



## **Chapter 3:** 最大似然估计和贝叶斯参数估计



## 要点:

- 重点掌握最大似然估计和贝叶斯参数估计的原理;
- 熟练掌握主成分分析和Fisher线性分析;
- 掌握隐马尔可夫模型;
- 了解维数问题;



## • 3.1 引言

### □ 贝叶斯框架下的数据收集

- 在以下条件下我们可以设计一个可选择的分类器：

- $P(\omega_i)$  (先验)

- $P(x | \omega_i)$  (类条件密度)

不幸的是，我们极少能够完整的得到这些信息！

### □ 从一个传统的样本中设计一个分类器

- 先验估计不成问题

- 对类条件密度的估计存在两个问题：1) 样本对于类条件估计太少了；2) 特征空间维数太大了，计算复杂度太高。



- 如果可以将类条件密度参数化，则可以显著降低难度。

- 例如： $P(x | \omega_i)$ 的正态性

$$P(x | \omega_i) \sim N(\mu_i, \Sigma_i)$$

- 用两个参数表示

将概率密度估计问题转化为参数估计问题。

- 估计

- 最大似然估计 (ML) 和贝叶斯估计；
  - 结果通常很接近, 但是方法本质是不同的。



- 最大似然估计将参数看作是确定的量，只是其值是未知! 通过最大化所观察的样本概率得到最优的参数——用分析方法。
- 贝叶斯方法把参数当成服从某种先验概率分布的随机变量，对样本进行观测的过程，就是把先验概率密度转化成为后验概率密度，使得对于每个新样本，后验概率密度函数在待估参数的真实值附近形成最大尖峰。
- 在这两种方法中，我们都用后验概率 $P(\omega_i | \mathbf{x})$ 表示分类准则!



## • 3.2 最大似然估计

### □ 最大似然估计的优点：

- 当样本数目增加时，收敛性质会更好；
- 比其他可选择的技术更加简单。

### 3.2.1 基本原理

假设有 $c$ 类样本，并且

- 1) 每个样本集的样本都是独立同分布的随机变量；
- 2)  $P(x | \omega_j)$  形式已知但参数未知，例如  $P(x | \omega_j) \sim N(\mu_j, \Sigma_j)$ ；
- 3) 记  $P(x | \omega_j) \equiv P(x | \omega_j, \theta_j)$ ，其中

$$\theta_j = (\mu_j, \Sigma_j)$$



- 使用训练样本提供的信息估计

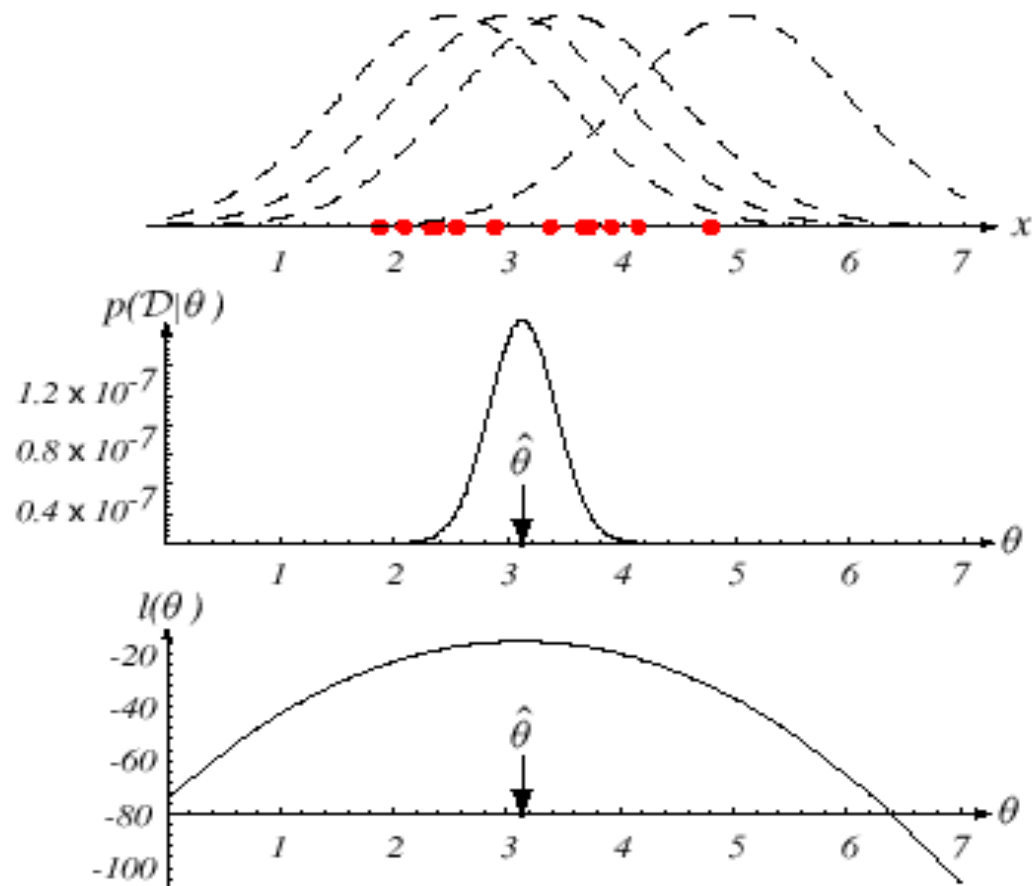
$\theta = (\theta_1, \theta_2, \dots, \theta_c)$ , 每个  $\theta_i$  ( $i = 1, 2, \dots, c$ ) 只和每一类相关。

- 假定D包括n个样本,  $x_1, x_2, \dots, x_n$

$$P(D | \theta) = \prod_{k=1}^{k=n} P(x_k | \theta) = F(\theta)$$

$P(D | \theta)$  被称为样本集D下的似然函数

- $\theta$ 的最大似然估计是通过定义最大化 $P(D | \theta)$ 的值  $\hat{\theta}$   
“ $\theta$ 值与实际观察中的训练样本最相符”



**FIGURE 3.1.** The top graph shows several training points in one dimension, known or assumed to be drawn from a Gaussian of a particular variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines. The middle figure shows the likelihood  $p(\mathcal{D}|\theta)$  as a function of the mean. If we had a very large number of training points, this likelihood would be very narrow. The value that maximizes the likelihood is marked  $\hat{\theta}$ ; it also maximizes the logarithm of the likelihood—that is, the log-likelihood  $l(\theta)$ , shown at the bottom. Note that even though they look similar, the likelihood  $p(\mathcal{D}|\theta)$  is shown as a function of  $\theta$  whereas the conditional density  $p(x|\theta)$  is shown as a function of  $x$ . Furthermore, as a function of  $\theta$ , the likelihood  $p(\mathcal{D}|\theta)$  is not a probability density function and its area has no significance. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.





## ■ 最优估计

- 令  $\theta = (\theta_1, \theta_2, \dots, \theta_p)^t$  并令  $\nabla_{\theta}$  为梯度算子 the gradient operator

$$\nabla_{\theta} = \left[ \frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_2}, \dots, \frac{\partial}{\partial \theta_p} \right]^t$$

- 我们定义  $l(\theta)$  为对数似然函数:  $l(\theta) = \ln P(D | \theta)$
- 新问题陈述:  
求解  $\theta$  为使对数似然最大的值

$$\hat{\theta} = \arg \max_{\theta} l(\theta)$$



对数似然函数 $l(\theta)$ 显然是依赖于样本集 $\mathbf{D}$ , 有:

$$l(\theta) = \sum_{k=1}^n \ln P(x_k | \theta)$$

最优求解条件如下:

$$\nabla_{\theta} l(\theta) = \sum_{k=1}^n \nabla_{\theta} \ln P(x_k | \theta)$$

令:

$$\nabla_{\theta} l(\theta) = 0$$

来求解.



### 3.2.3 高斯情况： $\mu$ 未知

- $P(x_k | \mu) \sim N(\mu, \Sigma)$   
(样本从一组多变量正态分布中提取)

$$\ln P(x_k | \mu) = -\frac{1}{2} \ln[(2\pi)^d |\Sigma|] - \frac{1}{2} (x_k - \mu)^t \Sigma^{-1} (x_k - \mu)$$

和  $\nabla_{\mu} \ln P(x_k | \mu) = \Sigma^{-1} (x_k - \mu)$

这里  $\theta = \mu$ ，因此：

- $\mu$ 的最大似然估计必须满足：

$$\sum_{k=1}^n \Sigma^{-1} (x_k - \hat{\mu}) = 0$$



- 乘  $\Sigma$  并且重新排序, 我们得到:

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k$$

即训练样本的算术平均值!

结论:

如果  $P(x_k | \omega_j)$  ( $j = 1, 2, \dots, c$ ) 被假定为  $d$  维特征空间中的高斯分布; 然后我们能够估计向量  $\theta = (\theta_1, \theta_2, \dots, \theta_c)^t$  从而得到最优分类!



### 3.2.3 高斯情况： $\mu$ 和 $\Sigma$ 均未知

- 未知  $\mu$  和  $\sigma$ ，对于单样本  $\mathbf{x}_k$

$$\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$$

$$l(\theta) = \ln P(x_k | \theta) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2} (x_k - \theta_1)^2$$

$$\nabla_{\theta} l = \begin{pmatrix} \frac{\sigma}{\sigma\theta_1} (\ln P(x_k | \theta)) \\ \frac{\sigma}{\sigma\theta_2} (\ln P(x_k | \theta)) \end{pmatrix} = 0$$

$$\begin{cases} \frac{1}{\theta_2} (x_k - \theta_1) = 0 \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} = 0 \end{cases}$$



对于全部样本，最后得到:

$$\left\{ \begin{array}{l} \sum_{k=1}^n \frac{1}{\hat{\theta}_2} (x_k - \theta_1) = 0 \end{array} \right. \quad (1)$$

$$\left\{ \begin{array}{l} -\sum_{k=1}^n \frac{1}{\hat{\theta}_2} + \sum_{k=1}^n \frac{(x_k - \hat{\theta}_1)^2}{\hat{\theta}_2^2} = 0 \end{array} \right. \quad (2)$$

联合公式 (1) 和 (2), 得到如下结果:

$$\mu = \sum_{k=1}^n \frac{x_k}{n} \quad ; \quad \sigma^2 = \frac{\sum_{k=1}^n (x_k - \mu)^2}{n}$$



### 3.2.4 偏差估计

- $\sigma^2$  的最大似然估计是有偏的（渐进无偏估计）

$$E\left[\frac{1}{n}\sum (x_i - \bar{x})^2\right] = \frac{n-1}{n} \cdot \sigma^2 \neq \sigma^2$$

- $\Sigma$  的一个基本的无偏估计是:

$$\underbrace{C = \frac{1}{n-1} \sum_{k=1}^{k=n} (x_k - \mu)(x_k - \hat{\mu})^t}_{\text{Sample covariance matrix}}$$



# 模型错误会怎么样？

达不到最优！





## • 3.3 贝叶斯估计

- 在最大似然估计中  $\theta$  被假定为固定值
- 在贝叶斯估计中  $\theta$  是随机变量

### 3.3.1 类条件密度

- 目标: 计算  $P(\omega_i | x, D)$

假设样本为  $D$ , 贝叶斯方程可以写成 :

$$P(\omega_i | x, D) = \frac{P(x | \omega_i, D).P(\omega_i | D)}{\sum_{j=1}^c P(x | \omega_j, D).P(\omega_j | D)}$$



- 先验概率通常可以事先获得，因此

$$P(\omega_i) = P(\omega_i | D)$$

- 每个样本只依赖于所属的类，有：

$$P(x | \omega_i, D) = P(x | \omega_i, D_i)$$

故：

$$P(\omega_i | x, D) = \frac{P(x | \omega_i, D_i).P(\omega_i)}{\sum_{j=1}^c P(x | \omega_j, D_j).P(\omega_j)}$$

即：只要在每类中，独立计算  $P(x | \omega_i, D_i)$  就可以确定x的类别。

因此，核心工作就是要估计  $P(x | D)$



### 3.3.2 参数的分布

- 假设  $p(x)$  的形式已知, 参数 $\theta$ 的值未知, 因此条件概率密度  $p(x|\theta)$  是知道的;
- 假设参数 $\theta$ 是随机变量, 先验概率密度函数 $p(\theta)$ 已知, 利用贝叶斯公式可以计算后验概率密度函数  $p(\theta | D)$  ;
- 希望后验概率密度函数 $p(\theta | D)$  在 $\theta$ 的真实值附件有非常显著的尖峰, 则可以使用后验密度 $p(\theta | D)$  估计  $\theta$  ;



### 3.3.2 参数的分布

➤ 注意到

$$\begin{aligned} p(x | D) &= \int p(x, \theta | D) d\theta \\ &= \int p(x | \theta) p(\theta | D) d\theta \end{aligned}$$

如果  $p(\theta | D)$  在某个值  $\hat{\theta}$  附近有非常显著的尖峰，  
则  $p(x | D) \approx p(x | \hat{\theta})$

即：如果条件概率密度具有一个已知的形式，  
则利用已有的训练样本，就能够通过  $p(\theta | D)$   
对  $p(x | D)$  进行估计。



## 3.4 贝叶斯参数估计: 高斯过程

### □ 单变量情形的 $p(\mu | D)$

$p(x | \mu) \sim N(\mu, \sigma^2)$ ,  $\mu$  是未知的。

假设  $p(\mu) \sim N(\mu_0, \sigma_0^2)$ ,  $\mu_0$  和  $\sigma_0^2$  已知

( $\mu_0$  是  $\mu$  最好的估计;  $\sigma_0^2$  是该估计的不确定性)

$$D = \{x_1, \dots, x_n\}, \quad p(\mu | D) = \frac{p(D | \mu) p(\mu)}{\int p(D | \mu) p(\mu) d\mu}$$

$$p(\mu | D) = \alpha \prod_{k=1}^n p(x_k | \mu) p(\mu)$$

$$= \alpha' \exp \left[ -\frac{1}{2} \left( \sum_{k=1}^n \left( \frac{\mu - x_k}{\sigma} \right)^2 + \left( \frac{\mu - \mu_0}{\sigma_0} \right)^2 \right) \right]$$

$$= \alpha'' \exp \left[ -\frac{1}{2} \left[ \left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - 2 \left( \frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2} \right) \mu \right] \right]$$



# 复制密度

$p(\mu | D) \sim N(\mu_n, \sigma_n^2)$  [reproducing density]

[称  $p(\mu)$ : conjugate prior]

$$\frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}, \quad \frac{\mu_n}{\sigma_n^2} = \frac{n}{\sigma^2} \hat{\mu}_n + \frac{\mu_0}{\sigma_0^2}$$

$$\mu_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$$

$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}$$

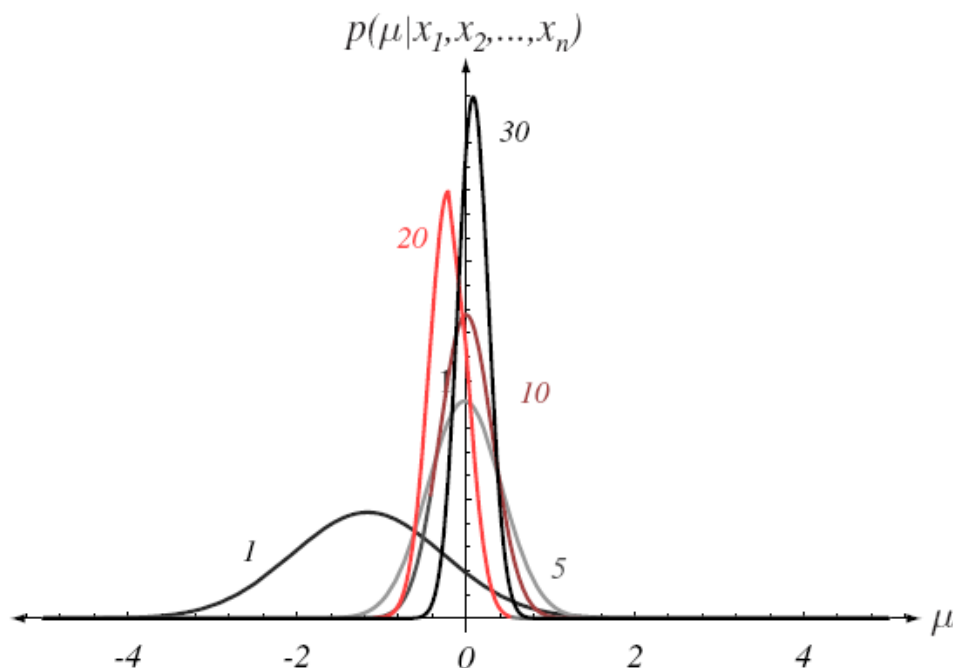
其中,  $\hat{\mu}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$



结论:

$\mu_n$  是  $\hat{\mu}_n$  和  $\mu_0$  的线性组合, 总是位于  $\hat{\mu}_n$  和  $\mu_0$  的连线上;  
当  $\sigma_0^2 \neq 0$  时,  $\mu_n$  将逼近  $\hat{\mu}_n$ , 否则  $\mu_n = \mu_0$ 。

贝叶斯学习





## □ 单变量情形的 $p(x|D)$

$$p(x|D) = \int p(x|\mu)p(\mu|D)d\mu$$
$$= \frac{1}{2\pi\sigma\sigma_n} \exp\left[-\frac{1}{2} \frac{(x-\mu_n)^2}{\sigma^2 + \sigma_n^2}\right] f(\sigma, \sigma_n)$$

$$\text{其中, } f(\sigma, \sigma_n) = \int_{-\infty}^{\infty} \exp\left[-\frac{1}{2} \frac{\sigma^2 + \sigma_n^2}{\sigma^2 \sigma_n^2} \left(\mu - \frac{\sigma_n^2 x + \sigma^2 n}{\sigma^2 + \sigma_n^2}\right)^2\right] d\mu$$

$$= \sqrt{2\pi \left( \frac{\sigma^2 \sigma_n^2}{\sigma^2 + \sigma_n^2} \right)}$$

$$p(x|D) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$$





多变量情形:

$$p(\mathbf{x} | \boldsymbol{\mu}) \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}), p(\boldsymbol{\mu}) \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad \text{其中仅}\boldsymbol{\mu}\text{未知.}$$

$$D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

$$p(\boldsymbol{\mu} | D) = \alpha \prod_{k=1}^n p(x_k | \boldsymbol{\mu}) p(\boldsymbol{\mu})$$

$$= \alpha' \exp \left[ -\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_n)^t \boldsymbol{\Sigma}_n^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_n) \right] \quad \text{复制密度}$$

$$\boldsymbol{\Sigma}_n^{-1} = n\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}_0^{-1}, \quad \boldsymbol{\Sigma}_n^{-1} \boldsymbol{\mu}_n = n\boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_n + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0$$

$$\text{其中, } \hat{\boldsymbol{\mu}}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$



利用  $(\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} = \mathbf{A}(\mathbf{A} + \mathbf{B})^{-1}\mathbf{B} = \mathbf{B}(\mathbf{A} + \mathbf{B})^{-1}\mathbf{A}$ , 得

$$\boldsymbol{\mu}_n = \boldsymbol{\Sigma}_0 \left( \boldsymbol{\Sigma}_0 + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \hat{\boldsymbol{\mu}}_n + \frac{1}{n} \boldsymbol{\Sigma} \left( \boldsymbol{\Sigma}_0 + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \boldsymbol{\mu}_0$$

$$\boldsymbol{\Sigma}_n = \boldsymbol{\Sigma}_0 \left( \boldsymbol{\Sigma}_0 + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \frac{1}{n} \boldsymbol{\Sigma}$$

利用  $p(\mathbf{x} | D) = \int p(\mathbf{x} | \boldsymbol{\mu}) p(\boldsymbol{\mu} | D) d\boldsymbol{\mu}$

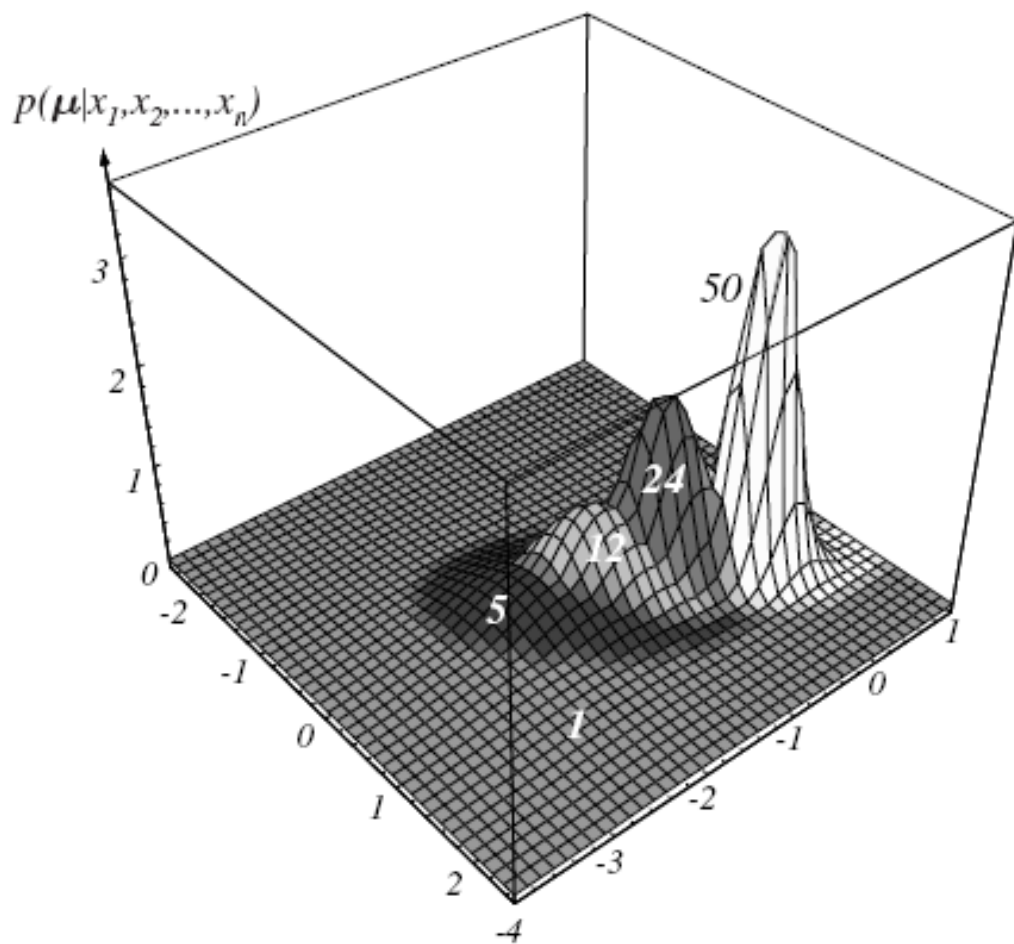
令  $\mathbf{y} = \mathbf{x} - \boldsymbol{\mu}$   $p(\mathbf{y} | \boldsymbol{\mu}) \sim N(\mathbf{0}, \boldsymbol{\Sigma})$

$$p(\boldsymbol{\mu} | D) \sim N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$$

$$\therefore p(\mathbf{x} | D) = p(\mathbf{y} + \boldsymbol{\mu} | D) \sim N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma} + \boldsymbol{\Sigma}_n)$$



# 多变量学习





## 3.5 贝叶斯参数估计：一般理论

- $p(x | D)$  的计算可推广于所有能参数化未知密度的情况中，基本假设如下：
  - 假定  $p(x | \theta)$  的形式已知，但是 $\theta$ 的值未知。
  - $\theta$ 被假定为满足一个已知的先验密度  $P(\theta)$
  - 其余的  $\theta$ 的信息 包含在集合 $D$ 中，其中 $D$ 是由 $n$ 维随机变量 $x_1, x_2, \dots, x_n$ 组成的集合，它们服从于概率密度函数 $p(x)$ 。

基本的问题是：

计算后验密度 $p(\theta | D)$ ，然后推导出  $p(x | D)$ 。



$$p(\mathbf{x} | D) = \int p(\mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | D) d\boldsymbol{\theta} \quad (49)$$

$$p(\boldsymbol{\theta} | D) = \frac{p(D | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int p(D | \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}} \quad (50)$$

$$p(D | \boldsymbol{\theta}) = \prod_{k=1}^n p(\mathbf{x}_k | \boldsymbol{\theta}) \quad (51)$$

问题:

$p(\mathbf{x} | D)$ 是否能收敛到 $p(\mathbf{x})$ , 计算复杂度如何?



# 递归贝叶斯学习

$$D^n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}, \quad p(D^n | \boldsymbol{\theta}) = p(\mathbf{x}_n | \boldsymbol{\theta}) p(D^{n-1} | \boldsymbol{\theta})$$

$$p(\boldsymbol{\theta} | D^n) = \frac{p(D^n | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int p(D^n | \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}}$$

$$= \frac{p(\mathbf{x}_n | \boldsymbol{\theta}) p(D^{n-1} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int p(\mathbf{x}_n | \boldsymbol{\theta}) p(D^{n-1} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}}$$

$$\text{因为: } p(\boldsymbol{\theta} | D^{n-1}) = \frac{p(D^{n-1} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int p(D^{n-1} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}}$$

$$\text{所以: } p(\boldsymbol{\theta} | D^n) = \frac{p(\mathbf{x}_n | \boldsymbol{\theta}) p(\boldsymbol{\theta} | D^{n-1})}{\int p(\mathbf{x}_n | \boldsymbol{\theta}) p(\boldsymbol{\theta} | D^{n-1}) d\boldsymbol{\theta}}$$

$$\text{令: } p(\boldsymbol{\theta} | D^0) = p(\boldsymbol{\theta})$$

该过程称为参数估计的递归贝叶斯方法，一种增量学习方法。



## 例1：递归贝叶斯学习

$$\text{假设: } p(x | \theta) \sim U(0, \theta) = \begin{cases} 1/\theta & 0 \leq x \leq \theta \\ 0 & \text{其他} \end{cases}$$

$$p(\theta) \sim U(0, 10), \quad D = \{4, 7, 2, 8\}$$

$$p(D^0 | \theta) = p(\theta) \sim U(0, 10)$$

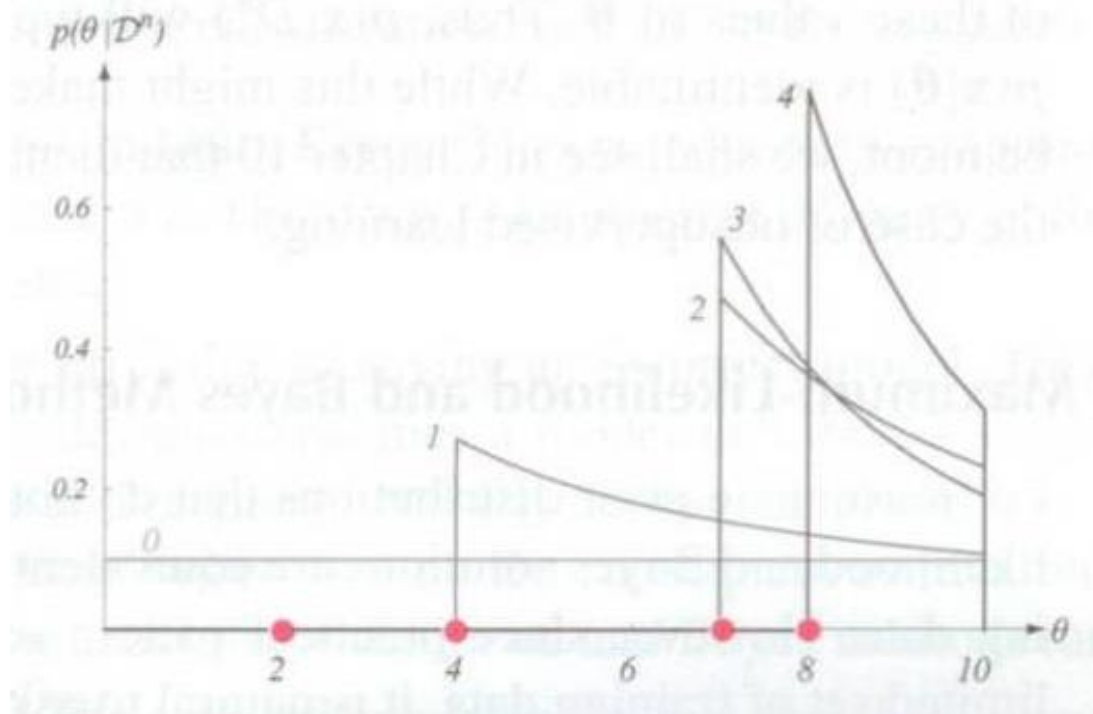
$$p(\theta | D^1) \propto p(x_1 | \theta) p(\theta | D^0) = \begin{cases} 1/\theta & \text{对于 } 4 \leq \theta \leq 10 \\ 0 & \text{其他} \end{cases}$$

$$p(\theta | D^2) \propto p(x_2 | \theta) p(\theta | D^1) = \begin{cases} 1/\theta^2 & \text{对于 } 7 \leq \theta \leq 10 \\ 0 & \text{其他} \end{cases}$$

$$p(\theta | D^n) \propto 1/\theta^n \quad \text{对于 } \max_x [D^n] \leq \theta \leq 10$$



例 1 贝叶斯推断

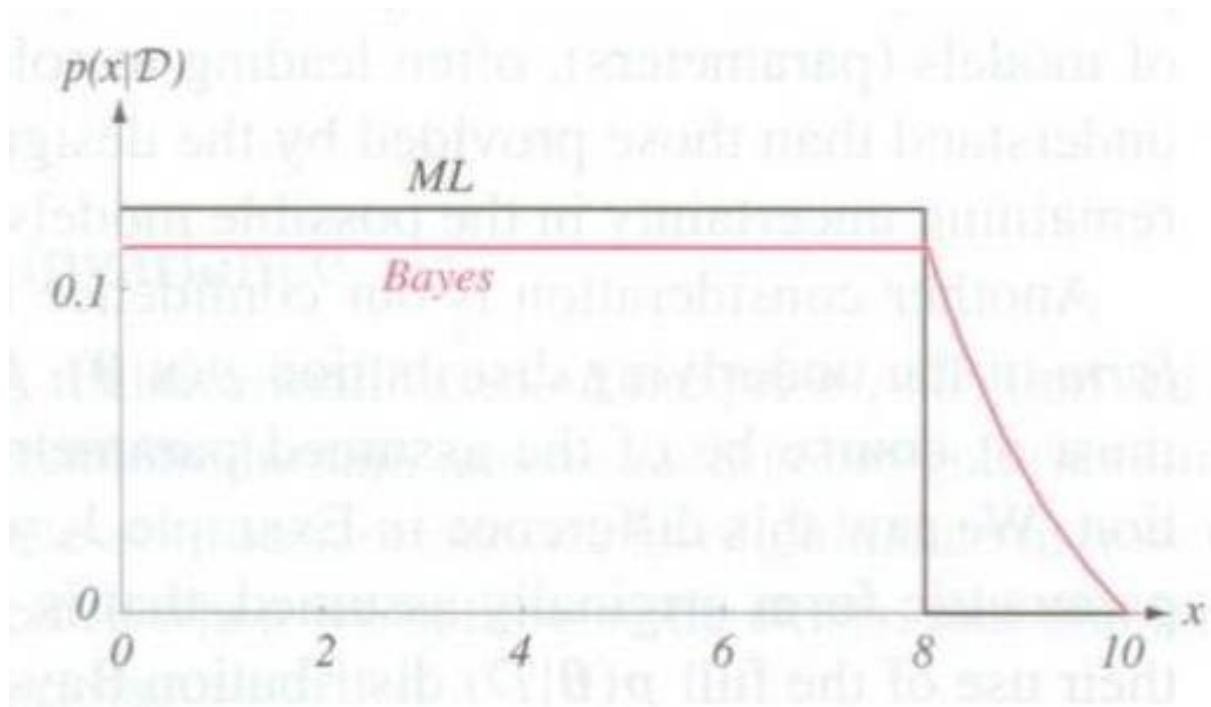


$$p(x | D) \sim U(0, 8)$$





# 例1: Bayes vs. ML

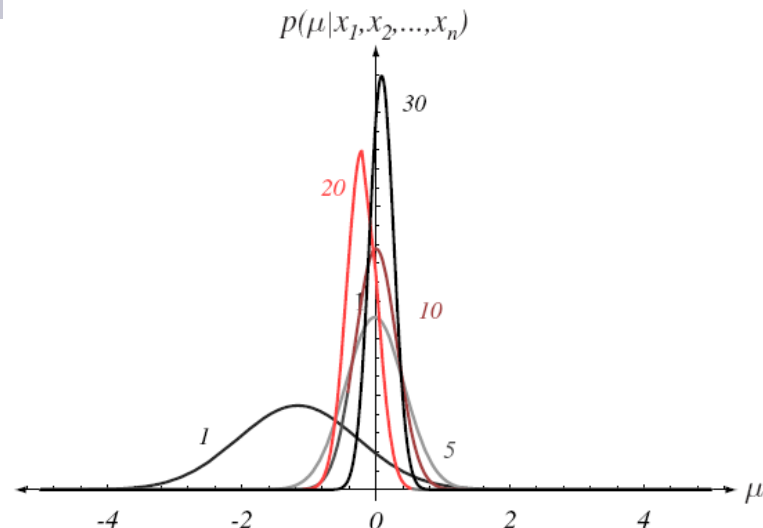


贝叶斯参数估计以来:  $p(\mathbf{x} | D) = \int p(\mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | D) d\boldsymbol{\theta}$

# 唯一性问题

■  $p(\mathbf{x}|\boldsymbol{\theta})$  是唯一的:

- 后验概率序列  $p(\boldsymbol{\theta}|D^n)$  收敛到 **delta** 函数;
- 只要训练样本足够多, 则  $p(\mathbf{x}|\boldsymbol{\theta})$  能唯一确定  $\boldsymbol{\theta}$ 。



在某些情况下, 不同  $\boldsymbol{\theta}$  值会产生同一个  $p(\mathbf{x}|\boldsymbol{\theta})$ 。

$p(\boldsymbol{\theta}|D^n)$  将在  $\boldsymbol{\theta}$  附近产生峰值, 这时不管  $p(\mathbf{x}|\boldsymbol{\theta})$  是否唯一,  $p(\mathbf{x}|D^n)$  总会收敛到  $p(\mathbf{x})$ 。

因此不确定性客观存在。



# 最大似然估计和贝叶斯参数估计的区别

	最大似然估计	贝叶斯参数估计
计算复杂度	微分	多重积分
可理解性	确定易理解	不确定不易理解
先验信息的信任程度	不准确	准确
例如 $p(\mathbf{x} \boldsymbol{\theta})$	与初始假设一致	与初始假设不一致



## 分类误差种类:

- 贝叶斯错误或不可分错误，例如  $P(x | \omega_i)$  之间相互重叠引起，固有问题；
- 模型错误，ML与Bays犯错一样；
- 估计错误，训练样本个数有限产生。



# Gibbs 算法

$$p(\mathbf{x} | D) = \int p(\mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | D) d\boldsymbol{\theta}$$

依据  $p(\boldsymbol{\theta} | D)$  来选择  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$

使得  $p(\mathbf{x} | D) \approx p(\mathbf{x} | \boldsymbol{\theta}_0)$  [Gibbs算法]

在较弱的假设条件下，Gibbs算法的误差概率至多是贝叶斯最优分类器的两倍。



## • 3.6 充分统计量

### ■ 统计量

□ 任何样本集  $D$  的函数；

■ 充分统计量即是一个样本集  $D$  的函数  $s$ ，其中  $s$  包含了有助于估计参数  $\theta$  的所有所有信息，即  $p(D|s, \theta)$  与  $\theta$  无关；

■ 满足上面,如果  $\theta$  是随机变量，则可以写成

$$p(\theta | s, D) = \frac{p(D | s, \theta) p(\theta | s)}{p(D | s)} = p(\theta | s)$$

反过来也成立。



## 因式分解定理:

- 一个关于参数  $\theta$  的统计量  $s$  是充分统计量当且仅当概率分布函数  $P(D|\theta)$  能够写成乘积形式:

$$P(D|\theta) = g(s, \theta) h(D)$$

其中  $g(.,.)$  和  $h(.)$  是两个函数。



## 例子：多维高斯分布

$$p(\mathbf{x} | \boldsymbol{\theta}) \sim N(\boldsymbol{\theta}, \boldsymbol{\Sigma})$$

$$p(D | \boldsymbol{\theta}) = \prod_{k=1}^n \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x}_k - \boldsymbol{\theta})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \boldsymbol{\theta}) \right]$$

$$= \exp \left[ -\frac{n}{2} \boldsymbol{\theta}^t \boldsymbol{\Sigma}^{-1} \boldsymbol{\theta} + \boldsymbol{\theta}^t \boldsymbol{\Sigma}^{-1} \left( \sum_{k=1}^n \mathbf{x}_k \right) \right] \times$$

$$\frac{1}{(2\pi)^{nd/2} |\boldsymbol{\Sigma}|^{n/2}} \exp \left[ -\frac{1}{2} \sum_{k=1}^n \mathbf{x}_k^t \boldsymbol{\Sigma}^{-1} \mathbf{x}_k \right]$$

$$\mathbf{s} = \sum_{k=1}^n \mathbf{x}_k \text{ and thus } \hat{\boldsymbol{\mu}}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k \text{ are sufficient for } \boldsymbol{\theta}$$





## 证明：必要性

假设 $\mathbf{s}$ 是关于 $\theta$ 的充分统计量，即

$P(D/\mathbf{s}, \theta)$ 不依赖于 $\theta$

$$\begin{aligned} P(D/\theta) &= \sum_{\mathbf{s}} P(D, \mathbf{s} | \theta) = \sum_{\mathbf{s}} P(D | \mathbf{s}, \theta) P(\mathbf{s} | \theta) \\ &= P(D | \mathbf{s}, \theta) P(\mathbf{s} | \theta) \end{aligned}$$

由定义  $= P(D | \mathbf{s}) P(\mathbf{s} | \theta)$

$$= h(D) P(\mathbf{s} | \theta) = h(D) g(\mathbf{s}, \theta)$$

注意到  $\mathbf{s} = \varphi(D)$ , 对于一个给定的样本，只有一个 $\mathbf{s}$ 与之对应。



充分性:

$$\mathbf{s} = \varphi(D), \quad \bar{D} = \{D \mid \varphi(D) = \mathbf{s}\}$$

$$\begin{aligned} P(D \mid \mathbf{s}, \boldsymbol{\theta}) &= \frac{P(D, \mathbf{s} \mid \boldsymbol{\theta})}{P(\mathbf{s} \mid \boldsymbol{\theta})} = \frac{P(D, \mathbf{s} \mid \boldsymbol{\theta})}{\sum_{D \in \bar{D}} P(D, \mathbf{s} \mid \boldsymbol{\theta})} \\ &= \frac{P(D \mid \boldsymbol{\theta})}{\sum_{D \in \bar{D}} P(D \mid \boldsymbol{\theta})} = \frac{g(\mathbf{s}, \boldsymbol{\theta})h(D)}{\sum_{D \in \bar{D}} g(\mathbf{s}, \boldsymbol{\theta})h(D)} = \frac{h(D)}{\sum_{D \in \bar{D}} h(D)} \end{aligned}$$

上式不依赖于  $\boldsymbol{\theta}$ ;

因此  $\mathbf{s}$  是关于  $\boldsymbol{\theta}$  的充分统计量。



## 核密度 (Kernel density)

- 把  $P(D|\theta)$  分解成  $g(s, \theta)h(D)$  不是唯一的:
  - 如果  $f(s)$  是一个函数,  $g'(s, \theta) = f(s)g(s, \theta)$  和  $h'(D) = h(D)/f(s)$  也是等价的分解;
- 这种二义性可以用定义核密度函数的方法来得到消除:

$$\bar{g}(s, \theta) = \frac{g(s, \theta)}{\int g(s, \theta) d\theta}$$



## 例子：多维高斯分布

$$p(\mathbf{x} | \boldsymbol{\theta}) \sim N(\boldsymbol{\theta}, \boldsymbol{\Sigma})$$

$$p(D | \boldsymbol{\theta}) = \exp \left[ -\frac{n}{2} \boldsymbol{\theta}^t \boldsymbol{\Sigma}^{-1} \boldsymbol{\theta} + \boldsymbol{\theta}^t \boldsymbol{\Sigma}^{-1} \left( \sum_{k=1}^n \mathbf{x}_k \right) \right] \times$$

$$\frac{1}{(2\pi)^{nd/2} |\boldsymbol{\Sigma}|^{n/2}} \exp \left[ -\frac{1}{2} \sum_{k=1}^n \mathbf{x}_k^t \boldsymbol{\Sigma}^{-1} \mathbf{x}_k \right]$$

$$= g(\hat{\boldsymbol{\mu}}_n, \boldsymbol{\theta}) h(D), \quad \hat{\boldsymbol{\mu}}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

$$g(\hat{\boldsymbol{\mu}}_n, \boldsymbol{\theta}) = \exp \left[ -\frac{n}{2} (\boldsymbol{\theta}^t \boldsymbol{\Sigma}^{-1} \boldsymbol{\theta} - 2 \boldsymbol{\theta}^t \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_n) \right]$$

$$\bar{g}(\hat{\boldsymbol{\mu}}_n, \boldsymbol{\theta}) = \frac{1}{(2\pi)^{d/2} \left| \frac{1}{n} \boldsymbol{\Sigma} \right|^{1/2}} \exp \left[ -\frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\mu}}_n)^t \left( \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} (\boldsymbol{\theta} - \hat{\boldsymbol{\mu}}_n) \right]$$



## 核密度与参数估计

- 对于最大似然估计情形，只需最大化  $g(\mathbf{s}, \boldsymbol{\theta})$ ，因为：  
 $P(D|\boldsymbol{\theta}) = g(\mathbf{s}, \boldsymbol{\theta}) h(D)$

- 对于贝叶斯估计情形：

$$p(\boldsymbol{\theta} | D) = \frac{p(D | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int p(D | \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}} = \frac{g(\mathbf{s}, \boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int g(\mathbf{s}, \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}}$$

- 如果我们对  $\boldsymbol{\theta}$  的先验概率不确定， $p(\boldsymbol{\theta})$  通常选择均匀分布，则  $p(\boldsymbol{\theta}|D)$  几乎等于核密度；
- 如果  $p(\mathbf{x}|\boldsymbol{\theta})$  可辨识时， $g(\mathbf{s}, \boldsymbol{\theta})$  通常在某个值处有明显的尖峰，并且如果  $p(\boldsymbol{\theta})$  在该值处连续并且非零，则  $p(\boldsymbol{\theta}|D)$  将趋近核密度函数。



# 充分统计量与指数族函数

$$p(\mathbf{x} | \boldsymbol{\theta}) = \alpha(\mathbf{x}) \exp[\mathbf{a}(\boldsymbol{\theta}) + \mathbf{b}(\boldsymbol{\theta})^t \mathbf{c}(\mathbf{x})]$$

$$\begin{aligned} p(D | \boldsymbol{\theta}) &= \exp\left[ n\mathbf{a}(\boldsymbol{\theta}) + \mathbf{b}(\boldsymbol{\theta})^t \sum_{k=1}^n \mathbf{c}(\mathbf{x}_k) \right] \prod_{k=1}^n \alpha(\mathbf{x}_k) \\ &= g(\mathbf{s}, \boldsymbol{\theta}) h(D) \end{aligned}$$

$$\mathbf{s} = \frac{1}{n} \sum_{k=1}^n \mathbf{c}(\mathbf{x}_k), \quad g(\mathbf{s}, \boldsymbol{\theta}) = \exp\left[ n\{\mathbf{a}(\boldsymbol{\theta}) + \mathbf{b}(\boldsymbol{\theta})^t \mathbf{s}\} \right]$$

$$h(D) = \prod_{k=1}^n \alpha(\mathbf{x}_k)$$



## • 3.7 维数问题

- 分类问题通常涉及50或100维以上的特征.
- 分类精度取决于维数和训练样本的数量
  - 考虑有相同协方差矩阵的两组多维向量情况:

$$p(\mathbf{x} | \omega_j) \sim N(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}), \quad j = 1, 2$$

如果它们的先验概率相同, 则贝叶斯误差概率为:

$$P(error) = \frac{1}{\sqrt{2\pi}} \int_{r/2}^{\infty} e^{\frac{-u^2}{2}} du$$

$$\text{其中: } r^2 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

$$\lim_{r \rightarrow \infty} P(error) = 0$$



- 如果特征是独立的，则有：

$$\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2)$$

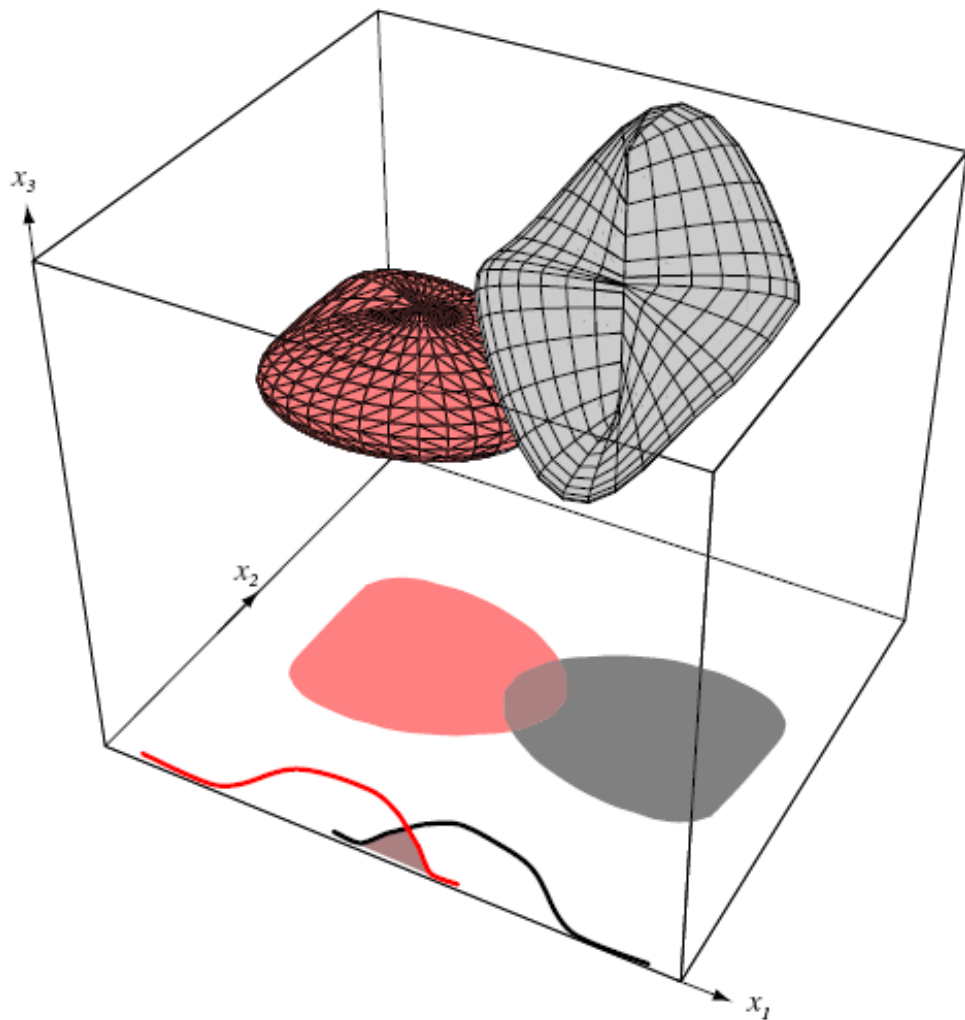
$$r^2 = \sum_{i=1}^d \left( \frac{\mu_{i1} - \mu_{i2}}{\sigma_i} \right)^2$$

- 最有用的特征是两类均值之间的距离大于标准方差的那些特征；
- 在实际观察中我们发现，当特征个数增加到某个临界点后会导致更糟糕的结果而不是好的结果：我们的模型有误，或者由于训练样本个数有限导致分布估计不精确，等等。





# 可分性与特征维数





# 学习过程的计算复杂度

$$\hat{\boldsymbol{\mu}}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k : O(nd)$$

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n)(\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n)^t : O(nd^2)$$

求解  $d \times d$  矩阵的逆 :  $O(d^3)$

求解  $d \times d$  行列式:  $O(d^3)$

$n > d$

$$g(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \hat{\boldsymbol{\mu}}_n)^t \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_n) - \frac{1}{2} \ln |\hat{\boldsymbol{\Sigma}}| + \ln P(\omega) - \frac{d}{2} \ln 2\pi$$

↓	↓	↓	↓	↓
$O(nd)$	$O(nd^2)$	$O(d^3)$	$O(n)$	$O(1)$



# 分类过程的计算复杂度

给定  $\mathbf{x}$

计算  $(\mathbf{x} - \hat{\boldsymbol{\mu}}_n)$ :  $O(d)$

将协方差矩阵的逆矩阵与差向量相乘:  $O(d^2)$

判别  $\max_i g_i(\mathbf{x})$ :  $O(c)$

整个分类问题的复杂度:  $O(d^2)$

➤ 分类阶段比学习阶段简单。



# 训练样本不足时的方法

## ■ 降维

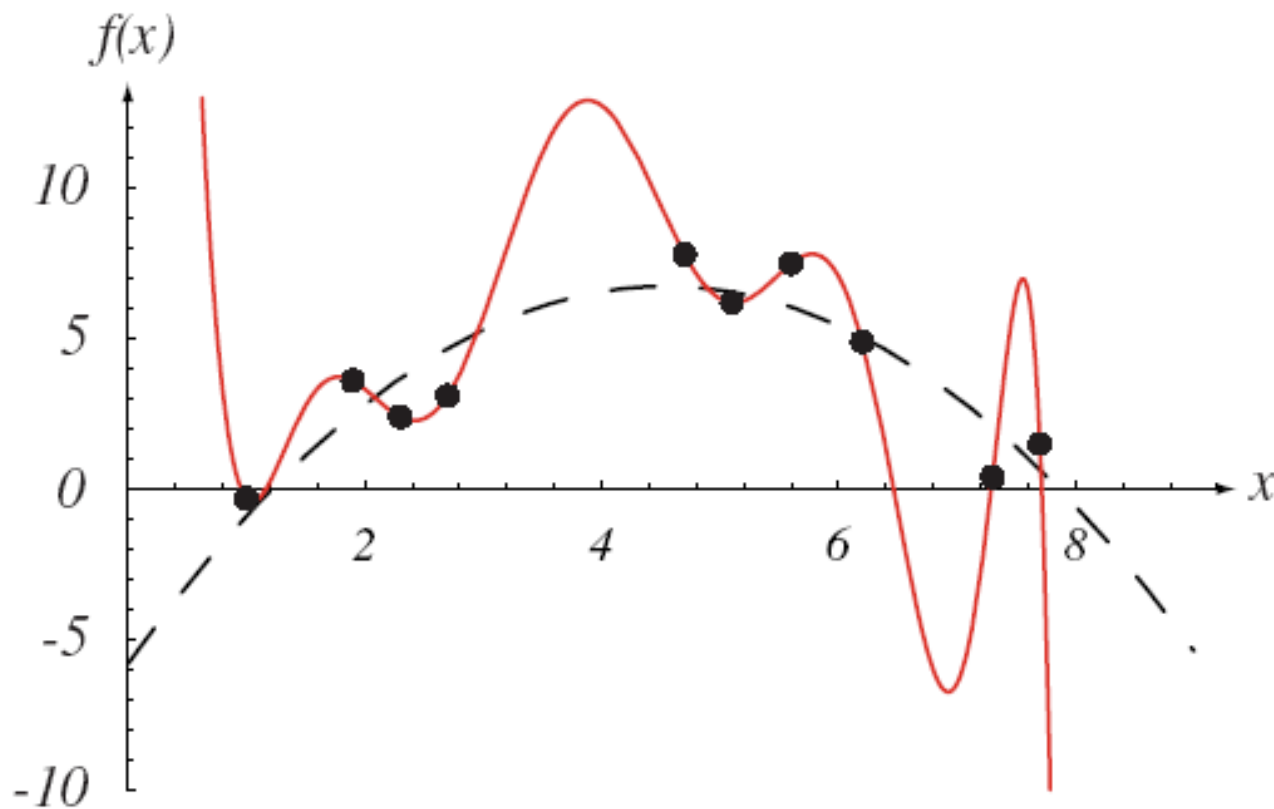
- 重新设计特征提取模块;
- 选择现有特征的子集;
- 将几个特征组合在一起;
- 假设各个类的协方差矩阵都相同, 将全部数据都归到一起;

## ■ 寻找协方差矩阵 $\Sigma$ 更好的估计;

- 如果有合理的先验估计  $\Sigma_0$ , 则可以用如下的伪贝叶斯估计  $\lambda\Sigma_0 + (1-\lambda)\hat{\Sigma}$  ;
- 设法将 $\Sigma_0$ 对角化: 阈值化或假设特征之间统计独立;



# 过拟合的概念



正确的拟合思想是：一开始用高阶的多项式曲线来拟合，然后依次去掉高阶项来逐渐简化模型，获得更光滑的结果。



## 缩并(Regularized Discriminant Analysis)

假设两类分布分别为 $N(\mu_1, \Sigma_1)$ 和 $N(\mu_2, \Sigma_2)$

$i$  为  $c$  个类中的任何一个下标,  $\Sigma$  是缩并后的协方差, 我们有:

$$\Sigma_i(\alpha) = \frac{(1-\alpha)n_i\Sigma_i + \alpha n\Sigma}{(1-\alpha)n_i + \alpha n}, \quad 0 < \alpha < 1$$

或将共同的协方差向单位矩阵缩并为

$$\Sigma(\beta) = (1-\beta)\Sigma + \beta\mathbf{I}, \quad 0 < \beta < 1$$



## • 3.8 成分分析与辨别函数

- 组合特征从而降低特征空间的维数
  - 线性组合通常比较容易计算和处理
  - 将高维数据投影到一个低维空间里去
  - 使用两种分类方法寻找理想一点的线性变换:
    - PCA (主成份分析) “在最小均方意义下的数据的最优表示的映射”
    - MDA (多类判别分析) “在最小均方意义下的数据的最优分类的映射”



# 主成分分析

给定  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , 找  $\mathbf{x}_0$  使得

$$J_0(\mathbf{x}_0) = \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{x}_0\|^2 \text{ 最小}$$

容易证明,  $\mathbf{x}_0 = \mathbf{m}$

$$\text{其中, } \mathbf{m} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

$$\begin{aligned} J_0(\mathbf{x}_0) &= \sum_{k=1}^n \|(\mathbf{x}_k - \mathbf{m}) - (\mathbf{x}_0 - \mathbf{m})\|^2 \\ &= \sum_{k=1}^n \|\mathbf{x}_0 - \mathbf{m}\|^2 + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2 \end{aligned}$$

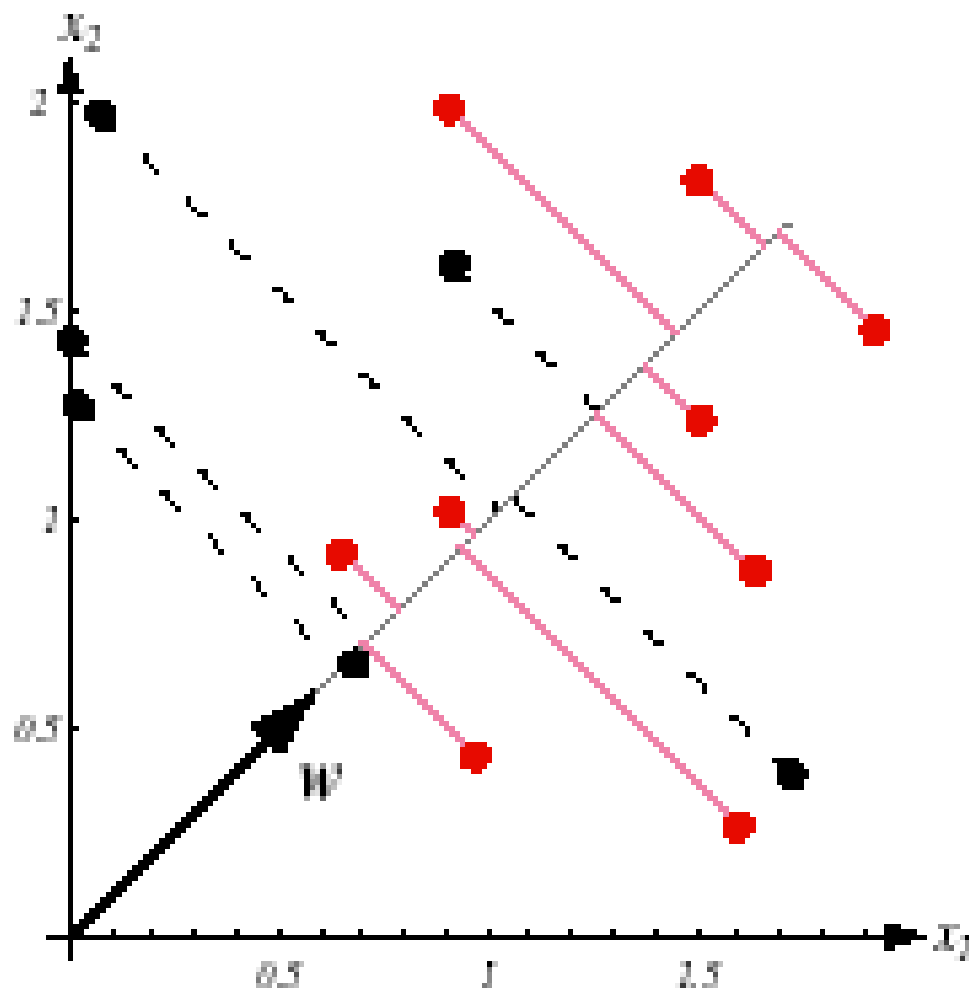
该项不依赖  $\mathbf{x}_0$

当  $\mathbf{x}_0 = \mathbf{m}$  使得 minimizes  $J_0(\mathbf{x}_0)$





# 沿直线投影:





# 对于通过样本均值直线的最佳投影

直线:  $\mathbf{x} = \mathbf{m} + a\mathbf{e}$

如果用  $\mathbf{m} + a_k\mathbf{e}$  表达  $\mathbf{x}_k$  , 那么通过最小化平方误差准则函数, 可以得到一组最优  $a_k$  的集合, 过程如下:

$$\begin{aligned} J_1(a_1, \dots, a_n; \mathbf{e}) &= \sum_{k=1}^n \|(\mathbf{m} + a_k\mathbf{e}) - \mathbf{x}_k\|^2 \\ &= \sum_{k=1}^n a_k^2 \|\mathbf{e}\|^2 - 2 \sum_{k=1}^n a_k \mathbf{e}^t (\mathbf{x}_k - \mathbf{m}) + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2 \end{aligned}$$

由于  $\|\mathbf{e}\| = 1$ , 通过对  $a_k$  求偏导数, 并且令结果为0, 得到:

$$a_k = \mathbf{e}^t (\mathbf{x}_k - \mathbf{m})$$



# 寻找最佳表达方向

寻找  $\mathbf{e}$ ，最小化下面过程：

$$\begin{aligned} J_1(\mathbf{e}) &= \sum_{k=1}^n a_k^2 - 2 \sum_{k=1}^n a_k^2 + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2 \\ &= - \sum_{k=1}^n \mathbf{e}^t \underbrace{(\mathbf{x}_k - \mathbf{m})(\mathbf{x}_k - \mathbf{m})^t}_{\mathbf{S}} \mathbf{e} + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2 \\ &= -\mathbf{e}^t \mathbf{S} \mathbf{e} + \sum_{k=1}^n \|\mathbf{x}_k - \mathbf{m}\|^2 \end{aligned}$$

其中,  $\mathbf{S}$  为散布矩阵 (Scatter matrix)

即是最大化  $\mathbf{e}^t \mathbf{S} \mathbf{e}$ , 满足  $\|\mathbf{e}\|^2 = 1$

Lagrange 方法: 最大化  $u = \mathbf{e}^t \mathbf{S} \mathbf{e} - \lambda(\mathbf{e}^t \mathbf{e} - 1)$

$$\nabla_{\mathbf{e}} u = 0 \Rightarrow \frac{\partial u}{\partial \mathbf{e}} = 2\mathbf{S}\mathbf{e} - 2\lambda\mathbf{e} = 0 \Rightarrow \mathbf{S}\mathbf{e} = \lambda\mathbf{e}$$



# 主成分分析 (PCA)

## —Principal component analysis

由1维推广到 $d'$ 维:

$$\text{投影空间: } \mathbf{x} = \mathbf{m} + \sum_{i=1}^{d'} a_i \mathbf{e}_i$$

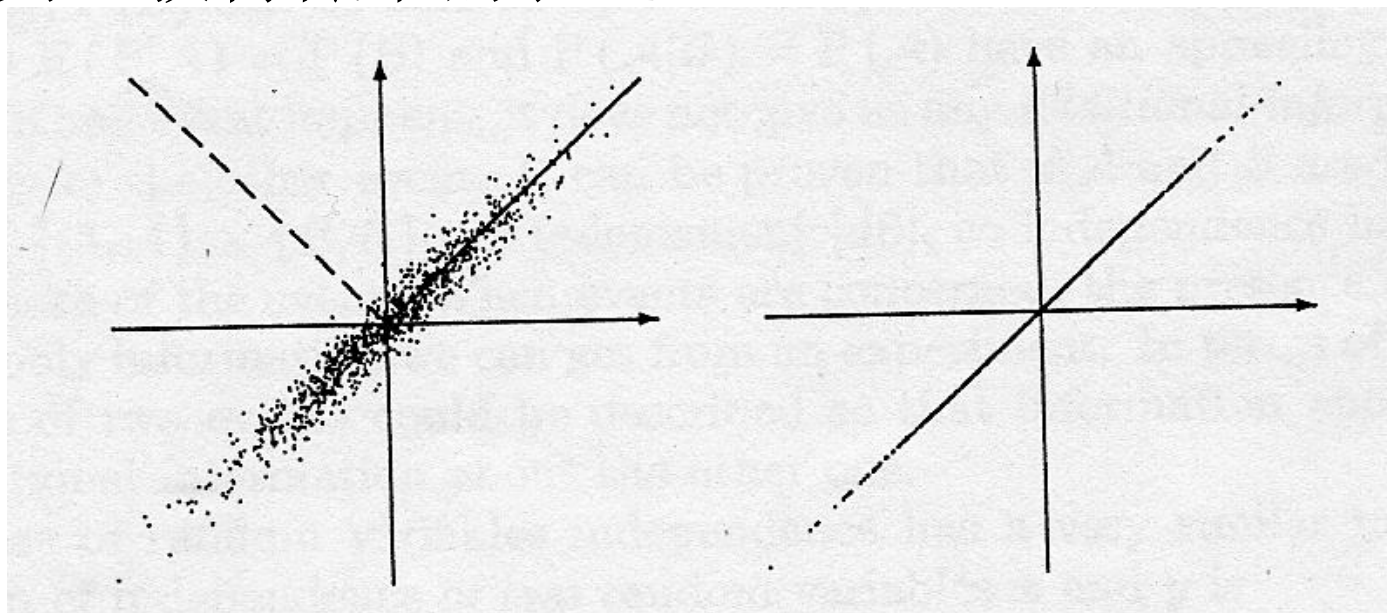
寻找  $\mathbf{e}_i$ ,  $i = 1, \dots, d'$  来最小化下式:

$$J_{d'}(\mathbf{e}_1, \dots, \mathbf{e}_{d'}) = \sum_{k=1}^n \left\| \left( \mathbf{m} + \sum_{i=1}^{d'} a_{ki} \mathbf{e}_i \right) - \mathbf{x}_k \right\|^2$$

$\Rightarrow \mathbf{e}_1, \dots, \mathbf{e}_{d'}$  是 $\mathbf{S}$ 的  $d'$  个具有最大特征值 的特征向量



- $L$ 个 $N$ 维空间的向量，构成 $N$ 维空间的 $L$ 个点。如果大多数点落在一个 $M$ 维超平面上，只要能找到 $M$ 维空间的坐标系，则可以将 $L$ 个向量投影到 $M$ 维空间，获得低维的表达。



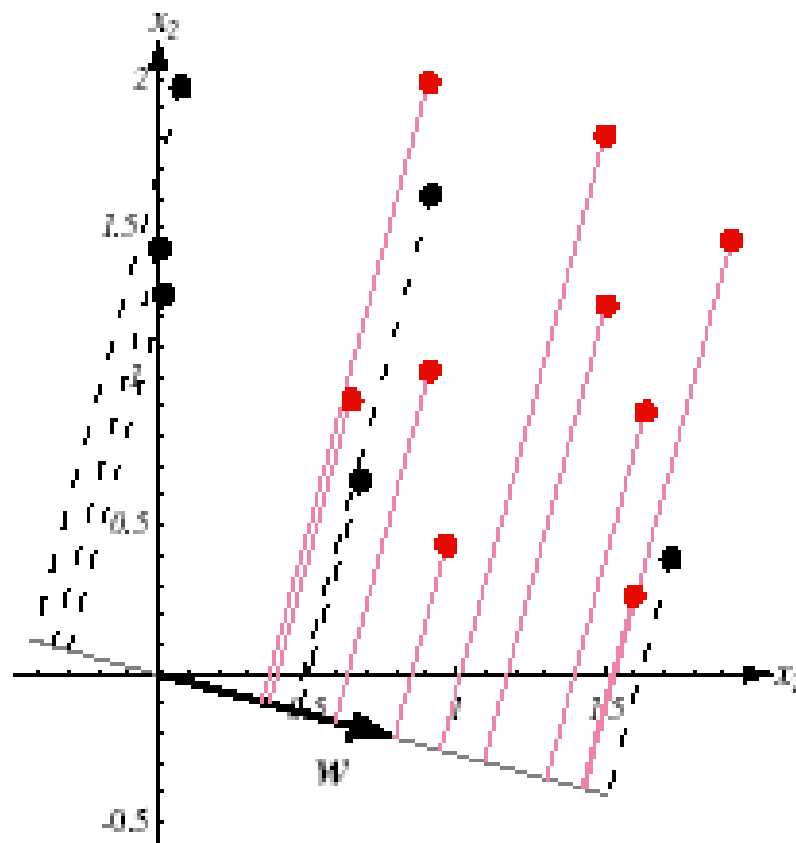
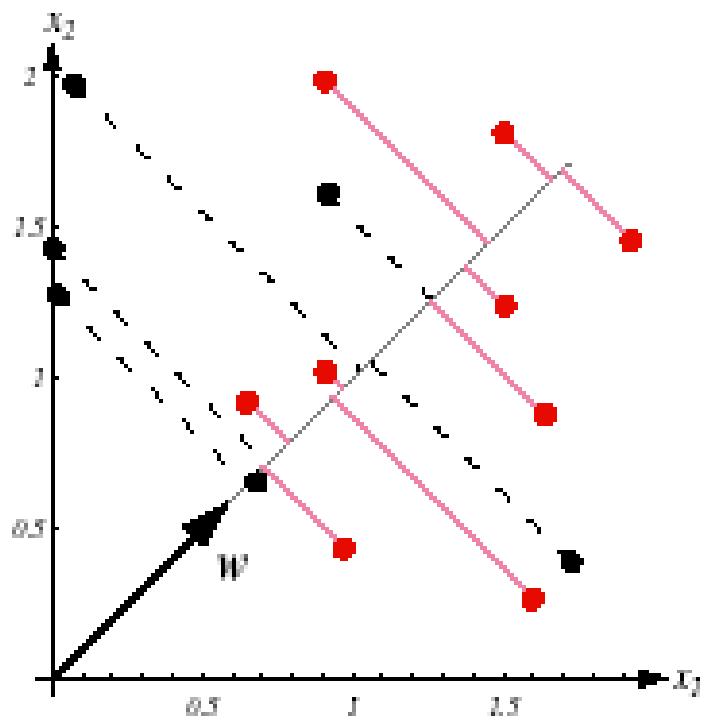
K-L变换

PCA

- K-L变换是压缩与特征提取的有效方法。



# Fisher 线性分类的概念



➤ 以“O”、“Q”为例，比较PCA与LDA的差别。



# Fisher 线性鉴别分析

## —Fisher Linear Discriminant Analysis

考虑两类样本的分类问题:

有一组 $n$ 个 $d$ 维样本:  $\mathbf{x}_1, \mathbf{x}_1, \dots, \mathbf{x}_n$

其中 $n_1$ 个样本属于 $\omega_1$ ,  $n_2$ 个样本属于 $\omega_2$ 。

找  $\mathbf{w}$ , 最大化分离  $\mathbf{y} = \mathbf{w}^t \mathbf{x}$

$$\mathbf{m}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in D_i} \mathbf{x}, i = 1, 2$$

$$\tilde{m}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in D_i} \mathbf{w}^t \mathbf{x} = \mathbf{w}^t \mathbf{m}_i, \quad \tilde{s}_i^2 = \sum_{\mathbf{x} \in D_i} (\mathbf{w}^t \mathbf{x} - \tilde{m}_i)^2$$

内类散度矩阵:  $\tilde{s}_1^2 + \tilde{s}_2^2$

$$\text{最大化 } J(\mathbf{w}) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$



# Fisher Linear Discriminant Analysis

$$\tilde{s}_i^2 = \sum_{x \in D_i} (\mathbf{w}^t \mathbf{x} - \mathbf{w}^t \mathbf{m}_i)^2 = \mathbf{w}^t \mathbf{S}_i \mathbf{w}$$

$$\mathbf{S}_i = \sum_{x \in D_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^t$$

$$\tilde{s}_1^2 + \tilde{s}_2^2 = \mathbf{w}^t \mathbf{S}_W \mathbf{w}, \mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$$

$$|\tilde{m}_1 - \tilde{m}_2|^2 = (\mathbf{w}^t \mathbf{m}_1 - \mathbf{w}^t \mathbf{m}_2)^2 = \mathbf{w}^t \mathbf{S}_B \mathbf{w}$$

$$\mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^t$$





$$J(\mathbf{w}) = \frac{\mathbf{w}^t \mathbf{S}_B \mathbf{w}}{\mathbf{w}^t \mathbf{S}_W \mathbf{w}}, \text{ 称为 generalized Rayleigh quotient}$$

(广义瑞利商),

使 $J(\mathbf{w})$ 达到最大的  $\mathbf{w}$  必须满足:

$$\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w} \text{ (generalized eigenvalue problem)}$$

$$\mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{w}$$

$\mathbf{S}_B \mathbf{w} = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^t \mathbf{w}$  总是位于  $(\mathbf{m}_1 - \mathbf{m}_2)$  方向上

$$\therefore \mathbf{w} = \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \text{ [忽略常数]}$$



## 对于正态分布的LDA

假设各类的协方差矩阵  $\Sigma$  相同

最佳判决边界的方程为:

$$\mathbf{w}^t \mathbf{x} + w_0 = 0 \quad P31$$

$$\text{其中, } \mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2)$$

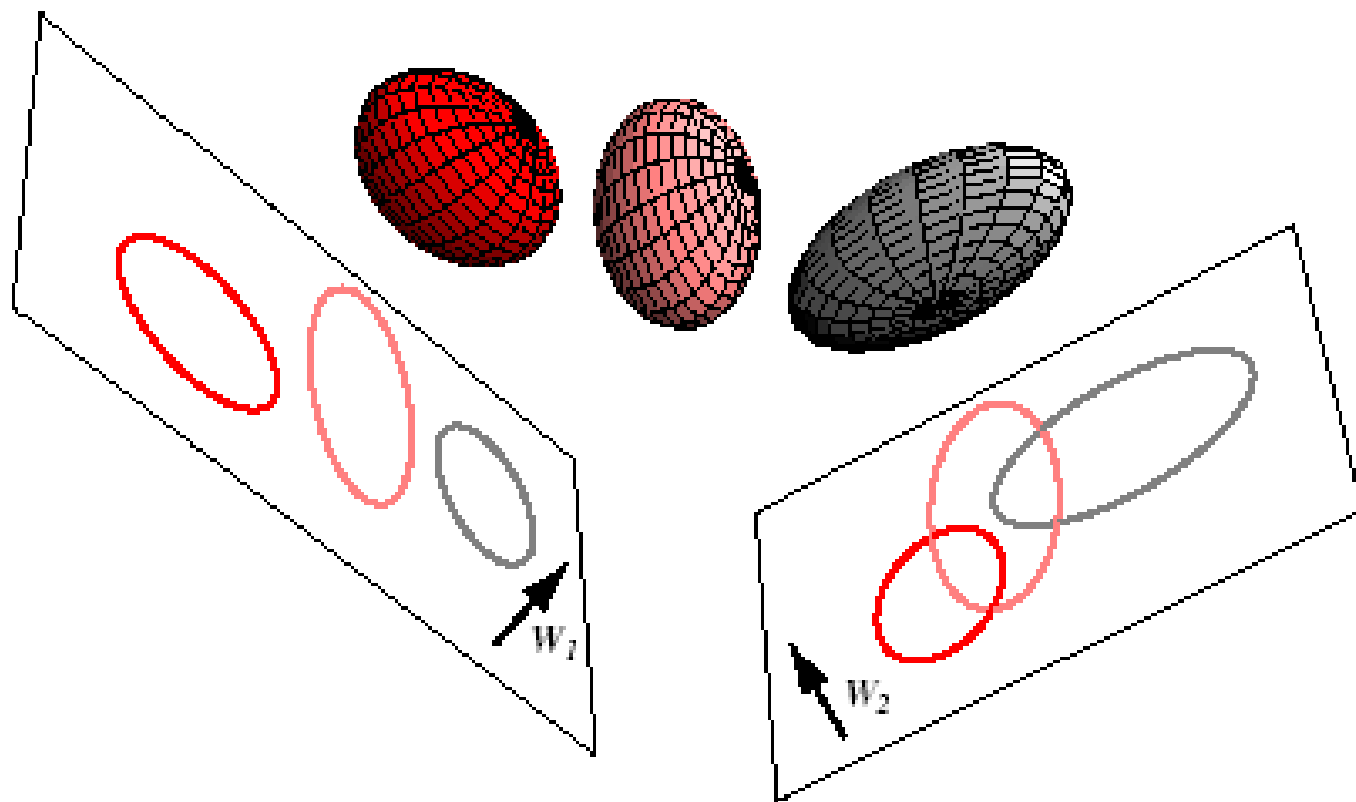
如果用样本估计  $\mu_1, \mu_2$ , 和  $\Sigma$ , 则

$$\mathbf{w} = \mathbf{S}_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$

[这就是Fisher linear discriminant analysis与贝叶斯准则的关系]



# 多重判别分析—MDA





# Multiple Discriminant Analysis

考虑  $c$  类的分类问题

将  $d$  维空间投影到  $(c-1)$  维子空间 ( $d > c$ )

$$y_i = \mathbf{w}_i^t \mathbf{x}, \quad i = 1, \dots, c-1 \Rightarrow \mathbf{y} = \mathbf{W}^t \mathbf{x}$$

$$\tilde{\mathbf{m}}_i = \frac{1}{n} \sum_{\mathbf{x} \in D_i} \mathbf{W}^t \mathbf{x}, \quad \tilde{\mathbf{m}} = \frac{1}{n} \sum_{i=1}^c n_i \tilde{\mathbf{m}}_i$$

$$\tilde{\mathbf{S}}_W = \sum_{i=1}^c \sum_{\mathbf{x} \in D_i} (\mathbf{W}^t \mathbf{x} - \tilde{\mathbf{m}}_i)(\mathbf{W}^t \mathbf{x} - \tilde{\mathbf{m}}_i)^t = \mathbf{W}^t \mathbf{S}_W \mathbf{W}$$

$$\mathbf{S}_W = \sum_{i=1}^c \mathbf{S}_i, \quad \mathbf{S}_i = \sum_{\mathbf{x} \in D_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^t, \quad \mathbf{m}_i = \frac{1}{n} \sum_{\mathbf{x} \in D_i} \mathbf{x}$$



$$\mathbf{m} = \frac{1}{n} \sum_{\mathbf{x}} \mathbf{x} = \frac{1}{n} \sum_{i=1}^c n_i \mathbf{m}_i$$

$$\mathbf{S}_T = \sum_{\mathbf{x}} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^t$$

$$= \sum_{i=1}^c \sum_{\mathbf{x} \in D_i} (\mathbf{x} - \mathbf{m}_i + \mathbf{m}_i - \mathbf{m})(\mathbf{x} - \mathbf{m}_i + \mathbf{m}_i - \mathbf{m})^t$$

$$= \sum_{i=1}^c \sum_{\mathbf{x} \in D_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^t + \sum_{i=1}^c \sum_{\mathbf{x} \in D_i} (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^t$$

$$= \mathbf{S}_W + \sum_{i=1}^c n_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^t = \mathbf{S}_W + \mathbf{S}_B$$

$$\tilde{\mathbf{S}}_B = \sum_{i=1}^c n_i (\tilde{\mathbf{m}}_i - \tilde{\mathbf{m}})(\tilde{\mathbf{m}}_i - \tilde{\mathbf{m}})^t = \mathbf{W}^t \mathbf{S}_B \mathbf{W}$$



我们希望寻找  $\mathbf{W}$  使得 between-class scatter 与 within-class scatter 的比最大。

用散度矩阵行列式度量的情形：

$$\therefore \text{let } J(\mathbf{W}) = \frac{|\tilde{\mathbf{S}}_B|}{|\tilde{\mathbf{S}}_W|} = \frac{|\mathbf{W}^t \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^t \mathbf{S}_W \mathbf{W}|}$$

$\mathbf{W}$  的列满足：

$$\mathbf{S}_B \mathbf{w}_i = \lambda_i \mathbf{S}_W \mathbf{w}_i$$

这是相对大特征值的广义特征向量。最优的  $\mathbf{W}$  是不唯一的, 对坐标轴进行适当地旋转和伸缩, 对准则函数和最后的分类器没有影响。



# 期望最大化 (EM)

- 将最大似然估计推广到允许包含丢失特征样本来学习特定分布的参数问题;
- 完整的样本集  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$
- $\mathbf{x}_k = \{ \mathbf{x}_{kg}, \mathbf{x}_{kb} \}$
- 把不同的特征分成两部分  $D_g$  和  $D_b$ 
  - $D$  是  $D_g$  和  $D_b$  的并集
- 组成函数

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^i) = E_{D_b} \left[ \ln p(D_g, D_b; \boldsymbol{\theta}) \mid D_g; \boldsymbol{\theta}^i \right]$$

左边是一个关于  $\boldsymbol{\theta}$  的函数,  $\boldsymbol{\theta}^i$  描述整个分布的真实参数。  
也可以看作参数向量  $\boldsymbol{\theta}^i$  是当前对整个分布的最好估计,  
而  $\boldsymbol{\theta}$  是在当前估计基础上进一步改善的估计。  
右边表示关于丢失的特征求对数似然函数的期望。

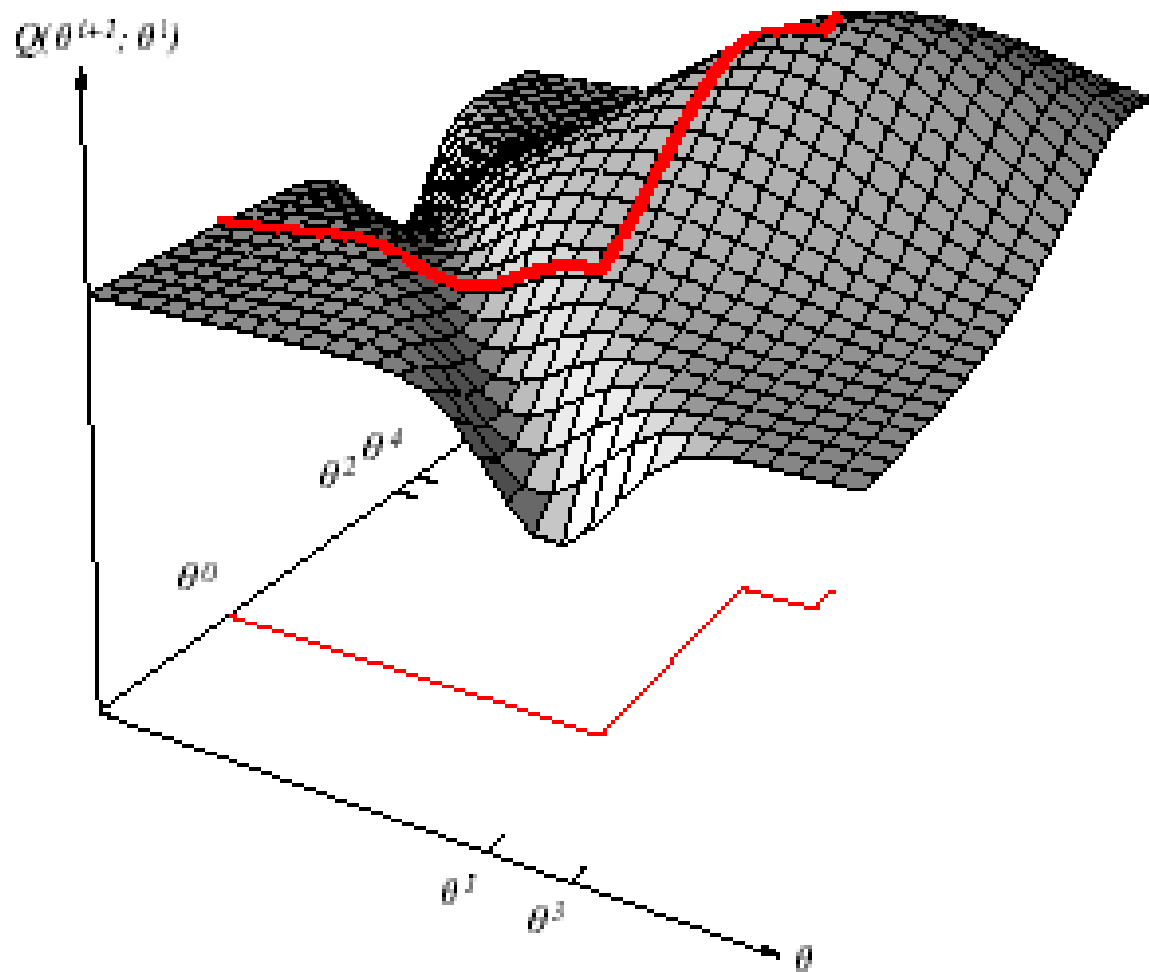


```
begin initialize  $\theta^0, T, i \leftarrow 0$   
  do  $i \leftarrow i + 1$   
    E step: Compute  $Q(\theta, \theta^i)$   
    M step:  $\theta^{i+1} \leftarrow \arg \max_{\theta} Q(\theta, \theta^i)$   
  until  $Q(\theta^{i+1}; \theta^i) - Q(\theta^i; \theta^{i-1}) \leq T$   
  return  $\theta \leftarrow \theta^{i+1}$   
end
```





# Expectation-Maximization (EM)





## Example: 2D 模型

$$D = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4\} = \left\{ \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} * \\ 4 \end{pmatrix} \right\}$$

$$D_b = x_{41}$$

假设 2D 高斯模型具有对角斜方差阵

$$\boldsymbol{\theta} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \sigma_1^2 \\ \sigma_2^2 \end{pmatrix}, \quad \boldsymbol{\theta}^0 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}$$



$$\begin{aligned} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^0) &= E_{x_{41}} \left[ \ln p(D_g, x_{41}; \boldsymbol{\theta}) \mid D_g, \boldsymbol{\theta}^0 \right] \\ &= \int_{-\infty}^{\infty} \left[ \sum_{k=1}^3 \ln p(\mathbf{x}_k \mid \boldsymbol{\theta}) + \ln p(\mathbf{x}_4 \mid \boldsymbol{\theta}) \right] \times \\ &\quad p(x_{41} \mid \boldsymbol{\theta}^0; x_{42} = 4) dx_{41} \\ &= \sum_{k=1}^3 \ln p(\mathbf{x}_k \mid \boldsymbol{\theta}) + \int_{-\infty}^{\infty} \ln p\left(\begin{pmatrix} x_{41} \\ 4 \end{pmatrix} \mid \boldsymbol{\theta}\right) \frac{p\left(\begin{pmatrix} x_{41} \\ 4 \end{pmatrix} \mid \boldsymbol{\theta}^0\right)}{K} dx_{41} \\ K &= \int_{-\infty}^{\infty} p\left(\begin{pmatrix} x'_{41} \\ 4 \end{pmatrix} \mid \boldsymbol{\theta}^0\right) dx'_{41} \end{aligned}$$



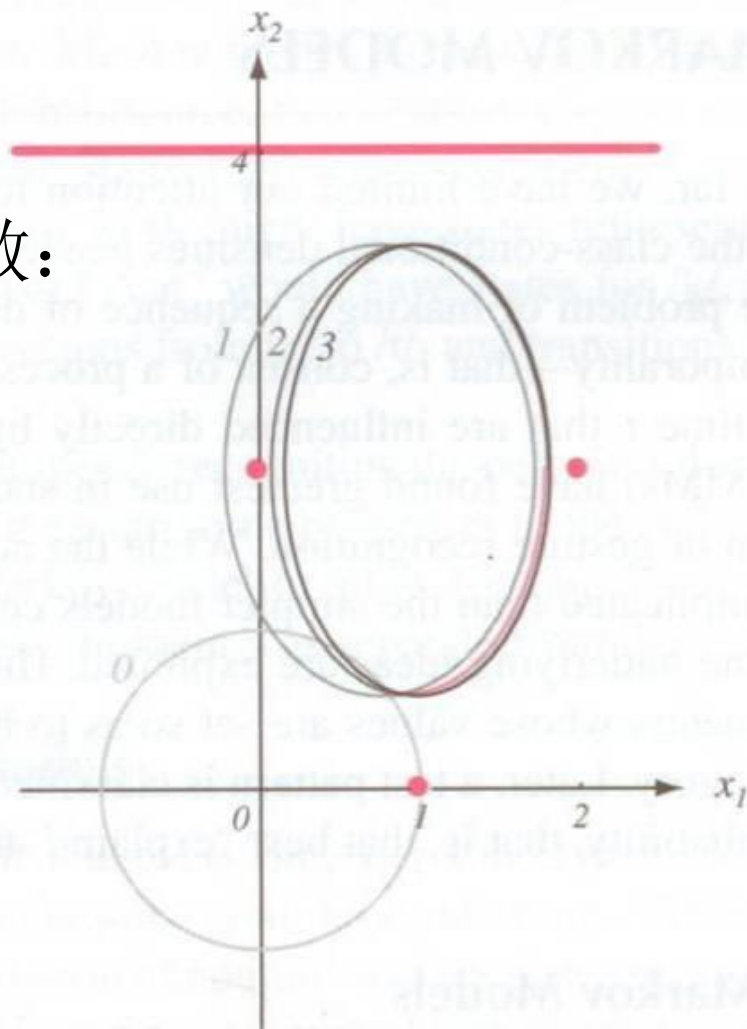
$$\begin{aligned} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^0) &= \sum_{k=1}^3 \ln p(\mathbf{x}_k | \boldsymbol{\theta}) + \\ &\quad \frac{1}{K} \int_{-\infty}^{\infty} \ln p\left(\begin{pmatrix} x_{41} \\ 4 \end{pmatrix} | \boldsymbol{\theta}\right) \frac{1}{2\pi} \exp\left[-\frac{1}{2}(x_{41}^2 + 4^2)\right] dx_{41} \\ &= \sum_{k=1}^3 \ln p(\mathbf{x}_k | \boldsymbol{\theta}) - \frac{1 + \mu_1^2}{2\sigma_1^2} - \frac{(4 - \mu_2)^2}{2\sigma_2^2} - \ln(2\pi\sigma_1\sigma_2) \\ \boldsymbol{\theta}^1 &= \begin{pmatrix} 0.75 \\ 2.0 \\ 0.938 \\ 2.0 \end{pmatrix} \end{aligned}$$



3 迭代后, 该算法在下列值收敛:

$$\mu = \begin{pmatrix} 1.0 \\ 2.0 \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} 0.667 & 0 \\ 0 & 2.0 \end{pmatrix}$$





## 广义期望最大化 (GEM)

- 代替最大化  $Q(\theta, \theta^i)$ , 我们在M步只需要找  $\theta^{i+1}$  使得

$$Q(\theta^{i+1}; \theta^i) > Q(\theta; \theta^i)$$

也能确保收敛。

- 收敛将没有那么快。
- 让用户自由选取计算更加简单的途径。
  - 有一种版本的**GEM**算法, 每次叠代时, 都计算未知特征的最大似然函数, 然后依此重新计算 $\theta$ 。



# 隐马尔可夫模型— Hidden Markov Model (HMM)

- 前面各章节,用一个 $n$ 维特征矢量确定一个对象的状态,并基于这个状态进行统计判决;
- 本节,用一个时间的(矢量)序列或空间的(矢量)阵列来描述对象的整体状态,并基于这个整体状态进行统计判决;
- 用于处理序列判决问题
  - 应用, 在语音和手势识别方面有用。
- 在  $t$  时刻发生的事件要收到 $t-1$ 时刻发生事件的直接影响。



# First Order Markov Models

## 一阶马尔可夫模型

假设在 $t$ 时刻的状态记为 $\omega(t)$ 。

有一个时间长度为 $T$ 的状态序列：

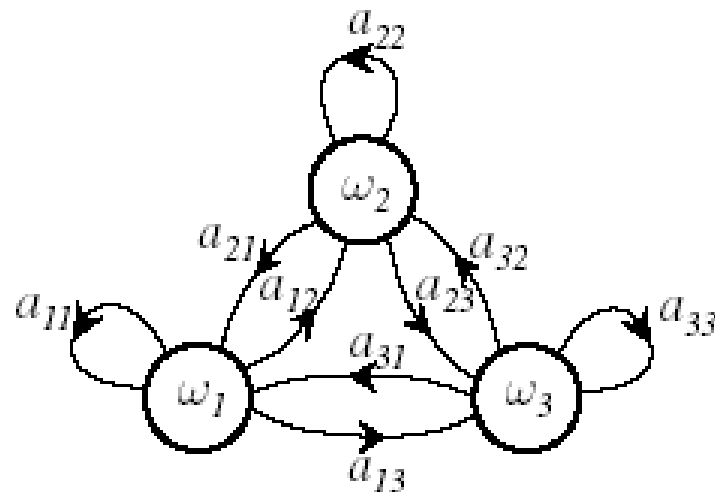
状态序列  $\boldsymbol{\omega}^T = \{\omega(1), \omega(2), \dots, \omega(T)\}$

e.g.,  $\boldsymbol{\omega}^6 = \{\omega_1, \omega_3, \omega_2, \omega_2, \omega_1, \omega_3\}$ ,

$$P(\boldsymbol{\omega}^6 | \boldsymbol{\theta}) = a_{13} a_{32} a_{22} a_{21} a_{13}$$

其中记 $P(\omega_j(t+1) | \omega_i(t)) = a_{ij}$ 为转移概率。

特定序列的概率为连续的转移概率相乘。







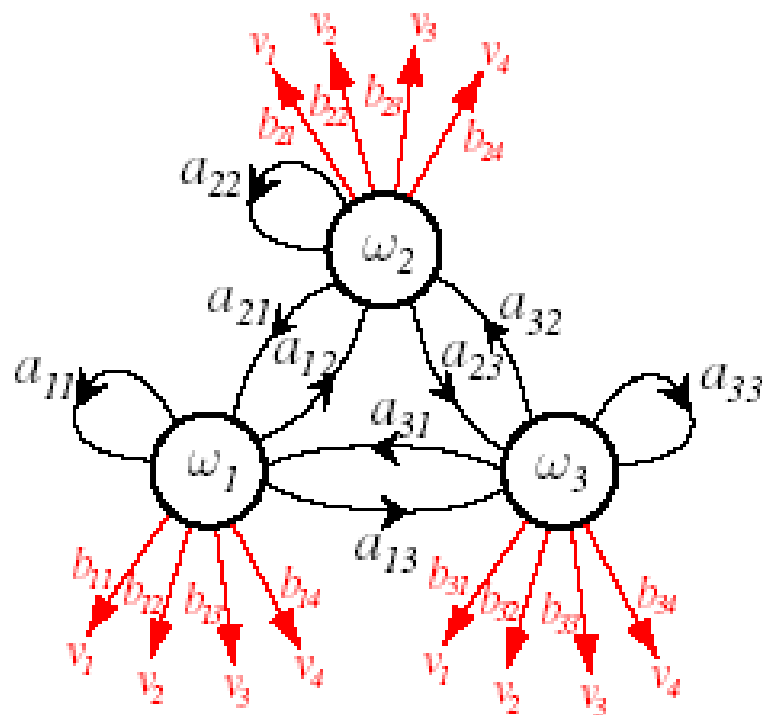
# First Order Hidden Markov Models

## 一阶隐马尔可夫模型

假定在任一时刻状态 $\omega(t)$ 都以某种概率产生

可见状态序列  $\mathbf{V}^T = \{v(1), v(2), \dots, v(T)\}$  中的一个。

e.g.,  $\mathbf{V}^6 = \{v_4, v_1, v_1, v_4, v_2, v_3\}$ ,  $P(v_k(t) | \omega_j(t)) = b_{jk}$





# Hidden Markov Model 概率

最终或吸收状态  $\omega_0 : a_{00} = 1$

转移概率:  $a_{ij} = P(\omega_j(t+1) | \omega_i(t))$

激发可见状态的概率:

$$b_{jk} = P(v_k(t) | \omega_j(t))$$

$$\sum_j a_{ij} = 1, \quad \sum_k b_{jk} = 1$$



## 一阶隐马尔可夫模型的例子:

下面给出一个天气的例子。设一天的天气可能是晴天、多云、雨天,其只随机(但有统计性)地依赖于已过去的一天的天气,天气状态的转移概率  $a_{ij}$  如下所示:

$$\mathbf{A} = (a_{ij}) = \begin{matrix} & \begin{matrix} \text{晴天} & \text{多云} & \text{雨天} \end{matrix} \\ \begin{matrix} \text{晴天} \\ \text{多云} \\ \text{雨天} \end{matrix} & \begin{pmatrix} 0.50 & 0.375 & 0.125 \\ 0.25 & 0.125 & 0.625 \\ 0.25 & 0.375 & 0.375 \end{pmatrix} \end{matrix}$$

可以看出,晴天后是雨天的概率是 0.125,雨天后是多云的概率是 0.375,等等。一片海藻的湿润情况是天气的反映,它们间的概率如下所示:

$$\mathbf{B} = (b_{jk}) = \begin{matrix} & \begin{matrix} \text{干燥} & \text{稍干} & \text{潮湿} & \text{浸水} \end{matrix} \\ \begin{matrix} \text{晴天} \\ \text{多云} \\ \text{雨天} \end{matrix} & \begin{pmatrix} 0.60 & 0.20 & 0.15 & 0.05 \\ 0.25 & 0.25 & 0.25 & 0.25 \\ 0.05 & 0.10 & 0.35 & 0.50 \end{pmatrix} \end{matrix}$$

晴天时海藻干燥的概率是 0.6,多云天气时海藻潮湿的概率是 0.25,等等。若以每天作为观测时刻,假设只能观测海藻的湿润情况,这种两层统计系统就是一个隐马尔可夫系统。



# Hidden Markov Model 的计算

## ■ 估值问题

- 利用给定的  $a_{ij}$  和  $b_{jk}$ , 计算某个特定观察序列  $\mathbf{V}^T$  的概率  $P(\mathbf{V}^T|\boldsymbol{\theta})$ 。

## ■ 解码问题

- 给定特定观察序列  $\mathbf{V}^T$ , 决定最有可能产生  $\mathbf{V}^T$  的隐状态序列  $\omega^T$ 。

## ■ 学习问题

- 已知HMM的大致结构（如隐状态和可见状态的数目），但  $a_{ij}$  和  $b_{jk}$  未知，如何从一组可见符号的训练集中，决定这些参数。

## ■ 运用HMM模型识别

- 利用各类的可见序列样本进行学习,产生代表每类的HMM参考模型; 对待识别可见序列,通过估值方法进行识别。



## Evaluation(估值问题)

$$P(\mathbf{V}^T) = \sum_{r=1}^{r_{\max}} P(\mathbf{V}^T | \boldsymbol{\omega}_r^T) P(\boldsymbol{\omega}_r^T), \quad \text{其中 } r_{\max} = c^T$$

$$P(\boldsymbol{\omega}_r^T) = \prod_{t=1}^T P(\omega(t) | \omega(t-1))$$

$$P(\mathbf{V}^T | \boldsymbol{\omega}_r^T) = \prod_{t=1}^T P(v(t) | \omega(t))$$

$$P(\mathbf{V}^T) = \sum_{r=1}^{r_{\max}} \prod_{t=1}^T P(v(t) | \omega(t)) P(\omega(t) | \omega(t-1))$$

计算复杂度  $O(c^T T)$



# HMM Forward(前向算法)

$$P(\mathbf{V}^T) = \sum_{r=1}^{r_{\max}} \prod_{t=1}^T P(v(t) | \omega(t)) P(\omega(t) | \omega(t-1))$$

$$\alpha_j(t) = P(\omega_j(t), \mathbf{V}^t)$$

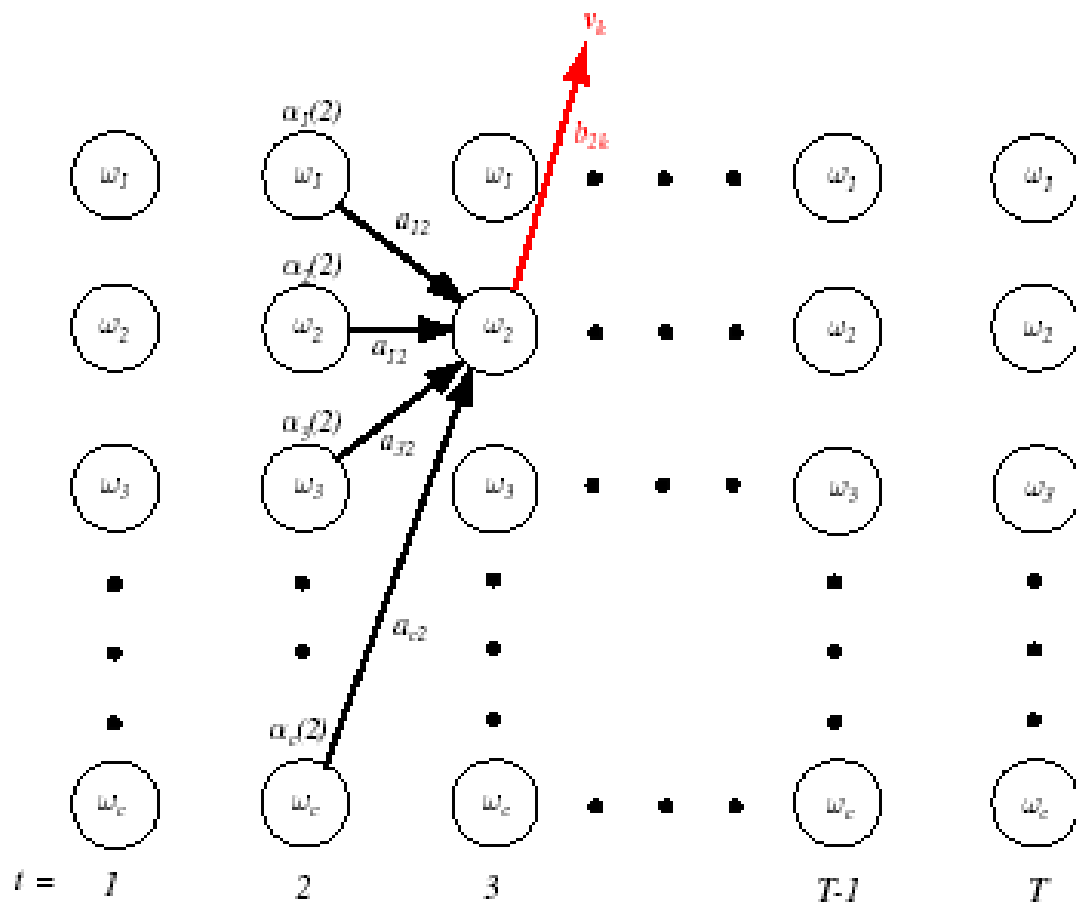
$$= P(v_k | \omega_j(t)) \sum_{i=1}^c P(\omega_j(t) | \omega_i(t-1)) P(\omega_i(t-1), \mathbf{V}^{t-1})$$

$$= b_{jkv(t)} \sum_{i=1}^c a_{ij} \alpha_i(t-1)$$

$$\alpha_j(0) = \begin{cases} 0, & j \neq \text{initial state} \\ 1, & j = \text{initial state} \end{cases}$$



# HMM Forward





# HMM Forward

initialize  $t \leftarrow 0, a_{ij}, b_{jk}, \mathbf{V}^T, \alpha_j(0) = 1$

for  $t \leftarrow t + 1, j = 0, 1, \dots, c$

$$\alpha_j(t) \leftarrow b_{jkv(t)} \sum_{i=1}^c \alpha_i(t-1) a_{ij}$$

until  $t = T$

return  $P(\mathbf{V}^T) = \alpha_0(T)$  for final state

end





# HMM Backward

$$P(\mathbf{V}^T) = \sum_{r=1}^{r_{\max}} \prod_{t=1}^T P(v(t) | \omega(t)) P(\omega(t) | \omega(t-1))$$

$$\beta_i(t) = P(\omega_i(t), \mathbf{V}^{T-t})$$

$$= \sum_{j=1}^c P(v_k | \omega_j(t+1)) P(\omega_j(t+1) | \omega_i(t)) P(\omega_j(t+1), \mathbf{V}^{T-(t+1)})$$

$$= \sum_{j=1}^c b_{jkv(t+1)} a_{ij} \beta_j(t+1)$$

$$\beta_i(T) = \begin{cases} 0, & i \neq 0 \\ 1, & i = 0 \end{cases}$$



# HMM Backward

initialize  $t \leftarrow T, a_{ij}, b_{jk}, \mathbf{V}^T, \beta_j(T)$

for  $t \leftarrow t - 1, i = 0, 1, \dots, c$

$$\beta_i(t) \leftarrow \sum_{j=0}^c \beta_j(t+1) a_{ij} b_{jkv(t+1)}$$

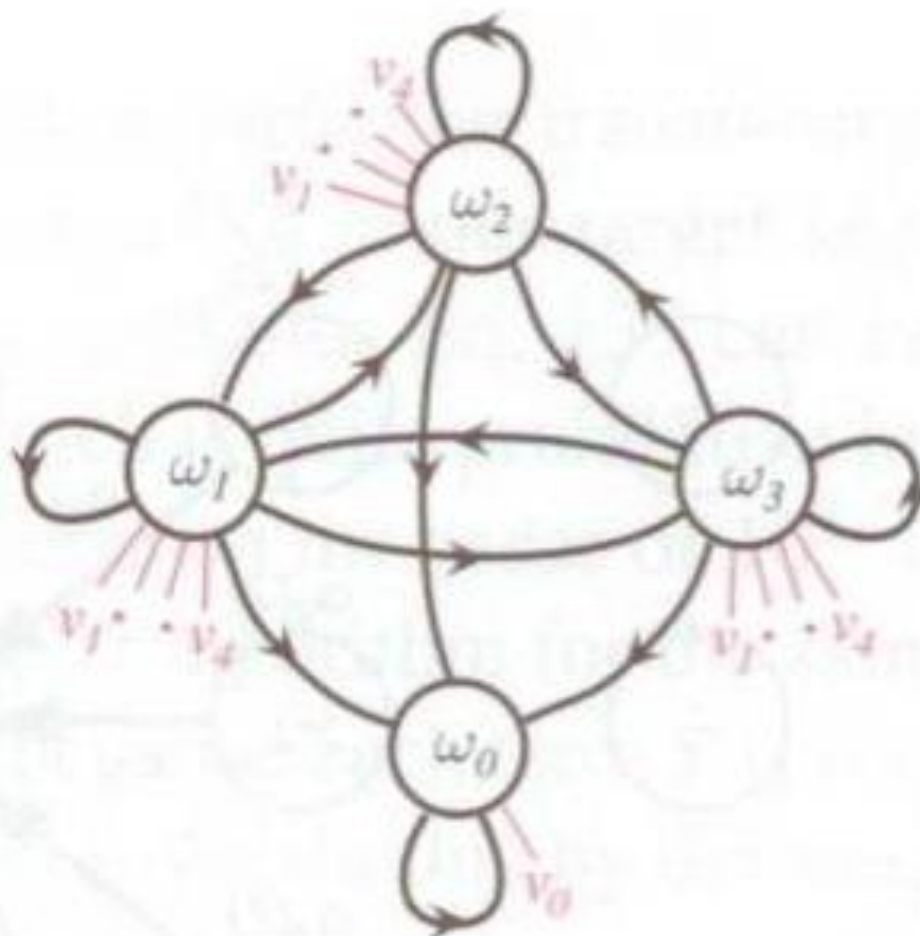
until  $t = 0$

return  $P(\mathbf{V}^T) = \beta_0(0)$  for initial state

end



# Example 3: Hidden Markov Model





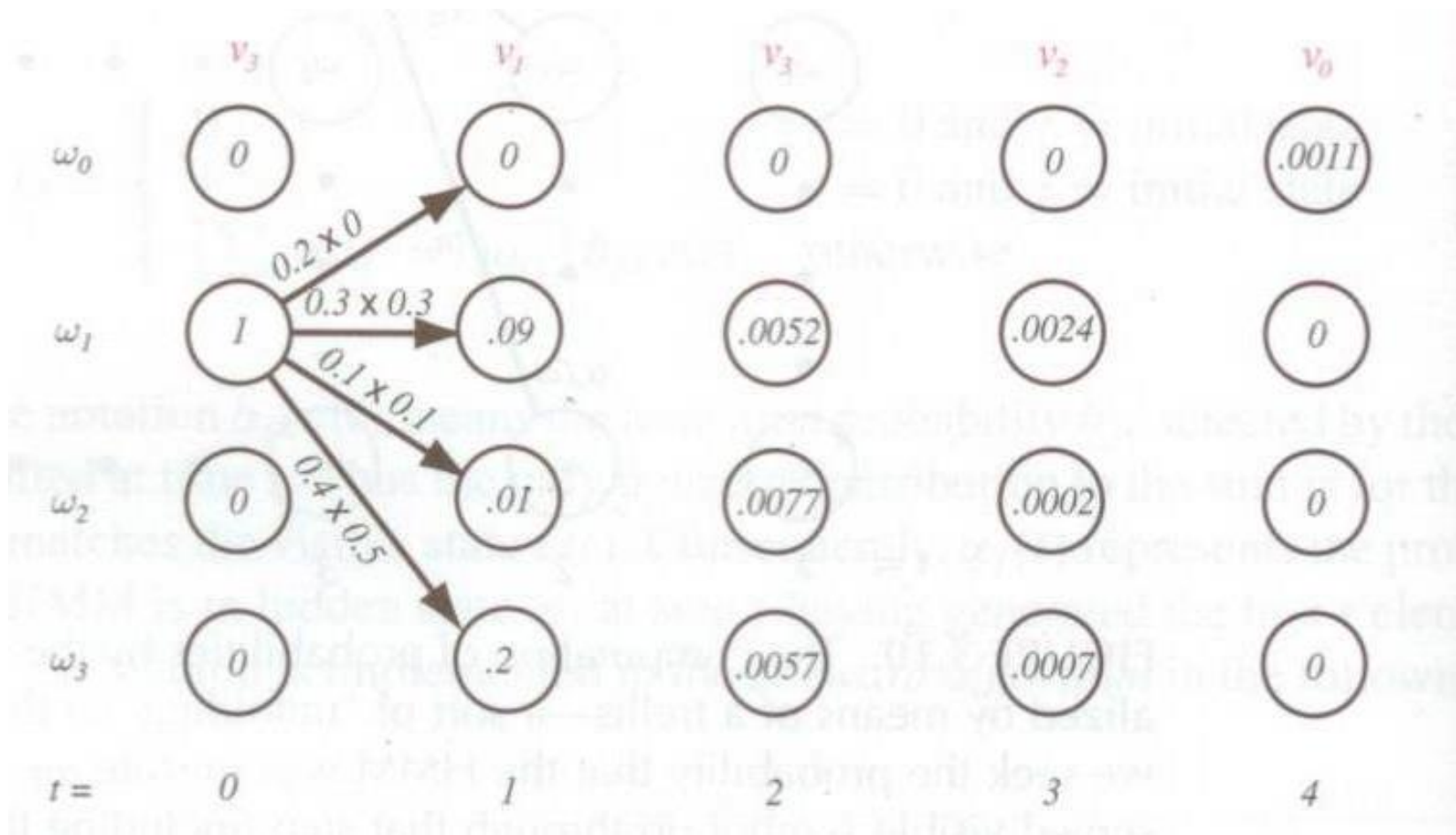
# Example 3: Hidden Markov Model

$$a_{ij} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0.2 & 0.3 & 0.1 & 0.4 \\ 0.2 & 0.5 & 0.2 & 0.1 \\ 0.8 & 0.1 & 0.0 & 0.1 \end{pmatrix}$$

$$b_{jk} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0.3 & 0.4 & 0.1 & 0.2 \\ 0 & 0.1 & 0.1 & 0.7 & 0.1 \\ 0 & 0.5 & 0.2 & 0.1 & 0.2 \end{pmatrix}$$

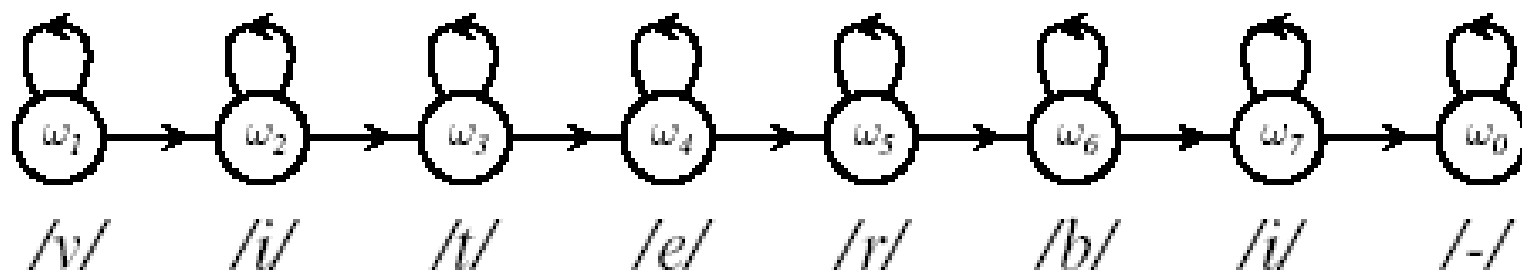


# Example 3: Hidden Markov Model





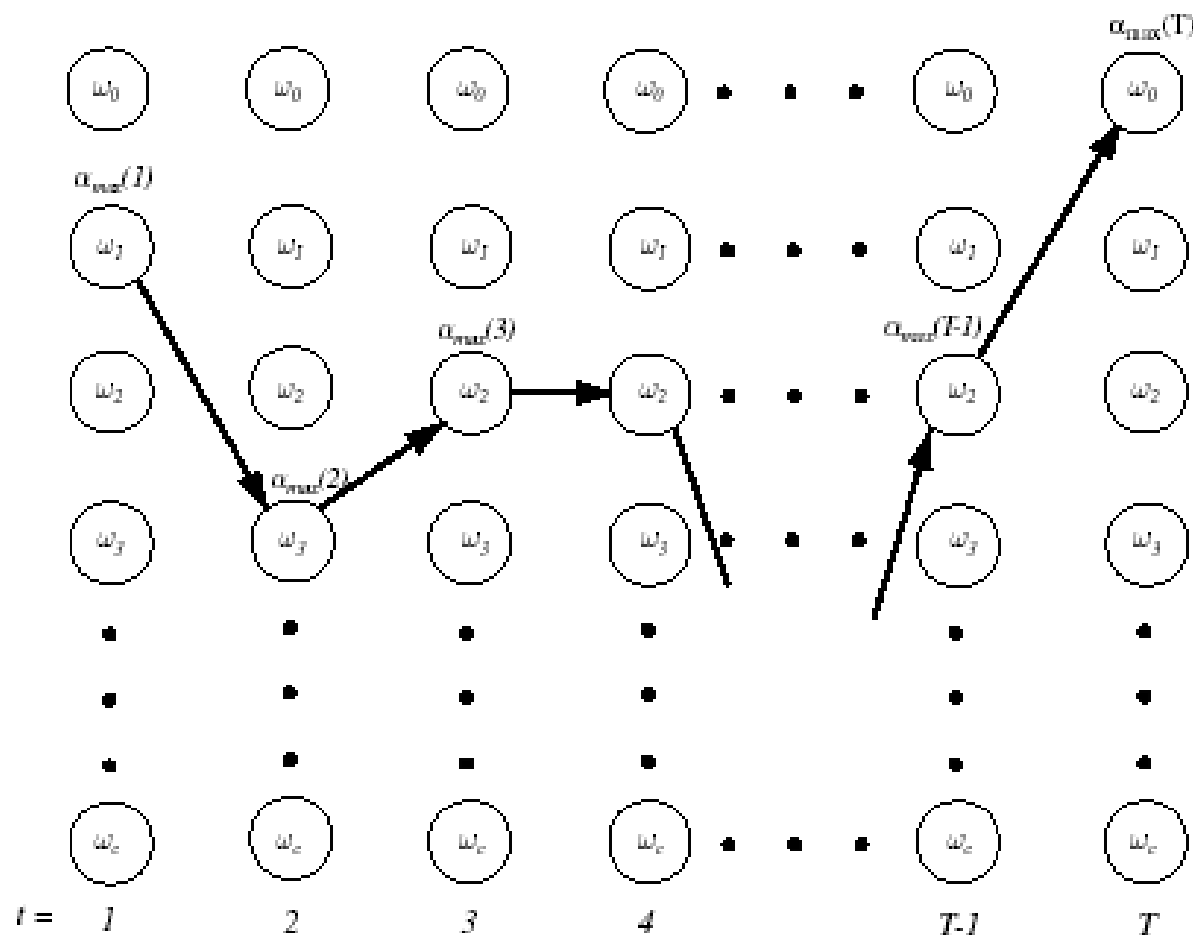
# Left-to-Right Models for Speech



$$P(\boldsymbol{\theta} | \mathbf{V}^T) = \frac{P(\mathbf{V}^T | \boldsymbol{\theta})P(\boldsymbol{\theta})}{P(\mathbf{V}^T)}$$



# HMM Decoding (解码)





# Problem of Local Optimization

- 局部最优：对于每个 $t$ 时刻，都寻找从前状态转移过来，并且产生可见状态 $v_k$ 的概率最大的概率。
- 算法本身不保证整个路径是合法的。



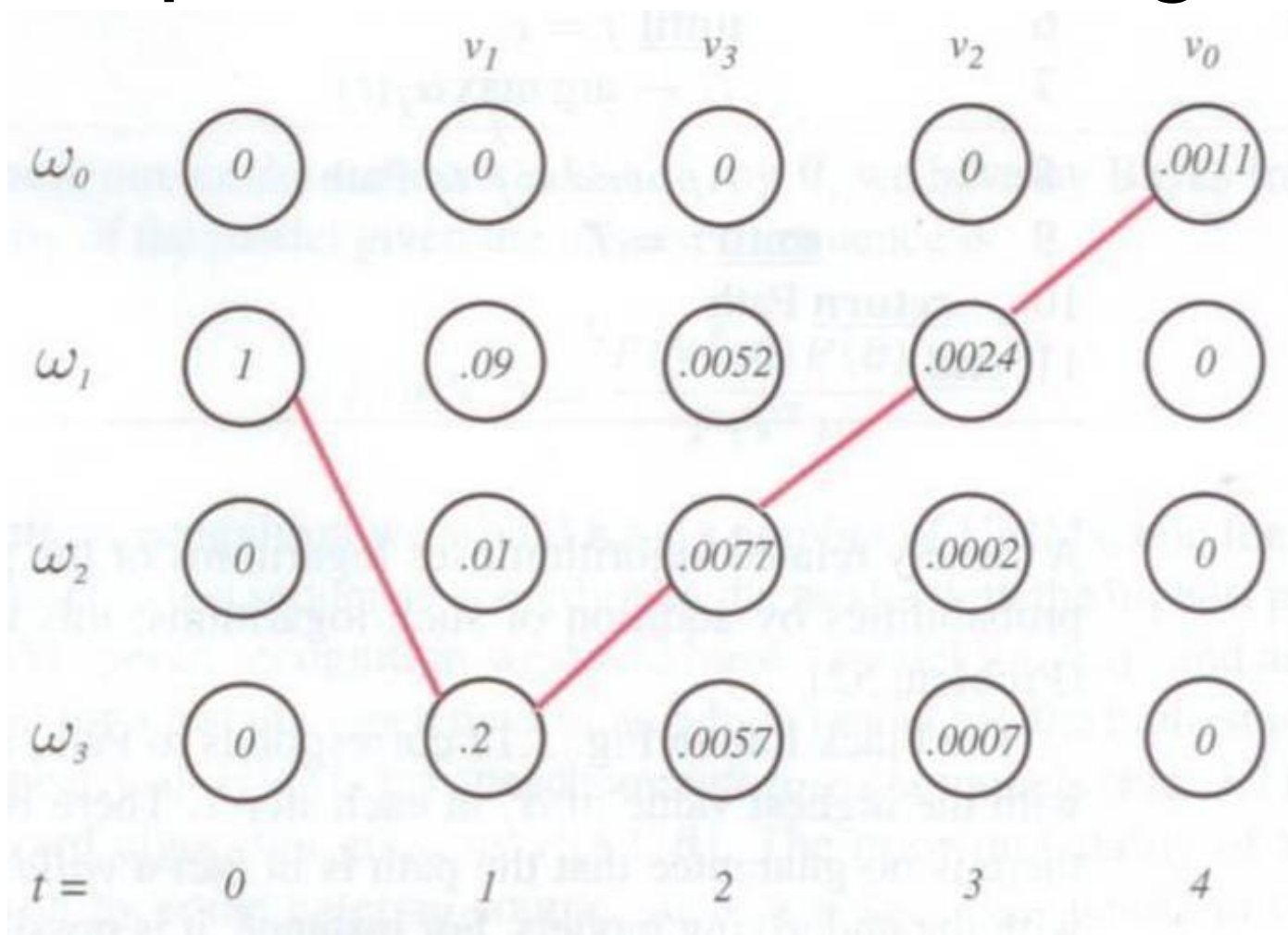


# HMM Decoding

```
initialize  $t \leftarrow 0, Path = \{ \}$   
for  $t \leftarrow t + 1$ ,  
     $j \leftarrow 0$   
    for  $j \leftarrow j + 1$   
         $\alpha_j(t) \leftarrow b_{j k_v(t)} \sum_{i=1}^c \alpha_i(t-1) a_{ij}$   
    until  $j = c$   
     $j' \leftarrow \arg \max_j \alpha_j(t)$   
    Append  $\omega_{j'}$  to  $Path$   
until  $t = T$   
return  $Path$   
end
```



# Example 4: HMM Decoding





# Forward-Backward Algorithm

- 从一组训练样本中,确定模型参数  $a_{ij}$  和  $b_{jk}$ 。
- “前向-后向算法”是“广义期望最大化算法”的具体实现
- 通过递归方法更新权重,以得到能够更好地解释训练样本序列的模型参数。



## Probability of Transition

$$\begin{aligned}\gamma_{ij}(t) &= P(\omega_i(t-1), \omega_j(t) \mid \mathbf{V}^T, \boldsymbol{\theta}) \\ &= \frac{P(\omega_i(t-1), \omega_j(t), \mathbf{V}^T \mid \boldsymbol{\theta})}{P(\mathbf{V}^T \mid \boldsymbol{\theta})} \\ &= \frac{\alpha_i(t-1)a_{ij}b_{jk}\beta_j(t)}{P(\mathbf{V}^T \mid \boldsymbol{\theta})}\end{aligned}$$



## Improved Estimate for $a_{ij}$

在任何时刻，序列中从状态  $\omega_i(t-1)$  到  $\omega_j(t)$  的转换估计值：

$$\sum_{t=1}^T \gamma_{ij}(t)$$

从  $\omega_i$  的任何转移的总预期数为：  $\sum_{t=1}^T \sum_k \gamma_{ik}(t)$

$a_{ij}$  的估计值为：

$$\hat{a}_{ij} = \frac{\sum_{t=1}^T \gamma_{ij}(t)}{\sum_{t=1}^T \sum_k \gamma_{ik}(t)}$$



同理，可以获得  $b_{jk}$  的估计值

$$\hat{b}_{jk} = \frac{\sum_{t=1, v(t)=v_k}^T \sum_i \gamma_{ij}(t)}{\sum_{t=1}^T \sum_i \gamma_{ij}(t)}$$



# Forward-Backward Algorithm (Baum-Welch Algorithm)

initialize  $a_{ij}, b_{jk}$ , training sequence  $\mathbf{V}^T$ , threshold  $\theta, z \leftarrow 0$

do  $z \leftarrow z + 1$

compute all  $\hat{a}_{ij}(z)$  from all  $a_{ij}(z-1)$  and  $b_{jk}(z-1)$

compute all  $\hat{b}_{jk}(z)$  from all  $a_{ij}(z-1)$  and  $b_{jk}(z-1)$

$a_{ij}(z) \leftarrow \hat{a}_{ij}(z)$

$b_{jk}(z) \leftarrow \hat{b}_{jk}(z)$

until  $\max_{i,j,k} [a_{ij}(z) - a_{ij}(z-1), b_{jk}(z) - b_{jk}(z-1)] < \theta$

return  $a_{ij} \leftarrow a_{ij}(z); b_{jk} \leftarrow b_{jk}(z)$

end