

聚类分析

相似性的测度

- 欧氏距离

设 $\mathbf{X}_1, \mathbf{X}_2$ 为两个 n 维模式样本, $\mathbf{X}_1 = [x_{11}, x_{12}, \dots, x_{1n}]^T, \mathbf{X}_2 = [x_{21}, x_{22}, \dots, x_{2n}]^T$

$$D(\mathbf{X}_1, \mathbf{X}_2) = \|\mathbf{X}_1 - \mathbf{X}_2\| = \sqrt{(\mathbf{X}_1 - \mathbf{X}_2)^T (\mathbf{X}_1 - \mathbf{X}_2)} = \sqrt{(x_{11} - x_{21})^2 + \dots + (x_{1n} - x_{2n})^2}$$

使特征数据标准化, 使其与变量的单位无关

- 马氏距离

$$D(\mathbf{X}) = \sqrt{(\mathbf{X} - \mathbf{M})^T \mathbf{C}^{-1} (\mathbf{X} - \mathbf{M})}$$

$$\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \text{ 为模式向量, } \mathbf{M} = \begin{bmatrix} m_1 \\ \vdots \\ m_n \end{bmatrix} \text{ 为均值向量, } \mathbf{C} \text{ 为协方差矩阵}$$

$$\mathbf{C} = E \{ (\mathbf{X} - \mathbf{M})(\mathbf{X} - \mathbf{M})^T \}$$

表示的概念是各分量上模式样本到均值的距离, 也就是在各维上模式的分散情况

$$d(\mathbf{X}, \mathbf{M}) = \sqrt{\left(\frac{x_1 - m_1}{\sigma_1}\right)^2 + \left(\frac{x_2 - m_2}{\sigma_2}\right)^2 + \dots + \left(\frac{x_n - m_n}{\sigma_n}\right)^2}$$

$$d(\mathbf{X}, \mathbf{Y}) = \sqrt{(\mathbf{X} - \mathbf{Y})^T \mathbf{C}^{-1} (\mathbf{X} - \mathbf{Y})} = \sqrt{\left(\frac{x_1 - y_1}{\sigma_1}\right)^2 + \left(\frac{x_2 - y_2}{\sigma_2}\right)^2 + \dots + \left(\frac{x_n - y_n}{\sigma_n}\right)^2}$$

- 明氏距离

$$D_m(\mathbf{X}_i, \mathbf{X}_j) = \left[\sum_{k=1}^n |x_{ik} - x_{jk}|^m \right]^{\frac{1}{m}}$$

$$m = 1 \text{ 时, } D_1(\mathbf{X}_i, \mathbf{X}_j) = \sum_{k=1}^n |x_{ik} - x_{jk}| \text{ 称为街区距离}$$

$m = 2$ 时, 为欧氏距离

- 汉明距离

设 $\mathbf{X}_i, \mathbf{X}_j$ 为 n 为二值 $(1, -1)$ 模式样本向量

$$D_h(\mathbf{X}_i, \mathbf{X}_j) = \frac{1}{2} \left(n - \sum_{k=1}^n x_{ik} \cdot x_{jk} \right)$$

- 角度相似性函数

$$S(\mathbf{X}_i, \mathbf{X}_j) = \frac{\mathbf{X}_i^T \mathbf{X}_j}{\|\mathbf{X}_i\| \cdot \|\mathbf{X}_j\|}$$

为模式向量 $\mathbf{X}_i, \mathbf{X}_j$ 之间夹角的余弦

聚类准则

- 阈值准则

- 函数准则

$$\text{误差平方和 } J = \sum_{j=1}^c \sum_{\mathbf{X} \in S_j} \|\mathbf{X} - \mathbf{M}_j\|^2$$

- $\{\mathbf{X}\}$ 为模式样本集
- $\{S_j, j = 1, 2, \dots, c\}$ 为模式类别
- c 为聚类类别数目
- $\mathbf{M}_j = \frac{1}{N_j} \sum_{\mathbf{X} \in S_j} \mathbf{X}$ 为属于 S_j 集的样本的均值向量
- N_j 为 S_j 中样本的数目

J 代表了分属于 c 个聚类类别的全部模式样本与其相应类别模式均值之间的误差平方和
适用于样本类数给定，各类样本密集且数目相差不多，而不同类间的样本又明显分开的情况

基于距离阈值的聚类算法

• 近邻聚类法

有 N 个待分类的模式 $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N\}$ ，要求按距离阈值 T 分类到以 $\mathbf{Z}_1, \mathbf{Z}_2, \dots$ 为聚类中心的模式类中

- 任取样本 \mathbf{X}_i 作为第一个聚类中心的初始值，如 $\mathbf{Z}_1 = \mathbf{X}_1$
- 计算样本 \mathbf{X}_2 到 \mathbf{Z}_1 的欧氏距离 $D_{21} = \|\mathbf{X}_2 - \mathbf{Z}_1\|$ ，若 $D_{21} > T$ 则定义一新的聚类中心 $\mathbf{Z}_2 = \mathbf{X}_2$ ，否则 $\mathbf{X}_2 \in$ 最近的聚类中心

很大程度上依赖于距离阈值 T 的大小、第一个聚类中心的位置选择、待分类模式样本的排列次序、以及样本分布的几何性质

• 最大最小距离算法

- 任取样本 \mathbf{X}_i 作为第一个聚类中心的初始值，如 $\mathbf{Z}_1 = \mathbf{X}_1$
- 选取离 \mathbf{Z}_1 最远的样本作为第二聚类中心 \mathbf{Z}_2
- 逐个计算各模式样本与已确定的所有聚类中心之间的距离，并选出其中的最小距离 $\min_k (\|\mathbf{X}_i - \mathbf{Z}_k\|)$
- 当最大值 $\max_i \min_k (\|\mathbf{X}_i - \mathbf{Z}_k\|)$ 达到 $\|\mathbf{Z}_1 - \mathbf{Z}_2\|$ 的一定分数比值（阈值 T ）以上，则选取新的聚类中心

层次聚类法

• 算法描述

- N 个初始模式样本自成一类 $G_1(0), G_2(0), \dots, G_N(0)$
计算各类之间的距离，得 $N \times N$ 维矩阵 $\mathbf{D}(0)$
- 找出 $\mathbf{D}(n)$ 中的最小元素，将对应的两类合并为一类，建立新的分类 $G_1(n+1), G_2(n+1), \dots$
- 计算合并后新类别的距离 $\mathbf{D}(n+1)$
- 当 $\mathbf{D}(n)$ 的最小分类超过阈值 T 或全部样本聚成一类时，算法停止

• 类间距离计算

◦ 最短距离法

设 H, K 为两个聚类，距离 $D_{HK} = \min \{D(\mathbf{X}_H, \mathbf{X}_K)\}, \mathbf{X}_H \in H, \mathbf{X}_K \in K$

若 K 由 I, J 合并而成，则 $D_{HK} = \min \{D_{HI}, D_{HJ}\}$

◦ 最长距离法

$D_{HK} = \max \{D(\mathbf{X}_H, \mathbf{X}_K)\}, \mathbf{X}_H \in H, \mathbf{X}_K \in K$

$D_{HK} = \max \{D_{HI}, D_{HJ}\}$

- 重心法
- 类平均距离法

动态聚类法

- **K-means**

设待分类的模式特征矢量集为 $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N\}$, 类的数目 K 事先取定。 S_j 为第 j 个聚类集, 聚类中心为 Z_j , 包含的样本数为 N_j

- 任选 K 个模式特征向量 $Z_1(1), Z_2(1), \dots, Z_K(1)$ 作为初始聚类中心
- 若 $D_j(k) = \min \{\|\mathbf{X} - \mathbf{Z}_i(k)\|\}, i = 1, 2, \dots, K$, 则 $\mathbf{X} \in S_j(k)$
- $Z_j(k+1) = \frac{1}{N_j} \sum_{\mathbf{X} \in S_j(k)} \mathbf{X}, j = 1, 2, \dots, K$

$$J_j = \sum_{\mathbf{X} \in S_j(k)} \|\mathbf{X} - \mathbf{Z}_j(k+1)\|^2, j = 1, 2, \dots, K$$

- 若 $Z_j(k+1) = Z_j(k), j = 1, 2, \dots, K$, 则结束
否则 $K = K + 1$

- **ISODATA**

- 确定控制参数、设置代表点

符号	内容
K	聚类期望数
θ_N	一个聚类的最少样本数
θ_C	类间距离控制参数
θ_S	标准偏差控制参数
L	每次迭代允许合并的最大聚类对数
I	允许迭代的次数
N_C	初始设定聚类数, 聚类中心为 $Z_i, i = 1, 2, \dots, N_C$

- 分类
若 $D_j = \min \{\|\mathbf{X} - \mathbf{Z}_i\|\}, i = 1, 2, \dots, N_C$, 则 $\mathbf{X} \in S_j$
- 撤销类内样本数过小类别
若有样本 $S_j, N_j < \theta_N$, 则舍去 S_j , 令 $N_C = N_C - 1$, 样本分配至其他类别

- 更新均值向量

$$\mathbf{Z}_j = \frac{1}{N_j} \sum_{\mathbf{X} \in S_j} \mathbf{X}, j = 1, 2, \dots, N_C$$

- 计算类内平均距离

$$\tilde{D}_j = \frac{1}{N_j} \sum_{\mathbf{X} \in S_j} \|\mathbf{X} - \mathbf{Z}_j\|, j = 1, 2, \dots, N_C$$

- 计算全部样本集到响应均值的平均距离

$$\tilde{D} = \frac{1}{N} \sum_{j=1}^{N_C} N_j \tilde{D}_j$$

- 入口选择, 判断分裂、合并

若迭代次数为 I 则合并

若 $N_C < \frac{K}{2}$, 则分裂

若迭代了偶数次, 或者 $N_C \geq 2K$, 则合并

其余分裂

- 分裂

- 求各聚类标准差

$$\sigma_j = [\sigma_{j1} \quad \sigma_{j2} \quad \cdots \quad \sigma_{jn}]^T$$

$$\sigma_{ji} = \sqrt{\frac{1}{N_j} \sum_{\mathbf{x}_k \in S_j} (x_{ki} - z_{ji})^2}$$

- 求 $\sigma_{j\max}, j = 1, 2, \dots, C$

- 执行分裂

若 $\sigma_{j\max} > \theta_S$, 且有 $\tilde{D}_j > \tilde{D}, N_j > 2\theta_N + 1$ 或者 $N_C \leq \frac{K}{2}$, 则把 S_j 分裂为两个聚类 $\mathbf{Z}_j^+, \mathbf{Z}_j^-$, $N_C = N_C + 1$

给定 $k, 0 < k < 1$, 令 $r_j = k\sigma_{j\max}$, $\mathbf{Z}_j^+ = \mathbf{Z}_j + r_j, \mathbf{Z}_j^- = \mathbf{Z}_j - r_j$

- 合并

- 计算聚类中心距离

$$D_{ij} = \|\mathbf{Z}_i - \mathbf{Z}_j\|, i = 1, 2, \dots, N_C - 1, j = i + 1, \dots, N_C$$

- 列出类间距离过近者

$$D_{i_1j_1} < D_{i_2j_2} < \cdots < D_{i_lj_l}, l \leq L$$

- 执行合并

将 $\mathbf{Z}_{i_l}, \mathbf{Z}_{j_l}$ 合并得到 $\mathbf{Z}_l = \frac{1}{N_{i_l} + N_{j_l}} [N_{i_l} \mathbf{Z}_{i_l} + N_{j_l} \mathbf{Z}_{j_l}]$, $N_C = N_C - 1$

- 结束步骤

若为最后一次迭代则中止, 否则迭代次数 +1

• DBSCAN

- 核心点: 在半径 Eps 内含有超过 MinPts 数目的点

边界点: 在半径 Eps 内点的数量小于 MinPts, 但是在核心点的邻居

噪音点: 任何不是核心点或边界点的点

- 直接密度可达: 给定一个对象集合 D , 如果 p 在 q 的 Eps 邻域内, 而 q 是一个核心对象, 则称对象 p 从对象 q 出发时是直接密度可达的

密度可达: 如果存在一个对象链 $p_1, p_2, \dots, p_n, p_1 = q, p_n = p$, 对于 $p_i \in D, 1 \leq i \leq n$, p_{i+1} 是从 p_i 关于 Eps 和 MinPts 直接密度可达的, 则对象 p 是从对象 q 关于 Eps 和 MinPts 密度可达的

密度相连: 如果存在对象 $O \in D$, 使对象 p 和 q 都是从 O 关于 Eps 和 MinPts 密度可达的, 那么对象 p 到 q 是关于 Eps 和 MinPts 密度相连

聚类结果的评价

- 聚类中心之间的距离：距离值大，通常可考虑分为不同类
- 聚类域中的样本数目：样本数目少且聚类中心距离远，可考虑是否为噪声
- 聚类域内样本的距离方差：方差过大的样本可考虑是否属于这一类

线性分类器

判别函数

- 定义
 - 直接用来对模式进行分类的准则函数
 - 若分属于 ω_1, ω_2 的两类模式可用一方程 $d(\mathbf{X}) = 0$ 来划分，那么称为判 $d(\mathbf{X}) = 0$ 别函数，或称判决函数、决策函数
- 前提条件
 - 样本所属类别未知
 - 样本类别数已知

线性判别函数

- 一般形式

$$d(\mathbf{X}) = w_1x_1 + w_2x_2 + \cdots + w_nx_n + w_{n+1} = \mathbf{W}_0^T \mathbf{X} + w_{n+1}$$

其中, $\mathbf{X} = [x_1, x_2, \cdots, x_n]^T$, \mathbf{W}_0 为权向量

增广形式: $d(\mathbf{X}) = w_1x_1 + w_2x_2 + \cdots + w_nx_n + w_{n+1} = \mathbf{W}^T \mathbf{X}$

其中, $\mathbf{X} = [x_1, x_2, \cdots, x_n, 1]^T$

- 分类方法

- 两类情况

$$d(\mathbf{X}) = \mathbf{W}^T \mathbf{X} \begin{cases} > 0 & \mathbf{X} \in \omega_1 \\ < 0 & \mathbf{X} \in \omega_2 \end{cases}$$

$d(\mathbf{X}) = 0$: 不可判, 可以 $\mathbf{X} \in \omega_1$ 或 $\mathbf{X} \in \omega_2$ 或拒绝

- $\omega_i/\overline{\omega_i}$ 两分法

$$d(\mathbf{X}) = \mathbf{W}^T \mathbf{X} \begin{cases} > 0 & \mathbf{X} \in \omega_i \\ < 0 & \mathbf{X} \in \overline{\omega_i} \end{cases} \quad i = 1, 2, \cdots, M$$

若只有 $d_i(\mathbf{X}) > 0$, 其余 $d(\mathbf{X}) < 0$, 则判为 ω_i 类

否则为失效区域 (IR)

- ω_i/ω_j 两分法

$$d_{ij}(\mathbf{X}) = \mathbf{W}_{ij}^T \mathbf{X}$$

$$d_{ji} = -d_{ij}$$

若 $d_{ij}(\mathbf{X}) > 0, \forall j \neq i, i, j = 1, 2, \cdots, M$, 则 $\mathbf{X} \in \omega_i$

否则为 IR 区

- ω_i/ω_j 两分法特例

当 $d_{ij}(\mathbf{X}) = d_i(\mathbf{X}) - d_j(\mathbf{X})$

若 $d_i(\mathbf{X}) = \max \{d_k(\mathbf{X}), k = 1, 2, \dots, M\}$, 则 $\mathbf{X} \in \omega_i$

除边界区外没有不确定区域

广义线性判别函数

- 非线性多项式函数

设 $\{\mathbf{X}\}$ 在模式空间 \mathbb{X} 中线性不可分

$$d(\mathbf{X}) = w_1 f_1(\mathbf{X}) + w_2 f_2(\mathbf{X}) + \dots + w_k f_k(\mathbf{X}) + w_{k+1} = \sum_{i=1}^{k+1} w_i f_i(\mathbf{X})$$

其中, $f_{k+1}(\mathbf{X}) = 1$

- 数学表达式

$$\mathbf{X}^* = [x_1^*, x_2^*, \dots, x_k^*, 1]^T = [f_1(\mathbf{X}), f_2(\mathbf{X}), \dots, f_k(\mathbf{X}), 1]^T$$

\mathbb{X}^* 空间的维数 k 高于 \mathbb{X} 的维数 n

$$d(\mathbf{X}) = \mathbf{W}^T \mathbf{X}^* = d(\mathbf{X}^*), \mathbf{W} = [w_1, w_2, \dots, w_k, w_{k+1}]^T$$

线性判别函数几何性质

- 模式空间与超平面

- 基本概念

模式空间: 以 n 维模式向量 \mathbf{X} 的 n 个分类为坐标变量的欧氏空间

超平面: $d(\mathbf{X}) = \mathbf{W}_0^T \mathbf{X} + w_{n+1} = 0$

- 法向量

$$\mathbf{W}_0^T (\mathbf{X}_1 - \mathbf{X}_2) = 0$$

\mathbf{W}_0 与超平面上任意向量正交, 称为超平面的法向量, 方向由超平面的负侧指向正侧

$$\text{单位法线向量 } \mathbf{U} = \frac{\mathbf{W}_0}{\|\mathbf{W}_0\|}$$

- $d(\mathbf{X})$ 取值

\mathbf{X} 在超平面上时, 取值为 0

否则, 将 \mathbf{X} 向超平面投影得到向量 \mathbf{X}_p , 构造向量 $\mathbf{R} = r \cdot \mathbf{U}$, 其中 r 为 \mathbf{X} 到超平面的代数距离

$$\mathbf{X} = \mathbf{X}_p + \mathbf{R}$$

$$d(\mathbf{X}) = r \|\mathbf{W}_0\| \begin{cases} > 0 & \text{超平面正侧} \\ < 0 & \text{超平面负侧} \end{cases}$$

- 权空间和权向量

- 权空间

以线性判别函数的权值为坐标变量的 $n + 1$ 维欧氏空间

- 规范化增广样本向量

对于两类问题: $\omega_1 = \{x_{11}, x_{12}, \dots, x_{1p}\}, \omega_2 = \{x_{21}, x_{22}, \dots, x_{2q}\}$, 令

$$\mathbf{X} = \begin{cases} X_{1i} & i = 1, 2, \dots, p \\ -X_{2i} & i = 1, 2, \dots, q \end{cases}$$

则称 \mathbf{X} 为规范化增广样本向量

若 $p + q$ 个样本的判别函数均大于零，则为解区

感知器算法

- 选择 N 个分属于 ω_1 和 ω_2 的模式识别样本构成训练样本集，编号为 $\mathbf{X}_1, \dots, \mathbf{X}_N$ ，任取权向量初始值 $\mathbf{W}(1)$ 开始迭代

- 第 k 次迭代时，输入一个样本 \mathbf{X}_i

$$\mathbf{W}(k+1) = \begin{cases} \mathbf{W}(k) & \mathbf{W}^T(k)\mathbf{X}_i > 0 \\ \mathbf{W}(k+1) = \mathbf{W}(k) + c\mathbf{X}_i & \mathbf{W}^T(k)\mathbf{X}_i \leq 0 \end{cases}$$

c 为校正增量系数， $c > 0$

- 若上一轮迭代未发生错误，则迭代结束

梯度法

- 选择 N 个分属于 ω_1 和 ω_2 的模式识别样本构成训练样本集，编号为 $\mathbf{X}_1, \dots, \mathbf{X}_N$ ，准则函数为 $J(\mathbf{W}, \mathbf{X})$ ，任取权向量初始值 $\mathbf{W}(1)$ 开始迭代

- 第 k 次迭代时，输入一个样本 \mathbf{X}_i

$$\nabla J(k) = \left. \frac{\partial J(\mathbf{W}, \mathbf{X}_i)}{\partial \mathbf{W}} \right|_{\mathbf{W}=\mathbf{W}(k)}$$

$$\mathbf{W}(k+1) = \mathbf{W}(k) - c\nabla J(k)$$

- 若对所有样本均有 $\nabla J = 0$ ，迭代结束

最小平方误差法 (LMSE)

选择准则函数 $J(\mathbf{W}, \mathbf{X}, \mathbf{B}) = \frac{1}{2} \|\mathbf{XW} - \mathbf{B}\|^2 = \frac{1}{2} \sum_{i=1}^N (\mathbf{W}^T \mathbf{X}_i - b_i)^2$ ，使得当 J 达到最小时，

$\mathbf{XW} = \mathbf{B}$ 可得到最小二乘近似解

- 选择 N 个分属于 ω_1 和 ω_2 的模式识别样本构成训练样本集，将属于 ω_2 的样本乘以 -1 ，写出增广样本矩阵 \mathbf{X}

- 求 \mathbf{X} 的伪逆矩阵 $\mathbf{X} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

- 设置初值 c 和 $\mathbf{B}(1)$ ， c 为正的校正向量， $\mathbf{B}(1)$ 各分量大于零

计算 $\mathbf{W}(1) = \mathbf{X}^\# \mathbf{B}(1)$

- 计算 $\mathbf{e}(k) = \mathbf{XW}(k) - \mathbf{B}(k)$

- 若 $\mathbf{e}(k) = 0$ ，则线性可分，解为 $\mathbf{W}(k)$ ，迭代结束

- 若 $\mathbf{e}(k) > 0$ ，则线性可分，进入下一步

- 若 $\mathbf{e}(k) < 0$ ，检查 $\mathbf{XW}(k)$ 。若 $\mathbf{XW}(k) > 0$ ，有解；否则无解，迭代结束

- 计算 $\mathbf{W}(k+1)$ 和 $\mathbf{B}(k+1)$

- 方法1：计算 $\mathbf{W}(k+1) = \mathbf{W}(k) + c\mathbf{X}^\# |\mathbf{e}(k)|$,

- $\mathbf{B}(k+1) = \mathbf{B}(k) + c[\mathbf{e}(k) + |\mathbf{e}(k)|]$

- 方法2：计算 $\mathbf{B}(k+1) = \mathbf{B}(k) + c[\mathbf{e}(k) + |\mathbf{e}(k)|]$, $\mathbf{W}(k+1) = \mathbf{X}^\# \mathbf{B}(k+1)$

然后返回上一步

势函数法

- 势函数

- I 型势函数：用对称的有限多项式展开, $K(\mathbf{X}, \mathbf{X}_k) = \sum_{i=1}^m \varphi_i(\mathbf{X}_k) \varphi_i(\mathbf{X})$

$\varphi_i(\mathbf{X})$ 在模式定义域内为正交函数集

- II 型势

- 函数：双变量的对称函数 $K(\mathbf{X}, \mathbf{X}_k) = K(\mathbf{X}_k, \mathbf{X})$

- 积累势函数初始值 $K_0(\mathbf{X}) = 0$

- 加入 \mathbf{X}_1

$$K_1(\mathbf{X}) = \begin{cases} K_0(\mathbf{X}) + K(\mathbf{X}, \mathbf{X}_1) & \mathbf{X}_1 \in \omega_1 \\ K_0(\mathbf{X}) - K(\mathbf{X}, \mathbf{X}_1) & \mathbf{X}_1 \in \omega_2 \end{cases}$$

- 加入 \mathbf{X}_{k+1}

$$K_{k+1}(\mathbf{X}) = K_k(\mathbf{X}) + r_{k+1} K(\mathbf{X}, \mathbf{X}_{k+1})$$
$$r_{k+1} = \begin{cases} 0 & \mathbf{X}_{k+1} \in \omega_1, K_k(\mathbf{X}_{k+1}) > 0 \\ 0 & \mathbf{X}_{k+1} \in \omega_2, K_k(\mathbf{X}_{k+1}) < 0 \\ 1 & \mathbf{X}_{k+1} \in \omega_1, K_k(\mathbf{X}_{k+1}) \leq 0 \\ -1 & \mathbf{X}_{k+1} \in \omega_2, K_k(\mathbf{X}_{k+1}) \geq 0 \end{cases}$$

贝叶斯判别准则

研究对象及相关概率

- 条件概率常用公式

- 概率乘法公式 $P(AB) = P(A)P(B|A)$
- 全概率公式 $P(B) = \sum P(A_i)P(B|A_i)$
- 贝叶斯公式 $P(A_i|B) = \frac{P(A_i)P(B|A_i)}{\sum P(A_i)P(B|A_i)}$

- 模式识别的概率

- 先验概率 $P(\omega_i)$
- 后验概率 $P(\omega_i|\mathbf{X})$
- 先验概率 $P(\mathbf{X}|\omega_i)$

$$P(\omega_i|\mathbf{X}) = \frac{P(\mathbf{X}|\omega_i)P(\omega_i)}{P(\mathbf{X})} = \frac{P(\mathbf{X}|\omega_i)P(\omega_i)}{\sum_{i=1}^M P(\mathbf{X}|\omega_i)P(\omega_i)}$$

贝叶斯决策

- 最小错误率贝叶斯决策

- $P(\omega_i|\mathbf{X}) = \max \{P(\omega_j|\mathbf{X})\}$, 则 $\mathbf{X} \in \omega_i$
- $l_{12}(\mathbf{X}) = \frac{P(\mathbf{X}|\omega_1)}{P(\mathbf{X}|\omega_2)} > \frac{P(\omega_2)}{P(\omega_1)}$, 则 $\mathbf{X} \in \begin{cases} \omega_1 \\ \omega_2 \end{cases}$

$$P(\mathbf{X}|\omega_1)P(\omega_1) \begin{matrix} > \\ < \end{matrix} P(\mathbf{X}|\omega_2)P(\omega_2), \text{ 则 } \mathbf{X} \in \begin{cases} \omega_1 \\ \omega_2 \end{cases}$$

- 最小风险贝叶斯决策

- 条件平均风险 $r_i(\mathbf{X}) = \sum_{j=1}^M L_{ij}(\mathbf{X})P(\omega_j | \mathbf{X})$

损失函数 $L_{ij}(\mathbf{X}) \begin{cases} \leq 0 & i = j \\ > 0 & i \neq j \end{cases}$ 指将自然属性是 ω_j 类的样本你决策为 ω_i 类时的是非代价

- $l_{12}(\mathbf{X}) = \frac{P(\mathbf{X} | \omega_1)}{P(\mathbf{X} | \omega_2)} > \frac{(L_{21} - L_{22})P(\omega_2)}{(L_{12} - L_{11})P(\omega_1)} = \theta_{12}$, 则 $\mathbf{X} \in \begin{cases} \omega_1 \\ \omega_2 \end{cases}$

- (0-1) 损失最小风险贝叶斯决策

- $P(\mathbf{X} | \omega_k)P(\omega_k) > P(\mathbf{X} | \omega_i)P(\omega_i)$, 则 $\mathbf{X} \in \omega_k$

- 正态分布模式的贝叶斯决策

- 相关知识

- 二次型 $\mathbf{X}^T \mathbf{A} \mathbf{X}$

- 正定二次型 $\forall \mathbf{X} \neq \mathbf{0}, \mathbf{X}^T \mathbf{A} \mathbf{X} > 0$

- 单变量正态分布 $p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right)$

- 3σ 规则

- 多变量正态随机变量

$$p(\mathbf{X}) = (\mathbf{X} - \mathbf{M})^T \mathbf{C}^{-1} (\mathbf{X} - \mathbf{M})$$

$$\mathbf{M} \text{ 为各维度的均值, 协方差矩阵 } \mathbf{C} = \begin{bmatrix} \sigma_{11}^2 & \cdots & \sigma_{1n}^2 \\ \vdots & & \vdots \\ \sigma_{n1}^2 & \cdots & \sigma_{nn}^2 \end{bmatrix}$$

- $p(\mathbf{X} | \omega_i) = \frac{1}{(2\pi)^{n/2} |\mathbf{C}_i|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{X} - \mathbf{M}_i)^T \mathbf{C}_i^{-1} (\mathbf{X} - \mathbf{M}_i)\right)$

- $d_i(\mathbf{X}) = \ln P(\omega_i) - \frac{1}{2} \ln |\mathbf{C}_i| - \frac{1}{2} \left\{ (\mathbf{X} - \mathbf{M}_i)^T \mathbf{C}_i^{-1} (\mathbf{X} - \mathbf{M}_i) \right\}$

贝叶斯分类器的错误概率

- 性能指标

$$\text{准确率 } A = \frac{TP+TN}{P+N}$$

$$\text{精确率 } P = \frac{TP}{TP+FP}$$

$$\text{虚警概率 } FA = \frac{FP}{TP+FP}$$

$$\text{漏警概率 } MA = \frac{FN}{TN+FN}$$

$$\text{召回率 } R = \frac{TP}{TP+FN}$$

- 先验概率相等时的错误概率

$$P(e) = \int_{r_{ij/2}}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right) dy$$

r_{ij} 为马氏距离

- 错误率估计

- 先验概率未知 $\hat{\varepsilon} = \frac{k}{N}$

- 先验概率已知 $\hat{\varepsilon}' = \sum_{i=1}^2 P(\omega_i) k_i / N_i$

聂曼-皮尔逊决策

- 相关知识

$$\Phi(\lambda) = \int_{-\infty}^{\lambda} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$$

$$\Phi(\lambda) = \int_{-\infty}^{\lambda} \varphi(x) dx$$

当 $\lambda < 0$ 时, $\Phi(\lambda) = 1 - \Phi(-\lambda)$

- $\frac{P(\mathbf{X} | \omega_1)}{P(\mathbf{X} | \omega_2)} > \mu$, 则 $\mathbf{X} \in \begin{cases} \omega_1 \\ \omega_2 \end{cases}$
 $P_2(e) = \int_{-\infty}^{t(\mu)} p(\mathbf{X} | \omega_2) d\mathbf{X}$

特征选择与特征提取

K-L 变换

- K-L 展开式

对随机变量集合 $\{\mathbf{X}\}$ 做离散正交展开, 正交向量系为 $\{\mathbf{u}_j\}$, 得到 $\mathbf{X} = \sum_{j=1}^n a_j \mathbf{u}_j$

$$\mathbf{u}_i^T \mathbf{u}_j = \begin{cases} 1 & j = i \\ 0 & j \neq i \end{cases}$$

将 $\{\mathbf{u}_j\}$ 的特征值排序 $\lambda_1 > \lambda_2 > \dots > \lambda_m > \dots > \lambda_n \geq 0$

$$\mathbf{X} = \sum_{j=1}^m a_j \mathbf{u}_j$$

- 利用 K-L 变换特征提取

- 求总体自相关矩阵 $\mathbf{R} = E[\mathbf{X}\mathbf{X}^T] \approx \frac{1}{N} \sum_{j=1}^N \mathbf{X}_j \mathbf{X}_j^T$
- 求出并取前 m 大的特征值 λ_i , 算出特征向量 \mathbf{u}'_j , 归一化得到 \mathbf{u}_j
- $\mathbf{X}^* = \mathbf{U}^T \mathbf{X}$

类别可分性测度

- 基于距离的可分性测度

- 类内距离

$$\overline{D^2} = E\{\|\mathbf{X}_i - \mathbf{X}_j\|^2\} = E\{(\mathbf{X}_i - \mathbf{X}_j)^T (\mathbf{X}_i - \mathbf{X}_j)\}$$

若样本相互独立, 则 $\overline{D^2} = 2 \sum_{k=1}^n \sigma_k^2$

- 类内散布矩阵: 该类分布的协方差矩阵
越小越好

- 类间距离

$$\overline{D_b^2} = \sum_{i=1}^c P(\omega_i) \|\mathbf{M}_i - \mathbf{M}_0\|^2 = \sum_{i=1}^c P(\omega_i) (\mathbf{M}_i - \mathbf{M}_0)^T (\mathbf{M}_i - \mathbf{M}_0)$$

其中, $P(\omega_i)$ 为先验概率, \mathbf{M}_i 为 ω_i 类的均值向量, $\mathbf{M}_0 = E\{\mathbf{X}\} = \sum_{i=1}^c P(\omega_i) \mathbf{M}_i$

- 类间散布矩阵

$$\mathbf{S}_b = \sum_{i=1}^c P(\omega_i) (\mathbf{M}_i - \mathbf{M}_0)(\mathbf{M}_i - \mathbf{M}_0)^T$$

$$\overline{D_b^2} = \text{tr} \{ \mathbf{S}_b \}$$

越大越好

- 多类情况距离

$$J_d = \frac{1}{2} \sum_{i=1}^c P(\omega_i) \sum_{j=1}^c P(\omega_j) \frac{1}{n_i n_j} \sum_{k=1}^{n_i} \sum_{l=1}^{n_j} D^2(\mathbf{X}_k^i, \mathbf{X}_l^j)$$

其中, \mathbf{X}_k^i 为 ω_i 类的第 k 个样本, n_i 为 ω_i 的样本数

$$J_d = \sum_{i=1}^c P(\omega_i) \left[\frac{1}{n_i} \sum_{k=1}^{n_i} (\mathbf{X}_k^i - \mathbf{M}_i)^T (\mathbf{X}_k^i - \mathbf{M}_i) + (\mathbf{M}_i - \mathbf{M}_0)^T (\mathbf{M}_i - \mathbf{M}_0) \right]$$

- 多类情况散布矩阵

$$\text{多类类间: } \mathbf{S}_b = \sum_{i=1}^c P(\omega_i) (\mathbf{M}_i - \mathbf{M}_0)(\mathbf{M}_i - \mathbf{M}_0)^T$$

$$\text{多类类内: } \mathbf{S}_w = \sum_{i=1}^c P(\omega_i) E \left\{ (\mathbf{X} - \mathbf{M}_i)(\mathbf{X} - \mathbf{M}_i)^T \right\}$$

$$\text{多类总体: } \mathbf{S}_t = E \left\{ (\mathbf{X} - \mathbf{M}_0)(\mathbf{X} - \mathbf{M}_0)^T \right\} = \mathbf{S}_b + \mathbf{S}_w$$

$$J_d = \text{tr}(\mathbf{S}_t)$$

- 基于概率分布的可分性测度

- 散度

$$\omega_i \text{ 类对 } \omega_j \text{ 类的散度 } J_{ij} = \int_X [p(\mathbf{X} | \omega_i) - p(\mathbf{X} | \omega_j)] \ln \frac{p(\mathbf{X} | \omega_i)}{p(\mathbf{X} | \omega_j)} d\mathbf{X}$$

越大越好

- $J_{ij} = J_{ji}$
- $J_{ij} \geq 0$
- 具有可加性
- 值越大, 错误率越小

- 两个正态分布模式类的散度

$$J_{ij} = (\mathbf{M}_i - \mathbf{M}_j)^T \mathbf{C}^{-1} (\mathbf{M}_i - \mathbf{M}_j)$$

$$\text{一维正态分布时 } J_{ij} = \frac{(m_i - m_j)^2}{\sigma^2}$$

特征选择

神经网络
