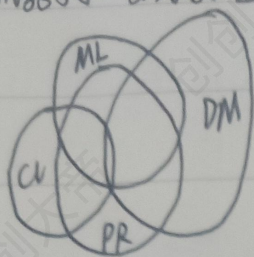


# 模式识别

绪论

模式识别 - 机器学习 - 数据挖掘 - 计算机视觉



CV: 模式识别的分支, 实现视觉信息高层理解

ML: 主要研究通用机器学习算法, 大部分针对分类

PR: 主要研究分类识别方法面向实际应用

DM: 针对各种数据中的信息提取和知识发现

② 模式的两个层次: 样本和类别

模式识别核心技术: 模式分类

识别对象表示: 特征

③ 评价性能: 测试数据集上的分类性能

④ 为什么需要结构方法: 1° 表示模式的内部结构 (字符的笔划)

2° 长度/大小不固定的模式

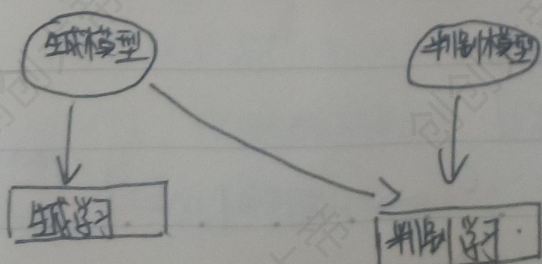
3° 相互关联的多个部件/部件同时分类

⑤ 生成模型: 表示每个类别内部结构或特征分布  $P(X|C)$

判别模型: 表示不同类别之间的区别, 一般为判别函数、边界函数或后验概率函数

生成学习: 得到每个类别的内部结构描述或分布函数, 不同类别分别学习

判别学习: 得到判别函数或边界函数或参数, 所有类别样本同时学习





MTWTFSS

## 第7章 贝叶斯决策理论

先验概率  $P(W_i)$   $\sum_{i=1}^C P(W_i) = 1$ 联合概率密度函数 (条件概率)  $P(X|W_i)$ 后验概率  $P(W_i|X) = \frac{P(X|W_i)P(W_i)}{P(X)} = \frac{P(X|W_i)P(W_i)}{\sum_{j=1}^C P(X|W_j)P(W_j)}$ if  $P(X|W_1)P(W_1) > P(X|W_2)P(W_2)$  选择  $W_1$ 

最小风险决策:

条件风险  $R(\alpha_i|X) = \sum_{j=1}^C \lambda(\alpha_i|W_j)P(W_j|X)$ 

二分类时

$$R(\alpha_1|X) = \lambda_{11}P(W_1|X) + \lambda_{12}P(W_2|X) \quad R(\alpha_2|X) = \lambda_{21}P(W_1|X) + \lambda_{22}P(W_2|X)$$

$$R(\alpha_1|X) < R(\alpha_2|X) \iff \lambda_{11}P(W_1|X) + \lambda_{12}P(W_2|X) < \lambda_{21}P(W_1|X) + \lambda_{22}P(W_2|X)$$

$$(\lambda_{21} - \lambda_{11})P(W_1|X) > (\lambda_{12} - \lambda_{22})P(W_2|X)$$

$$(\lambda_{21} - \lambda_{11})P(X|W_1)P(W_1) > (\lambda_{12} - \lambda_{22})P(X|W_2)P(W_2)$$

$$\text{选择 } W_1, \quad \frac{P(X|W_1)}{P(X|W_2)} > \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \frac{P(W_2)}{P(W_1)}$$

$$(4) \lambda(\alpha_i|W_j) = \begin{cases} 0 & i=j \\ 1 & i \neq j \end{cases} \quad 0-1 \text{ loss}$$

$$R(\alpha_i|X) = \sum_{j=1}^C \lambda(\alpha_i|W_j)P(W_j|X)$$

$$= \sum_j P(W_j|X)$$

$$= 1 - P(W_i|X) \Rightarrow \text{可以看出, 决策代价越大, 被选择概率越小.}$$

最小误差决策  $\iff$  最大后验概率 (MAP)

(5) 带拒绝的决策

$$C+1 \text{ class: } \lambda(\alpha_i|W_j) = \begin{cases} 0 & i=j \\ \lambda_s & i \neq j \\ \lambda_r & \text{reject} \end{cases} \quad \lambda_s > \lambda_r$$

$$\text{条件风险 } R(\alpha_i|X) = \sum_{j=1}^C \lambda(\alpha_i|W_j)P(W_j|X)$$

$$= \lambda_s [1 - P(W_i|X)]$$

$$R_i(X) = \begin{cases} \lambda_s [1 - P(W_i|X)] & i=1, \dots, C \\ \lambda_r & \text{reject} \end{cases}$$

$$\Rightarrow \arg \min_i R_i(X) = \begin{cases} \arg \max_i P(W_i|X) & \text{if } \max_i P(W_i|X) > 1 - \lambda_r \\ \text{reject} & \text{otherwise} \end{cases}$$

⑧ 判别函数:  $\arg \max g_i(x)$  最大后验概率决策

⑨ 例题:  $g(x) \equiv g_1(x) - g_2(x)$

$$g(x) = P(w_1|x) - P(w_2|x) = 0.$$

$$P(x|w_1)P(w_1) - P(x|w_2)P(w_2) = (\log P(x|w_1) + \log P(w_1)) - (\log P(x|w_2) + \log P(w_2))$$

$$= \ln \frac{P(x|w_1)}{P(x|w_2)} + \ln \frac{P(w_1)}{P(w_2)}$$

⑩ 概率密度估计方法: 1° 参数法  $P(x|w_i) = P(x|\theta_i)$

2° 非参数法 可表示任意概率分布. 无函数形式 (K-NN)

3° 半参数法 近似任意概率分布. 有函数形式 (GM, 高斯混合 EM 期望最大)

⑪ 高斯函数 (正态分布)

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] \quad (\text{在给定均值和方差的所有分布中, 正态分布熵最大})$$

⑫ 多元正态分布  $\rightarrow$  要考

$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left[-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right]$$

$$\text{均值: } \mu = E[x] = \int x p(x) dx \quad \mu_i = E[x_i]$$

$$\text{协方差矩阵: } \Sigma = E[(x-\mu)(x-\mu)^T] = \int (x-\mu)(x-\mu)^T p(x) dx.$$

⑬ 协方差矩阵的本征值和本征向量.

特征值      特征向量

⑭ PCA: 一种降维 (特征提取) 方法. 其将维数向量投影到低维子空间, 使子空间投影的重建误差最小. 其选择本征值最大的  $M$  个本征向量作为子空间的基.

⑮ 高斯密度下的判别函数:  $\rightarrow$  看例题

$$g_i(x) = \ln P(x|w_i) + \ln P(w_i)$$

$$P(x|w_i) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp\left[-\frac{1}{2}(x-\mu_i)^T \Sigma_i^{-1}(x-\mu_i)\right]$$

$$g_i(x) = -\frac{1}{2}(x-\mu_i)^T \Sigma_i^{-1}(x-\mu_i) - \frac{n}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(w_i)$$

情况 1:  $\Sigma_i = \sigma^2 I$  (去掉与类别无关项)

$$g_i(x) = -\frac{\|x-\mu_i\|^2}{2\sigma^2} + \ln P(w_i)$$

$$\text{展开 } \|x-\mu_i\|^2 \Rightarrow g_i(x) = w_i^T x + w_{i0} \quad w_i = \frac{1}{\sigma^2} \mu_i \quad w_{i0} = \frac{1}{2\sigma^2} \mu_i^T \mu_i + \ln P(w_i)$$

情况2:  $\Sigma_i = \Sigma$

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma^{-1}(x - \mu_i) + \ln P(W_i)$$

$$g_i(x) = W_i^T x + W_{i0}$$

$$W_i = \Sigma^{-1} \mu_i \quad W_{i0} = -\frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i + \ln P(W_i)$$

情况3:  $\Sigma_i = \Sigma_i$

$$g_i(x) = x^T W_i x + W_i^T x + W_{i0}$$

$$W_i = -\frac{1}{2} \Sigma_i^{-1} \quad W_i = \Sigma_i^{-1} \mu_i \quad W_{i0} = -\frac{1}{2} \mu_i^T \Sigma_i^{-1} \mu_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(W_i)$$

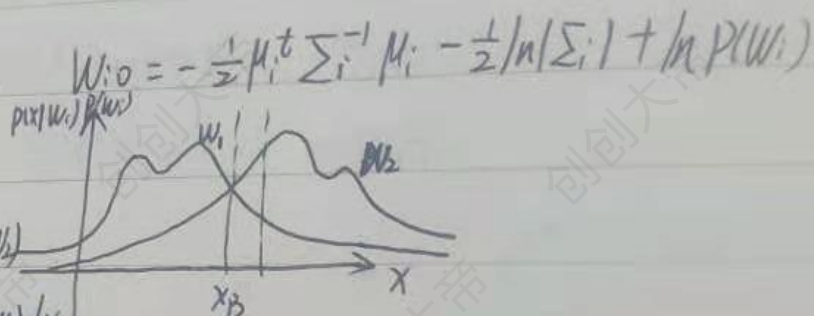
(14) 2类分错错误情况

$$P(\text{error}) = P(x \in R_2 | W_1) \cdot P(W_1) + P(x \in R_1 | W_2) \cdot P(W_2)$$

$$= \int_{R_2} P(x | W_1) P(W_1) dx + \int_{R_1} P(x | W_2) P(W_2) dx$$

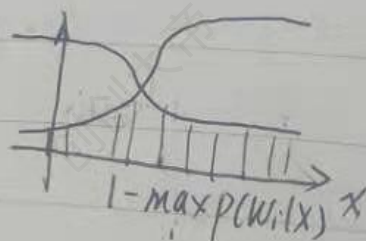
$$P(\text{correct}) = \sum_{i=1}^C \int_{R_i} P(x | W_i) P(W_i) dx$$

当决策面为  $x_B$  时为最小错误率决策。



(15) 用最大后验概率决策情况

$$P(\text{error}) = \int_x [1 - \max_i p(w_i | x)] \cdot P(x) dx$$





## 第章 贝叶斯决策、参数估计

① 贝叶斯决策: 1° 最小风险决策:  $\min R(\alpha_i | x) = \sum_{j=1}^c \lambda(\alpha_i | w_j) P(w_j | x)$ 2° 最小错误率 (最大后验概率)  $\max p(w_i | x)$  MAP② 离散特征变量:  $P(x | w_i) = P(x_1, x_2, \dots, x_d | w_i)$ ③ 独立=值特征  $P(x) = P(x_1, x_2, \dots, x_d) = \prod_{i=1}^d P(x_i)$ 

2类:  $P(x | w_1) = \prod_{i=1}^d p_i^{x_i} (1 - p_i)^{1 - x_i}$

$P(x | w_2) = \prod_{i=1}^d q_i^{x_i} (1 - q_i)^{1 - x_i}$

最大似然比  $\frac{P(x | w_1)}{P(x | w_2)} = \prod_{i=1}^d \left( \frac{p_i}{q_i} \right)^{x_i} \left( \frac{1 - p_i}{1 - q_i} \right)^{1 - x_i}$

④ 独立=值特征的决策面

$$g(x) = \log \frac{P(x | w_1) P(w_1)}{P(x | w_2) P(w_2)} = \sum_{i=1}^d \left[ x_i \ln \frac{p_i}{q_i} + (1 - x_i) \ln \frac{1 - p_i}{1 - q_i} \right] + \ln \frac{P(w_1)}{P(w_2)}$$

$$g(x) = \sum_{i=1}^d w_i x_i + w_0$$

$$w_i = \ln \frac{p_i (1 - q_i)}{q_i (1 - p_i)} \quad i = 1, \dots, d$$

$$w_0 = \sum_{i=1}^d \ln \frac{1 - p_i}{1 - q_i} + \ln \frac{P(w_1)}{P(w_2)}$$

例子1:  $P(w_1) = 0.5 \quad P(w_2) = 0.5$

$p_i = 0.8 \quad q_i = 0.5 \quad i = 1, 2, 3$

例子2:  $P(w_1) = 0.5 \quad P(w_2) = 0.5$

$p_1 = p_2 = 0.8 \quad p_3 = 0.5 \quad q_i = 0.5 \quad i = 1, 2, 3$

## ⑤ 贝叶斯分类

$$X = x_1, x_2, \dots, x_n \quad W = w(1), w(2), \dots, w(n)$$

$$P(w | x) = \frac{P(x | w) P(w)}{P(x)} = \frac{P(x | w) P(w)}{\sum_w P(x | w) P(w)} \quad W \text{ 类别数巨大, } P(x | w) \text{ 存储估计困难}$$

⑥ 贝叶斯融合:

$$P(w_i | e_1, \dots, e_k) = \frac{P(e_1, \dots, e_k | w_i) P(w_i)}{P(e_1, \dots, e_k)} \quad i = 1, \dots, c$$

$$P(e_1 = w_1, \dots, e_k = w_k | w_i) = \prod_{k=1}^k P(e_k = w_k | w_i)$$

## ⑦ 分类器设计

1° 给定分类器结构/函数形式, 从样本中估计参数

2° 统计生成模型: 密度估计

3° 统计判别模型: 判别函数

⑧统计模型的参数估计

- 1° ML 最大似然法: 假设参数为确定值, 最优估计: 似然度最大.
- 2° 贝叶斯估计: 假设参数为随机变量, 估计其分布.

⑨ ML 法:

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

最大似然:  $\max_{\theta} P(D|\theta) \leftrightarrow \nabla_{\theta} P(D|\theta) = 0$

可能有解析解 有可能需要迭代求解

⑩ Log-likelihood 求解 ML 问题

第一步变对数  $l(\theta) = \ln P(D|\theta) \quad l(\theta) = \sum_{k=1}^n \ln P(x_k|\theta)$

第二步求 ML (估计)  $\hat{\theta} = \arg \max_{\theta} l(\theta)$

$\nabla_{\theta} l = \sum_{k=1}^n \nabla_{\theta} \ln P(x_k|\theta) = 0$

$\frac{\partial l}{\partial \theta_j} = 0 \quad j = 1, \dots, p.$

高斯模型

当  $P(x_k|\mu) = \frac{1}{\sqrt{2\pi}} \frac{1}{|\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(x_k - \mu)^T \Sigma^{-1}(x_k - \mu)\right]$

$\nabla_{\theta} \ln P(x_k|\mu) = \Sigma^{-1}(x_k - \mu) = 0$

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

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

高斯估计  $\mu$  和  $\Sigma$ :

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

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

⑪ 高维密度贝叶斯估计  $\Rightarrow$  正则

$P(\mu|D) = \frac{P(D|\mu)P(\mu)}{\int P(D|\mu)P(\mu)d\mu} = \alpha \prod_{k=1}^n P(x_k|\mu) P(\mu)$   $\alpha$  为正则化因子.

⑫ 贝叶斯估计一般情况

1° 后验参数分布:  $P(\theta|D) = \frac{P(D|\theta)P(\theta)}{\int P(D|\theta)P(\theta)d\theta}$

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

2° 后验数据分布:  $P(x|D) = \int P(x|\theta) \cdot P(\theta|D) d\theta$



② 模型使用参数估计或最大后验估计

④ 最大似然估计与贝叶斯估计二者的区别和联系:

最大似然估计和贝叶斯估计最大区别在于估计的参数不同。最大似然估计要估计的参数  $\theta$  被当做是固定不变的。而贝叶斯估计则是有某种已知先验分布的随机变量。在贝叶斯估计中假设条件概率密度  $P(X|\theta)$  符合一定先验分布。参数  $\theta$  也符合一定先验分布。

同: 虽参数估计使得与条件概率密度相对简单, 但估计结果的准确性依赖于所假设的先验分布形式是否符合。查有和直实验数据所。

⑤ 增加特征有何好处和缺点?

答: 增加特征将提高模型判别性, 类别间有差异的特征有助于分类。

增加特征也可能导致分类性能变差, 因为模型估计误差。

⑥ 最大似然估计的计算复杂度

1. 参数估计复杂度  $\mu$  为  $O(dn)$   $d \rightarrow$  特征维度  
 $\Sigma$  为  $O(nd^2)$   $n \rightarrow$  样本数

2. 分类复杂度: 计算逆矩阵比较复杂, 一般为  $O(d^3)$

⑦ 过拟合: 特征维数高, 训练样本少导致模型参数估计不准确

解决方法: 特征降维, 参数共享/平滑。

⑧ 期望-最大法(EM)  $\rightarrow$  看作业题 (数据缺失情况下的参数估计)

EM算法保证数据的对数似然单调递增

⑨ 高斯混合的EM

$$P(X) = \sum_{k=1}^K \pi_k P(X|\theta_k) \quad \text{s.t.} \quad \sum_{k=1}^K \pi_k = 1$$

$$P(X|\theta_k) = N(X|\mu_k, \Sigma_k)$$

参数估计用 ML

$$\max LL = \log \prod_{i=1}^N P(X_i) = \sum_{i=1}^N \log \sum_{k=1}^K \pi_k P(X_i|\theta_k)$$
$$\nabla_{\pi_k} LL = 0 \quad \nabla_{\mu_k} LL = 0 \quad \nabla_{\Sigma_k} LL = 0$$

## ② 隐马尔可夫模型

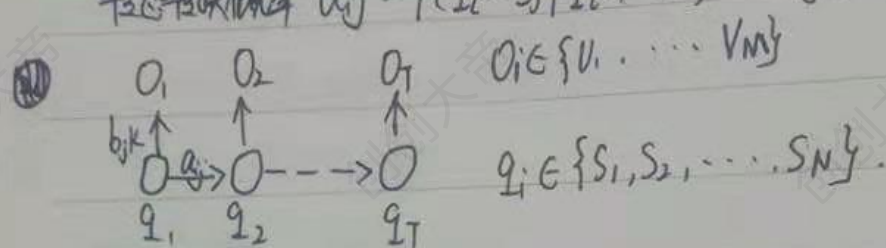
~~$P(O, q)$~~   $P(O|q)$

~~马尔可夫链~~  $P(q_1, q_2, \dots, q_T) = P(q_1)P(q_2|q_1)P(q_3|q_1, q_2) \dots P(q_T|q_1, \dots, q_{T-1})$

~~马尔可夫链~~  $P(q_t = s_j | q_{t-1} = s_i, q_{t+2} = s_k, \dots) = P(q_t = s_j | q_{t-1} = s_i)$

$P(q_1, q_2, \dots, q_T) = P(q_1) \cdot P(q_2|q_1)P(q_3|q_2) \dots P(q_T|q_{T-1})$

~~状态转移概率~~  $a_{ij} = P(q_t = s_j | q_{t-1} = s_i) \quad \sum_{j=1}^N a_{ij} = 1$



~~观察值可能分布~~  $b(k) = P(v_k \text{ at } t | q_t = s_j)$

④ 隐马尔可夫的计算  $\rightarrow$  看作业题

$$P(O|\lambda) = \sum_{\text{all } Q} P(O|Q, \lambda) P(Q|\lambda)$$

$$= \sum_{q_1, \dots, q_T} \pi_{q_1} b_{q_1}(O_1) a_{q_1 q_2} b_{q_2}(O_2) \dots a_{q_{T-1} q_T} b_{q_T}(O_T)$$

~~条件独立~~  $P(O|Q, \lambda) = \prod_{t=1}^T P(O_t | q_t, \lambda)$

$$P(Q|\lambda) = \pi_{q_1} a_{q_1 q_2} a_{q_2 q_3} \dots a_{q_{T-1} q_T}$$

$$= b_{q_1}(O_1) b_{q_2}(O_2) \dots b_{q_T}(O_T)$$

$$\text{复杂度 } \boxed{2TN^T}$$

④ 隐马尔可夫 编码和训练 (见书)



## 非参数方法

① 核密度估计  $p(x) = \frac{k}{nV}$ 

② 非参数核密度估计:

1° Parzen window: 固定局部区域体积,  $K$  变化2° KNN: 固定局部样本数,  $V$  变化.

③ Parzen 窗

$$p(u) = \begin{cases} 1 & |u| \leq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases} \quad \text{满足条件 } p(x) \geq 0 \quad \int p(u) du = 1$$

以  $x$  为中心, 体积为  $V_h = h_n^d$  的局部区域内样本数

$$k_n = \sum_{i=1}^n p\left(\frac{x - x_i}{h_n}\right)$$

核密度估计  $k_n / nV_h$ 大  $h$ : 低可变性, 未拟合小  $h$ : 高可变性, 过拟合窗宽  $h_n$  选择经验: 一般原则:  $n$  越大或密度越大,  $h_n$  越小

$$V_h = V / \sqrt{n}$$

④ K近邻估计  $\rightarrow$  老师讲会考固定局部区域样本数  $k$ , 体积  $V$  变化 $k$  窗越大更平滑

K-NN 族规则里没有核密度, 但要注意, 该规则是非参数核密度估计和贝叶斯决策规则的结合

最近邻规则的错误差率K近邻的快速计算  $O(dn)$ 

近邻搜索三种策略:

1° 暴力距离

2° 预结构化

3° 编辑

⑤ 度量学习: 距离度量的参数未被优化.

⑥ 对于图像变换的转换(平移, 缩放, 弯曲, 旋转), 距离将会有很大变化, 因此引出了切线距离

Tangent Distance 较好地解决了经过图像变换后距离度量问题, 其通过梯度下降求解切线

空间中优化求得被分类向量和原始图像及其经过变换后图像的空间中最近的点。

(7) 系列展开估计

$$\varphi\left(\frac{x-x_i}{h_n}\right) = \sum_{j=1}^m d_j \psi_j(x) \overset{\text{测试样本}}{\underset{\text{训练样本}}{X_j(x_i)}}$$

$$\sum_{i=1}^N \varphi\left(\frac{x-x_i}{h_n}\right) = \sum_{j=1}^m d_j \psi_j(x) \sum_{i=1}^N X_j(x_i)$$

$$p_n(x) = \sum_{j=1}^m b_j \psi_j(x) \quad b_j = \frac{a_j}{n h_n} \sum_{i=1}^N X_j(x_i) \quad b_j \text{ 随样本 } p_n(x) \text{ 随 } m \text{ 次计算}$$



## 线性判别函数

① 不需要有概率密度函数的确切函数形式，因此属于无参数估计方法。

② 模式分类的途径：

1° 估计类条件概率密度函数，利用贝叶斯公式求出后验概率，然后决策。概率密度参数估计和非参数估计

2° 直接估计后验概率，不需要估计类条件概率密度  $k$ -近邻

3° 直接计算判别函数，直接判别可用于分类的判别函数

③ 利用样本设计分类器基理想：

1° 给定一个判别函数，且已知该函数的参数形式

2° 采用样本来训练判别函数的参数

3° 对于新样本，采用判别函数对其进行分类，并按照一些准则来完成分类。

④ 多类情形  $\rightarrow$  作业题，估计误差

1° one vs all 存在不确信区域

2° one vs one 存在不确信区域

对  $j \neq i$  如果  $g_i(x) > g_j(x)$   $x$  应被分到  $W_i$  类，否则不决策  
 $x \in W_i: g_i(x) = \max_{j=1,2,\dots,C} g_j(x)$   
 线性机器将样本空间分为  $C$  个互不重叠的决策区域  $R_1, \dots, R_C$

⑤ 广义线性判别函数

数据点  $x$   $\xrightarrow{\text{非线性映射}} \text{数据点 } y$

在低维特征空间应用中应用线性判别函数方法。

$$g(x) = \sum_{i=1}^d a_i [y_i(x)]$$

变换函数

$$a = [a_1, a_2, \dots, a_d]^T \quad y = [y_1, y_2, \dots, y_d]^T$$

$$g(x) = a^T y$$

$a$  为广义权重向量， $g(x)$  对  $x$  非线性，对  $y$  线性。当特征空间维数足够高时， $g(x)$  可以逼近任意非线性判别函数。

⑥ 对线性判别函数采用齐次增广表示：线性齐次空间中增加一个维数，仍可保持欧式距离不变，分类结果与原来的决策面相同，但决策面经过坐标原点。  $\gamma = \frac{a^T y}{\|a\|}$

⑦ 规范化处理：规范化增广样本，解向量 ( $a^T y_i > 0$  的权重向量  $a$ )

⑧ 已知准则函数

$$a^T y_i > 0 \quad i=1, 2, \dots, n$$

在线性可分情形下， $a$  有无穷个，所以需引入准则：

$$J_p(a) = \sum_{i=1}^n y_i (-a^T y_i) \quad \text{称为全错分率}$$

$$\min_a J_p(a)$$

$$\frac{\partial J(a)}{\partial a} = - \sum_{y \in Y} y$$

$$a_{k+1} = a_k + \eta_k \sum_{y \in Y} y$$

⑨ 感知准则存在的问题

1° 每次错为样本点集合可能会改变

2° 倾向于保留离坐标原点远的

⑩ 固定增量单样本修正方法, 可变增量单样本修正方法

⑪ 感知准则收敛性定理证明

⑫ 学习准则: 线性准则, 平方准则, 松弛准则

$$J(a) = \sum_{y \in Y} (-a^T y)$$

$$J_2(a) = \sum_{y \in Y} (a^T y)^2$$

$$J_r(a) = \frac{1}{2} \sum_{y \in Y} \frac{(a^T y - b)^2}{\|y\|^2} \quad \text{归一化 + margin}$$

解决某些影响

$y$  为  $a^T y = b$  类

松弛准则特点: ① 线性准则为分段线性, 梯度不连续

② 平方准则梯度连续, 但目标函数过平滑, 收敛速度慢, 同时, 目标函数达到最优解时

③ 松弛准则避免上述缺点

$$\frac{\partial J_r(a)}{\partial a} = \sum_{y \in Y} \frac{a^T y - b}{\|y\|^2} y$$

$$a_{k+1} = a_k - \eta_k \sum_{y \in Y} \frac{a^T y - b}{\|y\|^2} y$$

$$\text{对于单样本更新: } a_{k+1} = a_k - \eta_k \frac{a^T y_k - b}{\|y_k\|^2} y_k$$

$$= a_k - \eta_k \frac{a^T y_k - b}{\|y_k\|} \cdot \frac{y_k}{\|y_k\|}$$

点到超平面距离

单位向量方向

$0 < \eta_k < 1$	靠近超平面
$\eta_k = 1$	到达超平面
$\eta_k > 1$	超过超平面

⑬ 松弛准则收敛性证明

⑭ 最小平方误差 (MSE) 准则函数

$$a^T y_i = b_i, \quad b_i > 0$$

$b_i$  为任意给定的正常数

$$e = Ya - b$$

$Y$  为矩阵  $a$  为向量  $b$  为向量

$$J_s(a) = \|e\|^2 = \|Ya - b\|^2 = \sum_{i=1}^n (a^T y_i - b_i)^2$$

$$\frac{\partial J_s(a)}{\partial a} = \sum_{i=1}^n 2(a^T y_i - b_i) y_i = 2Y^T(Ya - b)$$

令偏导为 0  $\Rightarrow a = Y^+ b$



对于  $Y'$  计算复杂度高, 可以采用梯度下降法,  $\alpha_{k+1} = \alpha_k + \eta_k Y^T (b - Y \alpha_k)$

其收敛条件是  $Y^T (b - Y \alpha) = 0$ .

$\alpha_{k+1} = \alpha_k + \eta_k (b_k - (\alpha_k)^T y^k) y^k$ . 考虑单个样本对误差的贡献.

⑮ Widrow-Hoff 方法.  $\eta_k$  随着  $k$  增加而逐渐减小.  $\eta_k = \eta_0 / k$ .

⑯ ~~Widrow-Hoff~~ 相对于感知器准则, 最小平方准则方法可能并不收敛于可分超平面 (即使该平面是存在的).

MSE 方法本质上是求样本至超平面的距离平方和

⑰ Ho-kashyap

$J_s(a, b) = \|Y a - b\|^2$   $a, b$  均未知, 进行优化.  $b$  可理解为 margin. 解决 ~~Widrow-Hoff~~ 方法得到最优解不一定在可分超平面上.

若  $e_k = Y \alpha_k - b_k$  全为 0, 此时  $b_k$  将不再更新. 如果  $e_k$  有部分元素小于 0, 则可证明该问题不是线性可分的.

⑱ 多类线性判别函数

- 1° MSE 准则扩展
- 2° 感知器准则扩展方法——修正修正法
- 3° Kelsner 构造

# 人工神经网络第1讲

① 神经网络应是非线性，否则多层网络计算能力并不比单层网络强

② 前馈网络与反馈网络比较

1° 前馈型网络“不存记忆”，结束输出为当前输入值的加权和再经激活

2° 反馈型网络要将以前的输出值循环返回到输入

3° 反馈型类似于“人类的短期记忆”，网络的输出状态都依赖于前输出

③ 无监督学习本质是抽取样本所隐含的统计特性

④ 梯度下降法

$$\Delta W_{ij} = -\eta \frac{\partial E}{\partial W_{ij}}$$

⑤ 线性单元训练：S规则

$$E(W) = \frac{1}{2} \sum (t_j^k - z_j^k)^2 = \frac{1}{2} \sum (t_j^k - \sum W_{ij} x_i^k)^2$$

$$\Delta W_{ij} = -\eta \frac{\partial E}{\partial W_{ij}} = \eta \sum (t_j^k - z_j^k) x_i^k = \eta \delta_j^k x_i^k \quad \delta_j = t_j^k - z_j^k$$

⑥ 考虑转移函数，转移函数是非线性函数且处处可微

单个样本的贡献： $\Delta W_{ij} = \eta \delta_j^k x_i^k$   $\delta_j^k = \frac{-\partial E}{\partial net_j^k} = f'(net_j^k) [t_j^k - z_j^k]$

解释：左边的转移函数点的输出乘以右边的指向结点，“经导数放大后的误差”

⑦ Sigmoid函数  $f(s) = \frac{1}{1+e^{-s}}$

$$f'(s) = \frac{e^{-s}}{(1+e^{-s})^2} = \frac{1}{1+e^{-s}} \left(1 - \frac{1}{1+e^{-s}}\right) = y(1-y)$$

$$\delta_j^k = y_j^k (1 - y_j^k) [t_j^k - y_j^k]$$

\* ⑧ BP算法基本原理：利用输出层的误差来估计输出层为前一层的误差，再用这个误差估计前一层的误差。如此一层一层地倒推回去，从而获得所有其它各层的误差估计

⑨ BP算法权重重新推导（看作业题）



人工神经网络第2讲

- ① 对分类问题, 采用 one-hot 编码
- ② 交叉熵准则:  $E_{ce}(W) = - \sum_{j=1}^K t_j^k \ln(z_j^k)$
- ③ 激活函数: 1° 非线性  
2° 有界连续可导  
3° 最好单调.
- ④ 隐含层数设定: 对分类问题, 隐含层的数决定了网络的表达能力, 决定了决策面的复杂度
- ⑤ 结点数: 太少, 不能建立复杂判别界面  
太多, 容易过拟合, 失去推广能力  
一般设置较大, 后续有无对结点数可能
  - 1° 压缩神经网络
  - 2° 稀疏连接神经网络
  - 3° dropout 技术.
- ⑥ 初始权重: 初始权重  $\Delta W_{ij} \neq 0$ , 否则不会学习
- ⑦ 正则化技术: 防止网络出现 overfitting 的一种有效方法
- ⑧ 学习率.
- ⑨ 附加冲量项.
- ⑩ 训练停止准则: 没有固定准则, 过度训练会 overfitting, 训练不够可能会出现识别精度不足
- ⑪ BP 算法中存在的问题:
  - 1° 完全难以训练: 网络的麻痹现象, 梯度消失, 局部最小
  - 2° 训练时间过长: 尤其对复杂问题需要很长时间训练, 可能选取了不恰当
- ⑫ 网络的麻痹现象: 由于在计算权重修正量时, 误差  $\delta_i$  正比于  $f'(net_i)$ , 当  $f'(net_i) \rightarrow 0$  时,  $\delta_i \rightarrow 0$  从而  $\Delta W_{ij} \rightarrow 0$ , 这相当于调节过程停顿.
- ⑬ 梯度消失: 由于  $\delta_i^k = f'(net_i^k) \delta_i^{k-1}$ , 由于  $f'(net_i^k)$  通常小于 1, 所以在多层神经网络向神经网络中, 越靠近输入层越容易出现梯度消失

DATE

MTWTFSS

⑭ 径向基函数网络：其可以处理连续的非线性数据并近似。收敛速度比通常多层神经网络快。

⑮ 径向基网络结构简化，采用聚类技术对数据进行聚类，每个隐含层结点代表一个聚类中心。  
经过聚类处理，还可防止 overfitting，增强网络的泛化能力，提高精度。

⑯ Hopfield网络：从初始状态按能量减小方向进行演化，直到达到稳定状态，稳定状态即为网络的解。

⑰ BM，一种随机 Hopfield 网络。BM 部分神经元外部相连，另一部分不与外部相连。

缺：网络结构复杂，训练代价大，局部极小。

⑱ RBM 层内结点不相连，信息双向流动。

⑲ 自组织映射



## 人工神经网络第3讲

①卷积神经网络降低网络权重数量：局部连接  $\rightarrow$  利用图像空间相关性

权重共享  $\rightarrow$  采用滤波器，每个滤波器可看成是学习一种特征

②卷积神经网络网络训练采用反向传播算法

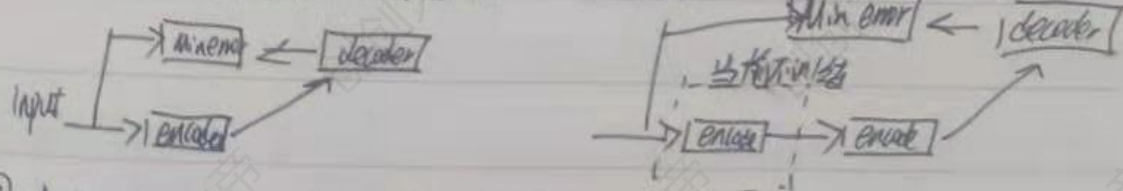
③卷积神经网络核心思想：1° 使用卷积后的特征是因为图像具有一种“静态性”的属性，这也就意味着在一个图像区域有用的特征极有可能在另一区域同样适用

2° 局部感受野权重共享以及时间或空间下采样这三种结构化思想结合起来获得了某种程度的位移、尺度、形变不变性

④自编码器：给定一个神经网络，假设其输出与输入是相同的，通过训练该网络，可以得到网络中的权重。

⑤自编码器是一种尽可能重构输入信号的神经网络

⑥自编码器使用逐层训练思想，每次都是一个层前向神经网络



⑦ Autoencoder 可用作特征提取器使用，只应用编码器，最后一层编码器的输出作为该样本的新特征。

⑧ Autoencoder 实质上是一种特征学习，层级越高，特征的语义性、抽象性、结构性越明显。

对于分类器，可以用 Autoencoder 学习，学习得到权重为初始权重，采用带有标签数据进行再次学习，进行微调。

⑨

## 特征提取与特征选择

④ SIFT 尺度不变特征转换: 1° 检测尺度空间内极值点

2° 检测角点并建立定位

3° 特征方向估计

4° 类与像素描述计算

② 特征生成:

③ 维数削减作用: 1° 低维数据可视化

2° 加速收敛性

3° 加速学习算法

4° 减少存储

④ PCA  $\rightarrow$  有作业题 计算流程

⑤ CCA 典型相关性分析: 探索两个变量之间的关联关系, 这两个变量来自一个相同体系

1° 数据降维: 用少量线性组合来解释两组变量之间的关联作用

2° 数据解释: 寻找特征值, 这些特征值在解释两个变量之间相互作用协方差

⑥ LDA 线性判别分析: 不同于PCA方差最大化理论, LDA算法的思想是将数据投影到低维空间后使得同类数据尽可能紧凑, 不同类的数据尽可能分散

⑦ LDA = 白化 + PCA

⑧ LDA的约束或假设: 1° 每个类是单模态高斯分布

2° 每一类的协方差矩阵都相同

3° 降维后维数不超过  $C-1$

4° 类别分离问题

⑨ 独立成分分析 ICA. ~~与PCA不同~~ 与PCA不同, ICA追求的是输出变量相互独立而非仅不相关, 因此需要利用数据分布的高阶统计信息而非仅二阶信息  
很多应用中, ICA提取的特征优于PCA



⑩ ICA会恢复幅值,但对应关系消失

⑪ 非线性打度,流形,保持邻居关系

解决方法: Kernel PCA

使用流形基于方法 (LLE, Isomap, ...)

⑫ Kernel PCA: 对于输入空间中的数据集  $X$ , 我们用一个非线性映射把  $X$  中的所有样本映射到一个高维甚至是无穷维的空间, ~~然后~~ 然后在这个高维空间进行 PCA 降维

⑬ LLE: 显然利用"局部线性"的假设,流形上的局部区域具有欧式空间的性质,那么在 LLE 中就假设某点  $X_i$  坐标可以由它周围的一些点的坐标线性组合求出。

保持局部邻域几何结构

权重对样本集的几何变换具有不变性。

主要步骤: 1° 寻找每个样本点的  $k$  个近邻点

2° 由每个样本点的近邻点计算出该样本点的局部重建权重矩阵

3° 由该样本点的局部重建权重矩阵和其近邻点计算出该样本点的输出值

⑭ ~~isomap~~ MDS 降维算法: 将高维坐标中的点投影到低维空间,保持点彼此之间的相似距离尽可能不变。

算法流程: 1° 计算原始空间中数据点矩阵

2° 计算内积矩阵  $B$

3° 对矩阵  $B$  进行特征值分解, 获得特征值矩阵  $\Lambda$  和特征向量矩阵  $V$

4° 取特征值矩阵最大前  $L$  项及其对应的特征向量  $Z = V_L A_L^{\frac{1}{2}}$

⑮ isomap 基于 MDP 算法, 不同之处在于 isomap 用图中两点的最短路径距离替代了 MDS 中欧氏空间的距离, 这样能更好的拟合流形数据。

算法流程 (1) 设置每个点最近邻点  $k$ , 构建连通图和邻接矩阵

(2) 通过图的最短路径构建原始空间中的距离矩阵

(3) 计算内积矩阵

14 对矩阵  $B$  进行特征值分解, 获得特征值矩阵  $\Lambda$  和特征向量矩阵  $V$

15 取特征值矩阵最大的前  $Z$  项及其对应特征向量  $Z = V_z \Lambda_z^{\frac{1}{2}}$

16 LPP 局部保持投影: 非线性降维方法的线性化。相同具有某种非线性关系的样本, 在降维后仍然保持这种非线性关系

17 特征选择: 降维 VS 特征选择

1° 所有原始特征被使用

2° 转化特征是原始特征线性组合

原始特征的若干子集被选择

18 特征选择不能一个个特征看, 不能忽略特征相关性

过滤法: 按照相关性或者相关性对各个特征进行评分, 设定阈值或者特选择用值个数, 选择特征

Wrapper: 包装法, 根据目标函数 (通常是通过验证集效果评分), 迭代选择若干特征, 或者排除若干特征

Embedded: 嵌入法, 用某些算法和模型进行训练, 得到各个特征的权值系数, 根据系数大小选择特征



## 模型选择

① 机器学习: 泛化 过拟合 维数灾难, 理论上保证 相关并不意味着因果

随着样本维数增加, 正确泛化难度从指数级增加.

② 模型集成: 若干单个模型集成.  $\begin{cases} \text{单一学习器错误率低于 } 0.5 \\ \text{每个分类器应当各不相同.} \end{cases}$

③ 层叠泛化: 多层结构, 第一层的输出作为第二层的输入

④ Bagging: 训练一组基分类器, 每个基分类器通过有放回地随机训练样本集来训练, 后通过投票进行统计.

⑤ 随机子空间: 对于每一个分类器, 从  $D$  个特征中选择  $d$  个子特征来构建一个训练集同时学习  $d$  个分类器

⑥ Adaboost: 从弱学习算法出发, 反复学习得到一系列弱分类器, 然后加权组合这些子分类器, 构成一个强分类器

1° 提高那些被前一轮弱分类器分类错误的样本权重, 降低已经被正确分类样本权重

2° 加大分类错误率较小的分类器的权重

## 数据聚类

① 聚类方法分类: 基于距离的聚类方法, 基于密度的聚类方法, 基于连通性的聚类方法

$$\textcircled{2} \nabla_{\theta} f_{\theta}(x_k) = P(w_i | x_k, \theta) \nabla_{\theta} \ln L(P(x_k | w_i, \theta; 1))$$

单样本  $x_k$  对 "似然函数关于  $\theta$  的梯度" 之贡献等于 " $x_k$  属于第  $i$  成分的概率" 乘以 " $x_k$  第  $i$  成分密度  $P(x | w_i, \theta)$  的对数关于  $\theta$  的梯度".

## ③ 基于混合密度的聚类

④ 谱聚类: 其本质是将聚类问题转化为图点划分的最优问题

拉普拉斯矩阵: 度矩阵减去邻接矩阵

子图相似度: 连接两子图所有边的权重之和

最小二方切割: 切开后两子图之间的相似性最小

归一化最小二方切割: 使用图的总权重体积来对切割进行归一化

⑤ 如果图  $G$  具有  $k$  个连通子图, 若每个连通子图为一个聚类, 那么使用其拉普拉斯矩阵的最小特征值的特征向量可以分离这些子图

⑥ 考虑拉普拉斯矩阵最小的特征值对应的特征向量, 并由这些特征向量组成新的特征正空间

⑦ 谱聚类法核心步骤: 1° 计算点之间相似度, 构建亲和度矩阵

2° 构建拉普拉斯矩阵

3° 求解拉普拉斯矩阵最小的特征值对应的特征向量

4° 由这些特征向量构成点的新特征, 使用  $k$ -means 等聚类方法完成最后的聚类

⑧ 未标准化的谱聚类: 计算  $L$  的  $k$  个最小特征值对应的特征向量, 并将原始数据点  $X$  转换为特征空间的数据点

⑨



## 机器学习与核方法

① VC维

② 硬边界-SVM

$$\arg\max_{w, b} \arg\min_{x_i \in D} \frac{|b + x_i \cdot w|}{\sqrt{\sum_{i=1}^d w_i^2}}$$

$$s.t. \forall x_i \in D: y_i(x_i \cdot w + b) \geq 0$$

$$\Downarrow$$

$$\arg\min_{w, b} \sum_{i=1}^d w_i^2$$

$$s.t. \forall x_i \in D: y_i(x_i \cdot w + b) \geq 1$$

③ soft-SVM

$$\{\vec{w}^*, b^*\} = \arg\min_{w, b} \sum_{i=1}^d w_i^2 + C \sum_{j=1}^N \xi_j$$

$$y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$$

$$y_N(\vec{w} \cdot \vec{x}_N + b) \geq 1 - \xi_N, \quad \xi_N \geq 0$$

3.1 hinge Loss

$$\arg\min_{f \in H} \frac{1}{n} \sum_{i=1}^n (1 - y_i f(x_i)) + \lambda \|f\|_H^2$$

$$\text{硬} \quad \max W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \quad s.t. \alpha_i \geq 0 \quad \sum_{i=1}^n \alpha_i y_i = 0$$

$$\text{软} \quad \max W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \quad s.t. C \geq \alpha_i \geq 0 \quad \sum_{i=1}^n \alpha_i y_i = 0$$

④

## 决策树法

① 熵  $H(X) = - \sum_{x \in X} P(X) \log_2 P(X)$

② 条件熵  $H(Y|X) = \sum_i P_i H(Y|X=x_i)$

信息增益: 得知特征X的信息而使类Y的信息不确定性减少的程度

$$g(D|A) = H(D) - H(D|A)$$

## ③ 信息增益算法 会考

输入: 训练数据集D和特征A:

输出: 特征A对训练数据集D的信息增益  $g(D, A)$

1) 计算数据集D的经验熵  $H(D)$

$$H(D) = - \sum_{k=1}^K \frac{|C_k|}{|D|} \log_2 \frac{|C_k|}{|D|}$$

2) 计算特征A对数据集D的经验条件熵  $H(D|A)$

$$H(D|A) = \sum_{i=1}^n \frac{|D_i|}{|D|} H(D_i) = - \sum_{i=1}^n \frac{|D_i|}{|D|} \sum_{k=1}^K \frac{|D_{i,k}|}{|D_i|} \log_2 \frac{|D_{i,k}|}{|D_i|}$$

3) 计算信息增益  $g(D, A) = H(D) - H(D|A)$

④ ID3 采用信息增益作为度量, 其偏向于具有大量值的属性. 在训练集中某个属性所取的不同值中, 出现越多的值也作为分裂属性

⑤ C4.5 采用信息增益率

$$\text{SplitInfo}(A) = - \sum_{j=1}^V \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)$$

$$\text{GainRatio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}(A)}$$

⑥ 分类与回归树 CART 使用二叉树. 基尼指数

$$\text{Gini}(P) = \sum_{k=1}^K P_k (1 - P_k) = 1 - \sum_{k=1}^K P_k^2$$

$$\text{Gini}(D, A) = \frac{|D_1|}{|D|} \text{Gini}(D_1) + \frac{|D_2|}{|D|} \text{Gini}(D_2)$$

⑦ 决策树过拟合. 剪枝

⑧ 随机森林: RF 使用了两次随机抽取. 1° 对训练样本的随机抽取 (bagging)  
2° 对变量的随机抽取

(A: 并对特征值随机抽取)



① 随机森林的两次随机过程 1°防止过拟合, 随机选择特征  
2°对特征空间扰动.

② 影响RF性能因素 1°分类强度, 每颗树的分类强度越大, RF分类性能越好  
2°树之间相关性 越小越好