{ik} \log \pi_{ik} + \sum_i \sum_k r_{ik} \log p(\mathbf{x}_i | \boldsymbol{\theta}_k) \right] - \log p(\boldsymbol{\pi}) + \sum_k \log p(\boldsymbol{\theta}_k)$$

□ M步:

- 类先验 $\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\alpha})$

$$\pi_k = \frac{r_k + \alpha_k - 1}{N + \sum_k \alpha_k - K}$$

- 第k个高斯模型的先验

$$p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \text{NIW}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k | \mathbf{m}_0, \kappa_0, \nu_0, \mathbf{S}_0)$$

$$\begin{aligned} \hat{\boldsymbol{\mu}}_k &= \frac{r_k \bar{\mathbf{x}}_k + \kappa_0 \mathbf{m}_0}{r_k + \kappa_0} \\ \bar{\mathbf{x}}_k &\triangleq \frac{\sum_i r_{ik} \mathbf{x}_i}{r_k} \\ \hat{\boldsymbol{\Sigma}}_k &= \frac{\mathbf{S}_0 + \mathbf{S}_k + \frac{\kappa_0 r_k}{\kappa_0 + r_k} (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T}{\nu_0 + r_k + D + 2} \\ \mathbf{S}_k &\triangleq \sum_i r_{ik} (\mathbf{x}_i - \bar{\mathbf{x}}_k)(\mathbf{x}_i - \bar{\mathbf{x}}_k)^T \end{aligned}$$



最大后验概率估计

■ 超参数

□ 类先验

$$\pi \sim \text{Dir}(\alpha)$$

$$\alpha_k = 1, k = 1, \dots, K$$

□ 第k个高斯模型的先验

■ $\kappa_0 = 0$

$$\hat{\mu}_k = \frac{\sum_i r_{ik} \mathbf{x}_i}{r_k}$$

$$\mathbf{S}_0 = \frac{1}{K^{1/D}} \text{diag}(s_1^2, \dots, s_D^2)$$

$$\hat{\Sigma}_k = \frac{\mathbf{S}_0 + \mathbf{S}_k}{\nu_0 + r_k + D + 2}$$

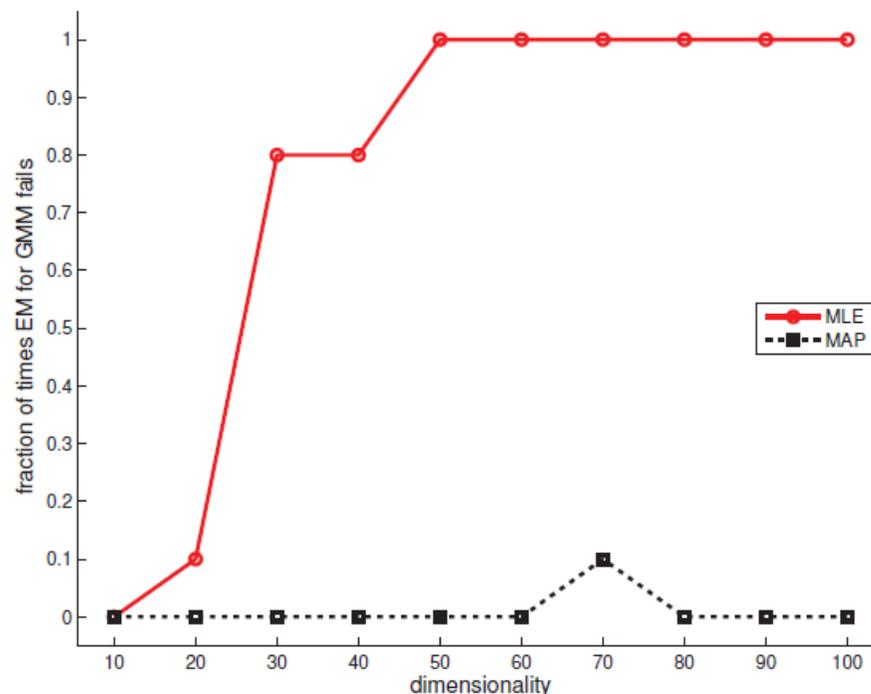
$$s_j = (1/N) \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2$$

■ $\nu_0 = D + 2$

MLE vs. MAP



- 比较极大似然估计和最大后验概率估计实验失败的次数
 - 使用合成数据
 - $N=100$
 - 5次重复实验



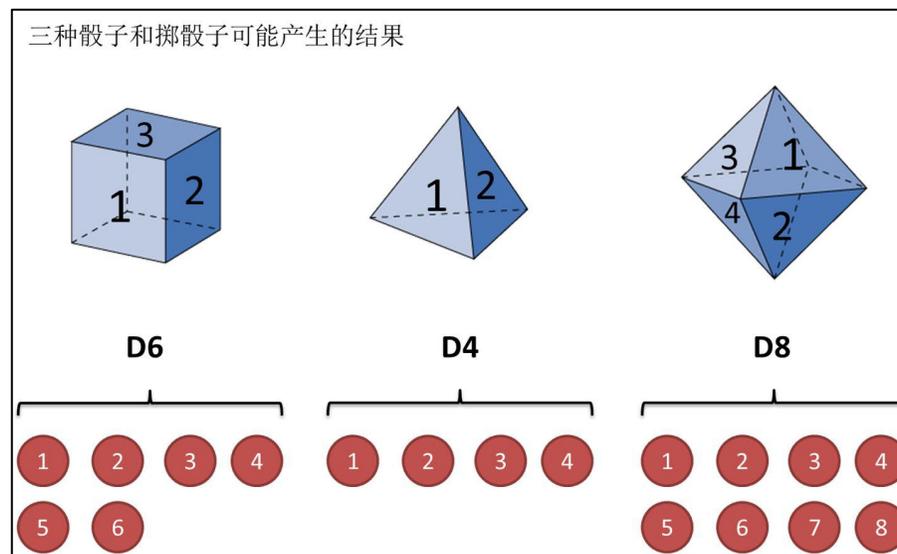
目录



- 生成分类器
- 单高斯模型
- 高斯混合模型
- 隐马尔可夫模型

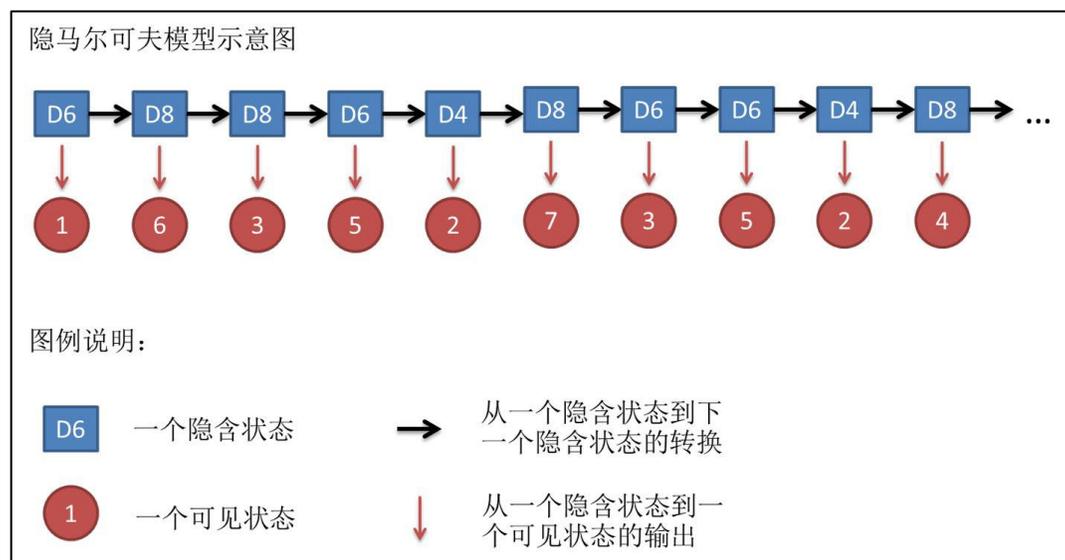
实例：扔骰子游戏

- 假设：有3个骰子 D6、D4、D8
 - 一开始，随机任选一个骰子；之后每次的骰子类型，依赖于上一次选的骰子
 - 掷骰子，得到1~8中的一个数字
- 重复N (N=10) 次，得到一组数字1 6 3 5 2 7 3 5 2 4



实例：扔骰子游戏

- 可见状态链 x ：投掷结果序列
- 隐含状态链 z ：选用骰子的序列



实例：扔骰子游戏

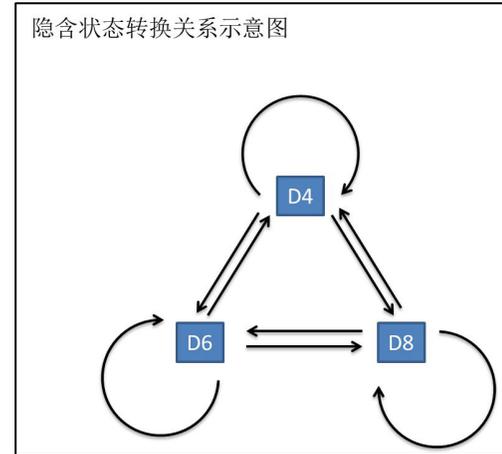
- 可用两种概率关系描述
 - 隐含状态之间：状态转移概率
 - 隐含与可见状态之间：观测概率
 - E.g., 一次投掷，获得1的概率

$$p(\mathbf{x} = 1) = \sum_z p(z)p(\mathbf{x}|z) = \frac{1}{3} \times \frac{1}{6} + \frac{1}{3} \times \frac{1}{4} + \frac{1}{3} \times \frac{1}{8}$$

- E.g., 两次投掷，获得1、6的概率

$$p(\mathbf{x} = 1\ 6) = p(\mathbf{x}_1 = 1, \mathbf{x}_2 = 6) \\ = \sum_{z_1} \sum_{z_2} p(z_1)p(\mathbf{x}_1 = 1|z_1)p(z_2|z_1)p(\mathbf{x}_2 = 6|z_2)$$

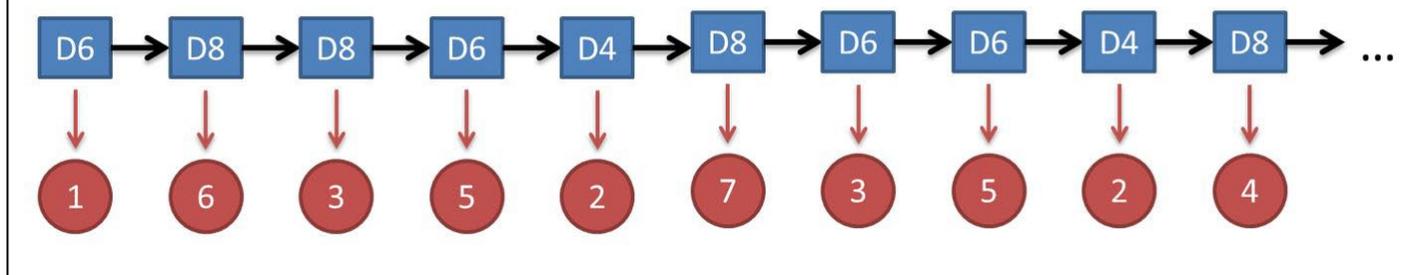
隐状态的切换



隐马尔可夫模型

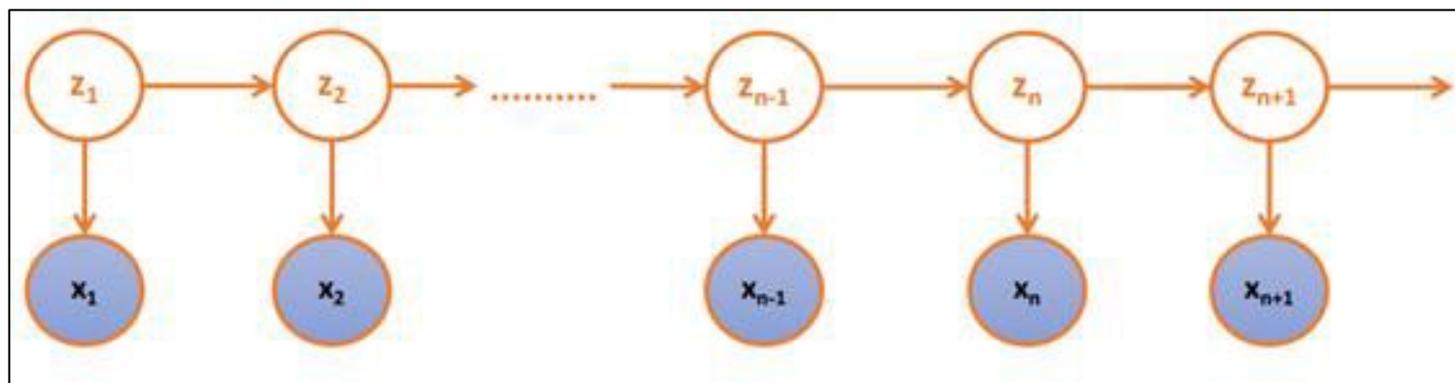
- 隐马尔可夫模型 (Hidden Markov Model, HMM)
 - 关于**时序**的概率模型
 - 描述由一个隐藏的马尔可夫链随机生成不可观测的**状态随机序列**，再由各个状态生成一个观测，从而产生**观测随机序列**的过程
 - 序列的每个位置可看作一个**时刻**

隐马尔可夫模型示意图



隐马尔可夫模型

- HMM的概率图模型表示
 - 隐变量 z , 观测变量 x





隐马尔可夫模型

- HMM模型参数: $\theta = (\pi, A, B)$

- 隐含状态的初始概率分布 π
- 隐含状态的状态转移概率矩阵A

$$A(i, j) = p(z_t = j | z_{t-1} = i)$$

- 观测概率矩阵B

$$p(\mathbf{x}_t | z_t = j)$$

- 依赖关系

- 初始状态概率向量 π 和状态转移概率矩阵A决定状态序列 z
- 观测概率矩阵B决定观测序列 x

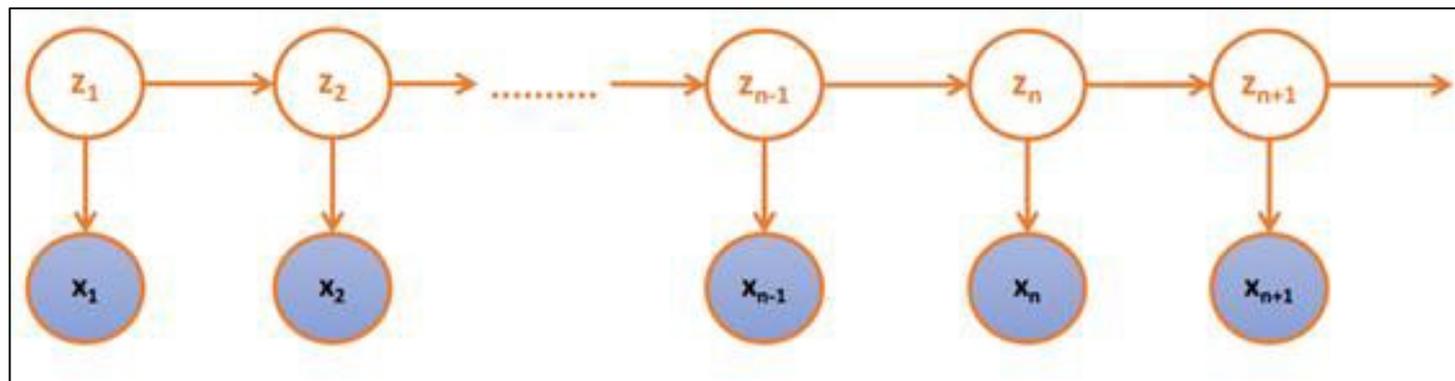
隐马尔可夫模型

■ 齐次马尔可夫性假设

- 隐马尔可夫链 t 时刻的状态只和 $t - 1$ 时刻的状态有关
- 与其他时刻的状态及观测无关，也与时刻 t 无关

$$p(z_{1:T}) = p(z_1)p(z_2|z_1)p(z_3|z_1, z_2) \dots p(z_T|z_1, \dots, z_{T-2}, z_{T-1})$$

$$= p(z_1)p(z_2|z_1)p(z_3|z_2) \dots p(z_T|z_{T-1})$$

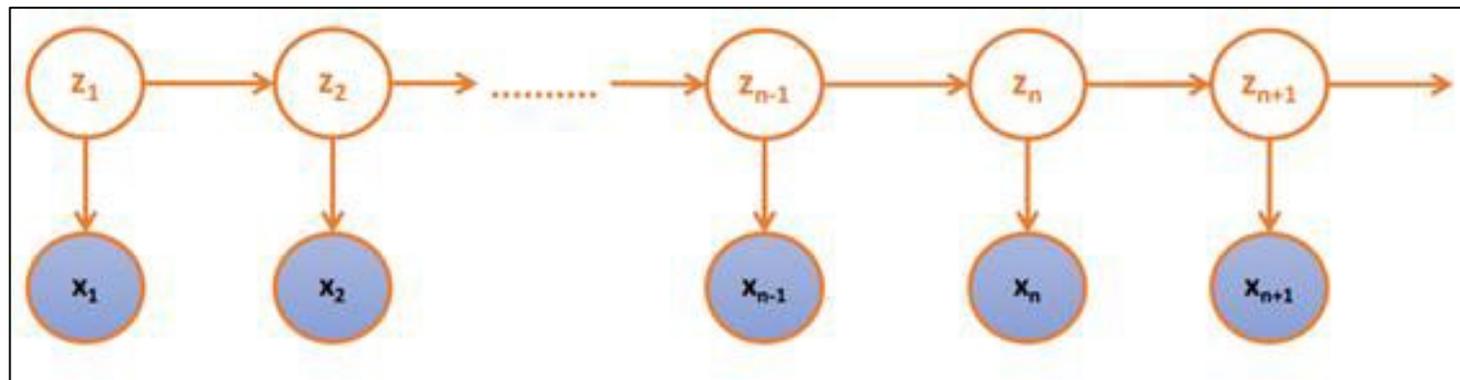


隐马尔可夫模型

■ 观测独立性假设

- 观测变量只和当前时刻的状态有关，与其他时刻的观测和状态均无关

$$p(x_{1:T} | z_{1:T}) = p(x_1 | z_1) p(x_2 | z_2) \dots p(x_T | z_T)$$





隐马尔可夫模型

■ 隐含状态链和观测序列的联合分布

$$p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T}) = p(\mathbf{z}_{1:T})p(\mathbf{x}_{1:T}|\mathbf{z}_{1:T}) = \left[p(z_1) \prod_{t=2}^T p(z_t|z_{t-1}) \right] \left[\prod_{t=1}^T p(\mathbf{x}_t|z_t) \right]$$

□ 离散状态空间

□ 离散或连续的观测变量

■ 离散观测变量：发射概率矩阵

$$p(\mathbf{x}_t = l | z_t = k, \boldsymbol{\theta}) = B(k, l)$$

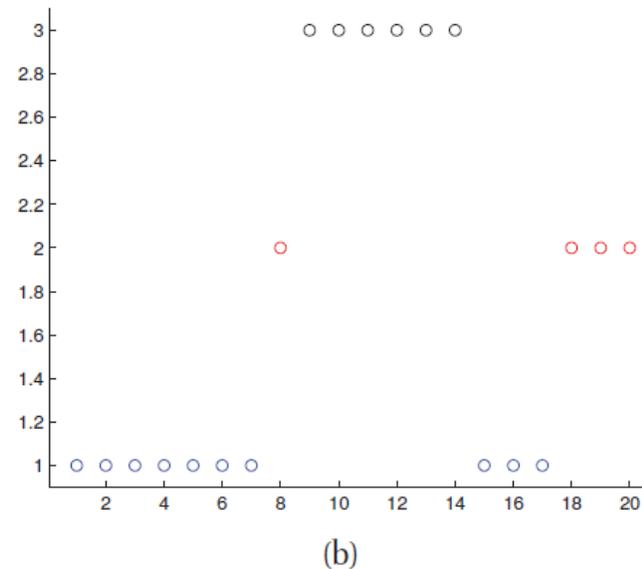
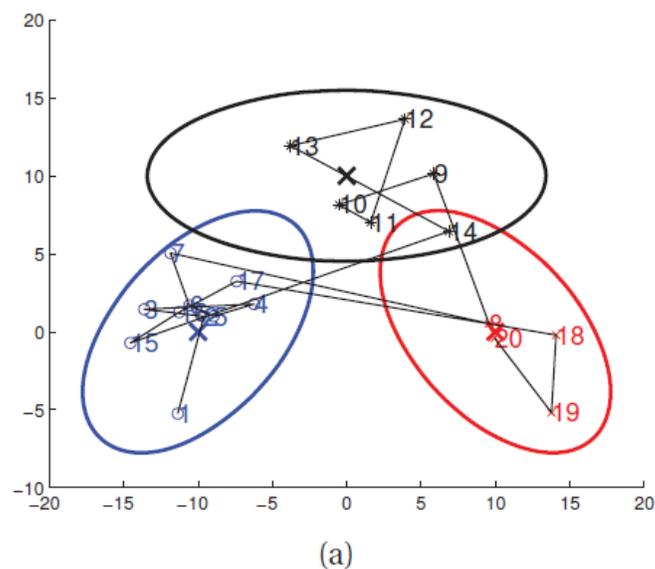
	X	Y	Z
z_1	0.4	0.1	0.5
z_2	0.1	0.5	0.4

■ 连续观测变量：高斯分布

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

隐马尔可夫模型

- 例子：
 - 从3个隐含状态采样的20个观测数据
 - 对应的隐含状态序列





隐马尔可夫模型的应用

- 在**时间序列**中建立观测与隐藏语义间的映射关系
- 语音识别
 - x : 从语音信号中提取的声学特征
 - z : 发出的音素或单词
- 行为识别
 - x : 从视频帧中提取的姿态或运动特征
 - z : 人的行为 (如跑、走、坐等)



模型推断

■ 概率计算问题

- 给定模型参数 θ 和观测序列 $x_{1:T}$ ，计算生成该观测序列的概率 $p(x_{1:T}|\theta)$
 - 前向算法: $p(z_t|x_{1:t})$ —— Filtering (在线学习)
 - 后向算法: $p(z_t|x_{1:T})$ —— Smoothing (离线学习)

■ 学习问题

- 给定观测序列 $x_{1:T}$ ，学习使生成该序列的概率最大的模型参数 θ
 $\max p(x_{1:T}|\theta)$

■ 预测问题 (解码问题)

- 给定模型参数 θ 和观测序列 $x_{1:T}$ ，求最有可能的隐含状态序列
 $\max p(z_{1:T}|x_{1:T}, \theta)$



概率计算问题

■ 前向算法

- 预测：一步前预测密度 (one-step-ahead predictive density)

$$p(z_t = j | x_{1:t-1}, \theta) = \sum_i p(z_t = j | z_{t-1} = i, \theta) p(z_{t-1} = i | x_{1:t-1}, \theta).$$

- 更新： t 时刻的信念状态 (belief state)

$$\begin{aligned} \alpha_t(j) &\triangleq p(z_t = j | x_{1:t}, \theta) = p(z_t = j | x_t, x_{1:t-1}, \theta) \\ &= \frac{1}{Z_t} p(x_t | z_t = j, x_{1:t-1}, \theta) p(z_t = j | x_{1:t-1}, \theta) \\ &= \frac{1}{Z_t} p(x_t | z_t = j, \theta) p(z_t = j | x_{1:t-1}, \theta), \end{aligned}$$

概率计算问题

■ 前向算法

□ 更新： t 时刻的信念状态

$$\begin{aligned} \alpha_t(j) &\propto p(x_t|z_t = j, \theta)p(z_t = j|x_{1:t-1}, \theta) \\ &= p(x_t|z_t = j, \theta) \sum_i p(z_t = j|z_{t-1} = i, \theta)p(z_{t-1} = i|x_{1:t-1}, \theta) \\ &= p(x_t|z_t = j, \theta) \sum_i p(z_t = j|z_{t-1} = i, \theta)\alpha_{t-1}(i) \end{aligned}$$

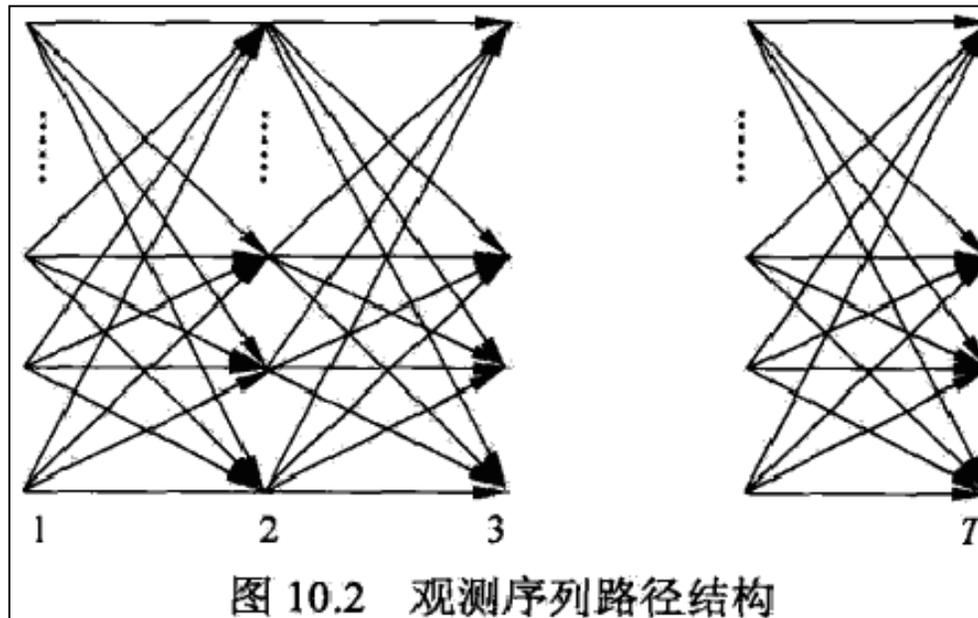
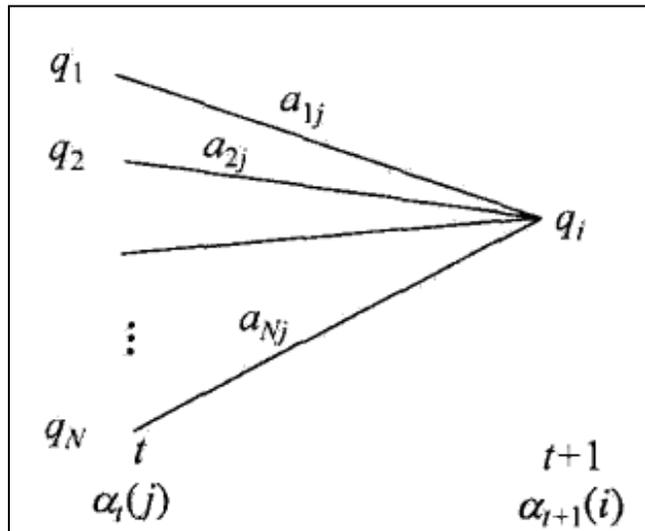
$$\Rightarrow \alpha_t \propto \psi_t \odot (\Psi^T \alpha_{t-1})$$

$$\psi_t(j) = p(x_t|z_t = j, \theta)$$

$$\Psi(i, j) = p(z_t = j|z_{t-1} = i)$$

概率计算问题

■ 前向算法的图形化解释





概率计算问题

- 前向算法：观测序列概率
 - 计算 $p(\mathbf{x}_{1:T}|\theta)$

$$\begin{aligned}\log p(\mathbf{x}_{1:T}|\theta) &= \log \left[p(x_1|\theta) \prod_{t=2}^T p(x_t|x_{1:t-1}, \theta) \right] \\ &= \log p(x_1|\theta) + \sum_{t=2}^T \log p(x_t|x_{1:t-1}, \theta) \\ &= \log \sum_{i=1}^K p(z_1 = i|\theta) p(x_1|z_1 = i, \theta) + \sum_{t=2}^T \log Z_t\end{aligned}$$

$$Z_t \triangleq p(x_t|x_{1:t-1}, \theta) = \sum_j p(z_t = j|x_{1:t-1}, \theta) p(x_t|z_t = j, \theta)$$

概率计算问题

■ 后向算法

□ 基本思想

$$\begin{aligned}
 p(z_t = j | x_{1:T}, \theta) &= p(z_t = j | x_{1:t}, x_{t+1:T}, \theta) \\
 &= \frac{p(z_t = j, x_{t+1:T} | x_{1:t}, \theta)}{p(x_{t+1:T} | x_{1:t}, \theta)} \\
 &= \frac{p(x_{t+1:T} | z_t = j, x_{1:t}, \theta) p(z_t = j | x_{1:t}, \theta)}{p(x_{t+1:T} | x_{1:t}, \theta)} \\
 &\propto p(z_t = j | x_{1:t}, \theta) p(x_{t+1:T} | z_t = j, x_{1:t}, \theta) \\
 &\propto \boxed{p(z_t = j | x_{1:t}, \theta)} \boxed{p(x_{t+1:T} | z_t = j, \theta)}.
 \end{aligned}$$

过去：信念状态
(已知)

未来：后向变量
(后向算法的计算目标)



概率计算问题

■ 后向算法

□ 定义

■ 后向变量

$$\beta_t(j) = p(x_{t+1:T} | z_t = j, \theta)$$

■ 平滑后的后验分布

$$\gamma_t(j) = p(z_t = j | x_{1:T}, \theta)$$

■ 直观上

$$\gamma_t(j) \propto \alpha_t(j)\beta_t(j)$$



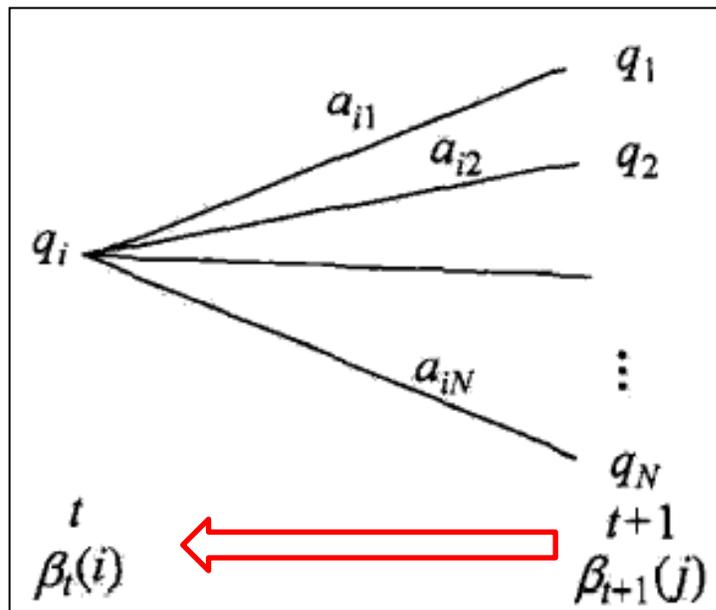
概率计算问题

- 后向算法
 - 计算后向变量

$$\begin{aligned}\beta_{t-1}(i) &= p(x_{t:T} | z_{t-1} = i, \theta) \\ &= \sum_j p(z_t = j, x_t, x_{t+1:T} | z_{t-1} = i, \theta) \\ &= \sum_j p(x_{t+1:T} | z_t = j, \cancel{x_t}, \cancel{z_{t-1}} = i, \theta) p(z_t = j, x_t | z_{t-1} = i, \theta) \\ &= \sum_j p(x_{t+1:T} | z_t = j, \theta) p(x_t | z_t = j, \cancel{z_{t-1}} = i, \theta) p(z_t = j | z_{t-1} = i, \theta) \\ &= \sum_j \beta_t(j) \psi_t(j) \Psi(i, j) \quad \longrightarrow \quad \boxed{\beta_{t-1} = \Psi(\psi_t \odot \beta_t)}\end{aligned}$$

概率计算问题

- 后向算法的图形化解释
 - 从 $t + 1$ 逆推



- 最后时刻 T , 初始化

$$\beta_T(i) = p(x_{T+1} | z_T = i, \theta) = p(\emptyset | z_T = i, \theta) = 1$$



概率计算问题

- 后向算法：观测序列概率
 - 计算 $p(\mathbf{x}_{1:T}|\theta)$

$$\begin{aligned} p(\mathbf{x}_{1:T}|\theta) &= \sum_j p(z_1 = j|\theta)p(\mathbf{x}_{1:T}|z_1 = j, \theta) \\ &= \sum_j p(z_1 = j|\theta)p(\mathbf{x}_1|z_1 = j, \theta)p(\mathbf{x}_{2:T}|\mathbf{x}_1, z_1 = j, \theta) \\ &= \sum_j p(z_1 = j|\theta)p(\mathbf{x}_1|z_1 = j, \theta)\beta_1(j) \end{aligned}$$

初始状态 × 当前观测 × 未来证据



学习问题

- 估计模型参数 $\theta = (\pi, A, B)$

- 隐含状态的初始概率分布 π

- 状态转移概率矩阵 $A = \Psi$

$$A(i, j) = p(z_t = j | z_{t-1} = i)$$

- 观测概率矩阵 B

- 离散观测：概率表

- 连续观测：类条件概率密度函数 $p(x_t | z_t = j, \theta)$

- 目标： $\arg \max_{\theta} p(x_{1:T} | \theta)$



学习问题

- 如果 z 可观测
 - 监督学习方法
 - 极大似然估计
- 如果 z 不可观测
 - 非监督学习方法
 - Baum-Welch 算法



极大似然估计

■ 估计 π, A

□ 给定观测数据集

$$\mathcal{D} = \{(\mathbf{z}_i, \mathbf{x}_i), i = 1, 2, \dots, N | \mathbf{z}_i = z_{i,1:T_i}, \mathbf{x}_i = x_{i,1:T_i}\}$$

□ 生成1个长度为 T 的状态序列数据 $z_{1:T}$ 的似然函数

$$\begin{aligned} p(z_{1:T}|\theta) &= p(z_1)p(z_2|z_1)\cdots p(z_T|z_{T-1}) \\ &= \pi(z_1)\Psi(z_1, z_2)\cdots\Psi(z_{T-1}, z_T) \\ &= \prod_{j=1}^K (\pi_j)^{\mathbb{I}[z_1=j]} \prod_{t=2}^T \prod_{j=1}^K \prod_{k=1}^K (\Psi_{j,k})^{\mathbb{I}[z_{t-1}=j, z_t=k]} \end{aligned}$$



极大似然估计

- 估计 π, A
 - 生成 N 个状态序列数据的对数似然函数
 - 记 N 个状态序列数据的集合为 \mathcal{D}_z

$$\mathcal{D}_z = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$$

- 模型参数 θ 生成状态序列集 \mathcal{D}_z 的对数似然函数为**所有样本似然的和**

$$\begin{aligned} \log p(\mathcal{D}_z | \theta) &= \sum_{i=1}^N \log p(\mathbf{z}_i | \theta) \\ &= \sum_{i=1}^N \log \left[\prod_{j=1}^K (\pi_j)^{\mathbb{I}[z_{i1}=j]} \prod_{t=2}^{T_i} \prod_{j=1}^K \prod_{k=1}^K (\Psi_{j,k})^{\mathbb{I}[z_{i,t-1}=j, z_{it}=k]} \right] \\ &= \sum_{i=1}^N \sum_{j=1}^K \mathbb{I}[z_{i1}=j] \log \pi_j + \sum_{i=1}^N \sum_{t=2}^{T_i} \sum_{j=1}^K \sum_{k=1}^K \mathbb{I}[z_{i,t-1}=j, z_{it}=k] \log \Psi_{j,k} \end{aligned}$$



极大似然估计

■ 估计 π, A

- 定义初始状态计数 N_j^1 、状态转移计数 N_j^k

$$N_j^1 = \sum_{i=1}^N \mathbb{I}[z_{i1} = j], \quad N_{jk} = \sum_{i=1}^N \sum_{t=2}^{T_i} \mathbb{I}[z_{i,t-1} = j, z_{it} = k]$$

- 参数估计

$$\hat{\pi}_j = \frac{N_j^1}{\sum_j N_j^1}, \quad \hat{A}_{jk} = \frac{N_{jk}}{\sum_k N_{jk}}$$



极大似然估计

■ 估计 B

- 如果观测变量 x 是**离散**的
 - 假设观测变量 x 服从**多项分布**

$$B(k, l) = p(x_t = l | z_t = k, \theta), k = 1, 2, \dots, K, l = 1, 2, \dots, L$$

■ 参数估计

$$\hat{B}_{kl} = \frac{N_{kl}^X}{N_k}$$

$$N_{kl}^X = \sum_{i=1}^N \sum_{t=1}^{T_i} \mathbb{I}[z_{i,t} = k, x_{i,t} = l]$$

$$N_k = \sum_{i=1}^N \sum_{t=1}^{T_i} \mathbb{I}[z_{i,t} = k]$$

- N_{kl}^X 表示在整个数据集D中，状态 k 生成观测变量的出现次数
- N_k 表示在整个数据集D中，状态 k 出现的次数

极大似然估计

■ 估计 B

□ 如果观测变量 x 是**连续的**

■ 假设观测变量 x 服从**高斯分布**

$$p(x_t | z_t = k, \theta) = \mathcal{N}(x_t | \mu_k, \Sigma_k)$$

■ 参数估计

$$\hat{\mu}_k = \frac{\bar{\mathbf{x}}_k}{N_k}, \quad \hat{\Sigma}_k = \frac{(\overline{\mathbf{x}\mathbf{x}})_k^T - N_k \hat{\mu}_k \hat{\mu}_k^T}{N_k}$$

$$\bar{\mathbf{x}}_k \triangleq \sum_{i=1}^N \sum_{t=1}^{T_i} \mathbb{I}(z_{i,t} = k) \mathbf{x}_{i,t}$$

$$(\overline{\mathbf{x}\mathbf{x}})_k^T \triangleq \sum_{i=1}^N \sum_{t=1}^{T_i} \mathbb{I}(z_{i,t} = k) \mathbf{x}_{i,t} \mathbf{x}_{i,t}^T$$

□ $\hat{\mu}_k$ 表示在整个数据集 D 中，由状态 k 生成的观测变量 x_t 的平均值

□ $\hat{\Sigma}_k$ 表示在整个数据集 D 中，由状态 k 生成的观测变量 x_t 的经验协方差矩阵



Baum-Welch算法

- 多数实际情况：
 - 隐状态 z 不可观测，MLE难以直接求解

- EM算法

- E步

$$\begin{aligned} Q(\theta, \theta^{old}) &= \sum_{z_{1:T}} p(z_{1:T} | x_{1:T}, \theta^{old}) \log p(z_{1:T}, x_{1:T} | \theta) \\ &= \sum_{k=1}^K \mathbb{E} [N_k^1] \log \pi_k + \sum_{j=1}^K \sum_{k=1}^K \mathbb{E} [N_{jk}] \log A_{jk} \\ &\quad + \sum_{i=1}^N \sum_{t=1}^{T_i} \sum_{k=1}^K p(z_t = k | \mathbf{x}_i, \theta^{old}) \log p(\mathbf{x}_{i,t} | \phi_k) \end{aligned}$$

初始分布 状态转移 观测概率

Baum-Welch算法



■ E步

- 期望计数可以通过前向-后向算法求解

$$\begin{aligned}\mathbb{E}[N_k^1] &= \sum_{i=1}^N p(z_{i1} = k | \mathbf{x}_i, \boldsymbol{\theta}^{old}) \\ \mathbb{E}[N_{jk}] &= \sum_{i=1}^N \sum_{t=2}^{T_i} p(z_{i,t-1} = j, z_{i,t} = k | \mathbf{x}_i, \boldsymbol{\theta}^{old}) \\ \mathbb{E}[N_j] &= \sum_{i=1}^N \sum_{t=1}^{T_i} p(z_{i,t} = j | \mathbf{x}_i, \boldsymbol{\theta}^{old})\end{aligned}$$

$$\begin{aligned}\gamma_{i,t}(j) &\triangleq p(z_t = j | \mathbf{x}_{i,1:T_i}, \boldsymbol{\theta}) \\ \xi_{i,t}(j, k) &\triangleq p(z_{t-1} = j, z_t = k | \mathbf{x}_{i,1:T_i}, \boldsymbol{\theta})\end{aligned}$$



Baum-Welch算法

■ M步

- 对Q函数分别关于 π, A 求偏导

$$\hat{A}_{jk} = \frac{\mathbb{E}[N_{jk}]}{\sum_{k'} \mathbb{E}[N_{jk'}]}, \quad \hat{\pi}_k = \frac{\mathbb{E}[N_k^1]}{N}$$

- 假设观测变量 x 离散，服从多项分布

$$\hat{B}_{jl} = \frac{\mathbb{E}[M_{jl}]}{\mathbb{E}[N_j]}$$

$$\mathbb{E}[M_{jl}] = \sum_{i=1}^N \sum_{t=1}^{T_i} \gamma_{i,t}(j) \mathbb{I}(x_{i,t} = l) = \sum_{i=1}^N \sum_{t: x_{i,t}=l} \gamma_{i,t}(j)$$



Baum-Welch算法

■ M步

- 假设观测变量 x 连续，服从高斯分布

$$\hat{\boldsymbol{\mu}}_k = \frac{\mathbb{E}[\bar{\mathbf{x}}_k]}{\mathbb{E}[N_k]}, \quad \hat{\boldsymbol{\Sigma}}_k = \frac{\mathbb{E}[(\bar{\mathbf{x}\mathbf{x}})_k^T] - \mathbb{E}[N_k] \hat{\boldsymbol{\mu}}_k \hat{\boldsymbol{\mu}}_k^T}{\mathbb{E}[N_k]}$$

$$\mathbb{E}[\bar{\mathbf{x}}_k] = \sum_{i=1}^N \sum_{t=1}^{T_i} \gamma_{i,t}(k) \mathbf{x}_{i,t}$$
$$\mathbb{E}[(\bar{\mathbf{x}\mathbf{x}})_k^T] = \sum_{i=1}^N \sum_{t=1}^{T_i} \gamma_{i,t}(k) \mathbf{x}_{i,t} \mathbf{x}_{i,t}^T$$

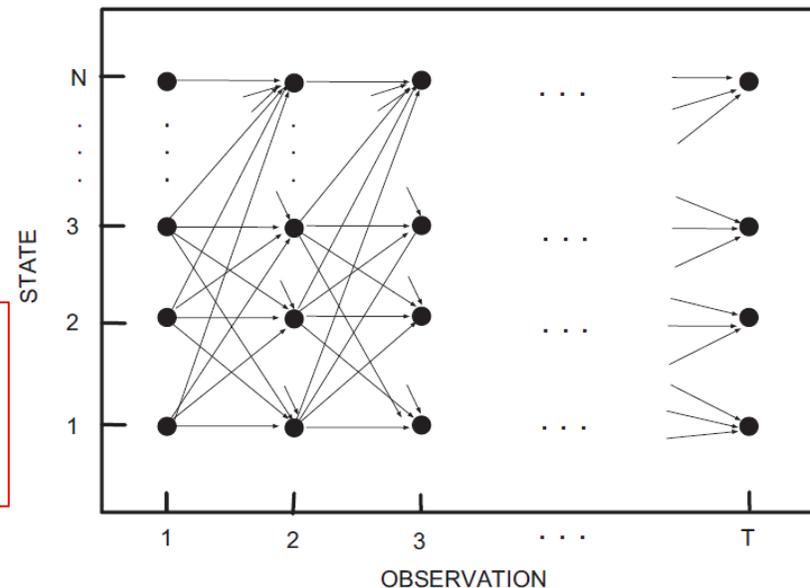
预测问题

- 给定模型参数和观测数据，找到最有可能的状态序列

$$\mathbf{z}^* = \arg \max_{\mathbf{z}_{1:T}} p(\mathbf{z}_{1:T} | \mathbf{X}_{1:T})$$

- 等价于最大路径问题

$$\log \pi_1(z_1) + \log \phi_1(z_1) + \sum_{t=2}^T [\log \psi(z_{t-1}, z_t) + \log \phi_t(z_t)]$$



■ Viterbi算法

- 用**动态规划**解概率最大路径，一个路径对应一个状态序列
- 最优路径：如果最优路径在时刻 t 通过结点 z_t^* ，那么这一路径从结点 z_t^* 到终点 z_T^* 的部分路径，对于从 z_t^* 到 z_T^* 的所有可能的部分路径而言，必须是最优的
- 为了找出最优路径的各个结点，从终结点开始，由后向前逐步求得结点 $z_T^*, z_{T-1}^*, \dots, z_1^*$ ，得到最优路径

预测问题

■ 定义

- $\delta_t(j)$: 在时刻 t 隐藏状态为 j 的所有可能的状态转移路径 z_1, \dots, z_t 中的概率最大值

$$\delta_t(j) \triangleq \max_{z_1, \dots, z_{t-1}} p(\mathbf{z}_{1:t-1}, z_t = j | \mathbf{x}_{1:t})$$



$$\delta_t(j) = \max_i \delta_{t-1}(i) \psi(i, j) \phi_t(j)$$

- 在保证时刻 t 隐藏状态为 j 的前提下， $t-1$ 时刻的最优状态

$$a_t(j) = \operatorname{argmax}_i \delta_{t-1}(i) \psi(i, j) \phi_t(j)$$



预测问题

- 初始化

$$\delta_1(j) = \pi_j \phi_1(j)$$

- 计算t时刻最优隐含状态

$$z_T^* = \arg \max_i \delta_T(i)$$

$$z_t^* = a_{t+1}(z_{t+1}^*)$$

Thanks!



Questions?