

10.3 分类

① 分类的定义

将已知类别的训练数据输入机器学习算法模型进行训练，用于将未知类别的测试数据划分成已知类别。核心操作位

② 原理

输入：设空间 $X \subseteq \mathbb{R}^N$ ($N \geq 1$)

输出：设输出空间 $Y = \{c_1, c_2, \dots, c_k\}$ $k \in \mathbb{C}(\mathbb{Z}^+)$ c_k 表示类别。

设 $D = X \times Y$ (笛卡尔乘积) 是一个未知数据集

目标函数： $f: X \rightarrow Y$ ， f 为分类函数，将 X 映射到 Y

训练集： $S = \{(x_i, y_k) \mid x_i \in X, y_k = c_k \in Y, i \in [1, m], k \in [1, k]\}$

其中 $(x_i, y_k) \in D$

设假设函数集 H 是 $X \rightarrow Y$ 的函数集合，目标：通过训练得到 $\hat{f} \in H$ ：满足

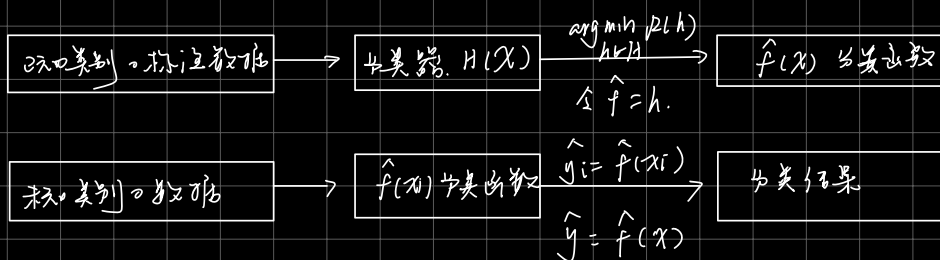
是： $\hat{f}: X \rightarrow Y$ ，并且 \hat{f} 和 f 有最小的估计误差和泛化误差

$$\hat{f} = \arg \min_{h \in H} P(h) = \arg \min_{h \in H} [E(h(X) \neq f(X))]$$

并用 \hat{f} 对未知数据进行分类。

给定具有 m' 个未知类别的 $T \subseteq X$ ， $T = \{x_j \mid x_j \in X, j \in [1, m']\}$

故 $\hat{f}(T) = \{(x_j, y_k) \mid j \in [1, m'], y_k = h(x_j) = c_k \in Y, h \in [1, k]\}$



② 线性与非线性 : 主要区别 线性为类器 非线性为类器

1° 线性为类器. 特点 : 线性判别函数. 线性决策边界.

(i) 线性判别函数 : $y(x) = w \cdot x + b = \sum_{i=0}^N w_i x_i + b$

其中 $w \in \mathbb{R}^N$: $w = (w_1, \dots, w_N)$

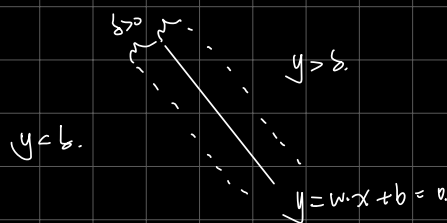
$x \in \mathcal{X}$: $x = (x_1, \dots, x_N)^T$

$b \in \mathbb{R}$: 偏差.

(ii) 线性决策边界

对 $b > 0$: (目标函数 $y(x) = w \cdot x + b$ 令 $y(x) = 0$)

$|y(x)| < b$ 则 $y(x) \rightarrow 0$



2° 非线性为类器 : 特点 : 若干非线性决策边界 (可分非凸集)

$$y(x) = w_0 + w_1 x + w_2 x^2 + \dots + w_N x^N = \sum_{i=1}^N w_i x^i$$

判别函数 : $y(x) = \sum_{i=1}^N w_i \phi(x_i) = w^T \Phi(x)$

其中 $\Phi(x) = (x, x^2, \dots, x^N)$ 核函数.

④ 类别划分

类别个数 = 决策个数 , $y_k(x) = w_k \cdot x + b$

决策-位置 取关于超平面位置.

10.4. 回归

① 定义

回归用于估计因变量与自变量之间的关系，并给出预测值

② 符号

空间: \mathcal{R}^N ($N \geq 1$) 非负实数集

输入空间: $\mathcal{X} \subseteq \mathcal{R}^N$

其中 $\mathcal{D} = \mathcal{X} \times \mathcal{Y}$ 是未知分布

输出空间: \mathcal{Y}

目标函数: $f: \mathcal{X} \rightarrow \mathcal{Y}$

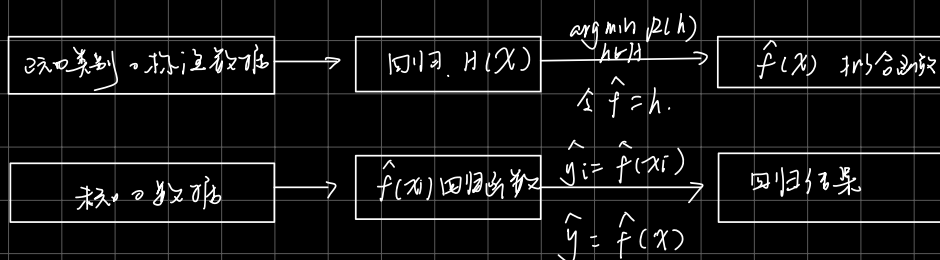
输入数据: $S = \{(x_i, y_i) \mid (x_i, y_i) \in \mathcal{D}, i = [1, m]\}$

估计函数: $\hat{f} = \argmin_{h \in H} p(h) = \argmin_{h \in H} E_{\mathcal{D}} [L(h(x), f(x))]$

其中 $h \in H$, h 是候选函数

未知数据: $T = \{x_j \mid x_j \in \mathcal{X}, j \in [1, m']\}$

输出数据: $\hat{f}(T) = h(T) = \{y_j = h(x_j) \mid j \in [1, m'], h(x_j) \in \mathcal{R}^N\}$



③ 线性与非线性

1° 线性回归函数是线性的组合. $y(x) = \sum_{i=0}^N w_i \cdot x^i + c$ (c 为常数) $c \sim N(\mu, \sigma^2)$

$$p(y|x, \theta) = N(y | \mu(x), \sigma^2(x)) \quad \mu(x) = w^T x \quad \sigma^2(x) = \sigma^2$$

$$\theta = (w, \sigma^2) \quad \text{K:} \quad p(y|x, \theta) = N(y | w^T x, \sigma^2)$$

∴ 非线性性: $y(x) = \sum_{i=0}^n w_i x^i \Rightarrow \sum_{i=0}^n w_i \phi_i(x) = w^T \cdot \Phi(x)$ (Φ 为核函数)

$$p(y|x, \theta) = N(y | w^T \Phi(x), \sigma^2)$$

10.5.

(不同)	函数	输出	训练目标函数	计算方式	目的
分类	分类函数	离散	$\arg \max_{h \in H} p(h)$	与训练样本距离最小	划分成不同类别
回归	回归函数	连续	$\arg \max_{\theta} p(y x, \theta)$	后验概率最大	估计(预测)数值

(相同) . 分类和回归都是统计学习方法, 训练和预测流程类似.

10.6. 聚类

① 定义:

将未标注数据通过训练寻找数据之间内在关系, 将数据分成不同类别.

② 表示:

空间: \mathcal{X}^N ($N \geq 1$) 非负实数集

输入空间: $\mathcal{X} \subseteq \mathcal{X}^N$

其中 $\mathcal{D} = \mathcal{X} \times \mathcal{Y}$ 为未知分布.

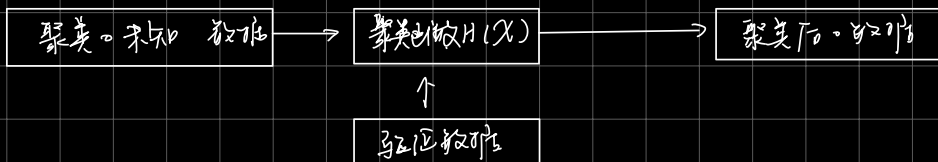
输出空间: \mathcal{Y}

最佳聚类函数: $f: \mathcal{X} \rightarrow \mathcal{Y}$ 且 $f \in H$ (H 为假设函数集)

输入数据 (未标注): $T = \{x_i | x_i \in \mathcal{X}, i \in [1, m]\}$

输出结果: $y' = \{c_1, c_2, \dots, c_k\}$ $k \leq K$. K 为假设值

$$y' = f(T) = \{y_i = h(x_i) | x_i \in \mathcal{X}, i \in [1, m], y_k = c_k \in y', k \in [1, k]\}$$



② 聚类类型

{	硬性	硬性划分类别, 在该类别计算的均值.
	软性	按照概率分布分配数据值. 计算也是 (软)

1° 基于连通性聚类 (层次聚类) : 根据对象之间距离

(i) 单链接聚类 $\min \{d(x, y) : x \in A, y \in B\}$

(ii) 全链接聚类 $\max \{d(x, y) : x \in A, y \in B\}$

(iii) 平均链接聚类 $\frac{1}{|A| \cdot |B|} \sum_{x \in A} \sum_{y \in B} d(x, y)$

2° 基于质心聚类

先设个聚类个数, 寻找聚类中心, 将数据分配到聚类中心, 使得该聚类中数据与聚类中心平方距离最小.

输入数据: $T = \{x_1, x_2, \dots, x_n\}$

设 $k = K$, $C = \{c_1, \dots, c_K\}$ K 个类别

设 μ_k 是第 k 个聚类中的均值, 则目标函数为

$$f = \arg \min_C \sum_{k=1}^K \sum_{x \in C_k} \|x - \mu_k\|^2$$

算法流程: (迭代次数 t)

(i) $t=0$. 随机出初始 K 个聚类中心.

repeat

(ii) 将输入数据分配到 K 个聚类, 使得平方欧氏距离最小.

$$\forall k, 1 \leq k \leq K: C_k^{(t)} = \{x \mid \|x - \mu_k^{(t)}\|^2 \leq \|x - \mu_l^{(t)}\|^2\}$$

(iii) 更新: 计算新均值 (质心)

$$\mu_k^{(t+1)} = \frac{1}{|C_k^{(t)}|} \sum_{x \in C_k^{(t)}} x$$

until 质心不再变化.

(iv) 算法终止.

3° 基于分布聚类 : 寻找属于同一分布概率最大的数据对象 (期望和方差最大化)

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

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

$$z_i^* = \underset{k}{\operatorname{argmax}} r_{ik} = \underset{k}{\operatorname{argmax}} \log p(x_i | z_i = k, \theta) + \log p(z_i = k | \theta)$$

4° 基于密度聚类 : 根据欧式距离定义为

$$\rho_i = \sum_j \chi(d_{ij} - d_c) \chi(x) = \begin{cases} 1 & x \leq 0 \\ 0 & x > 0 \end{cases}$$

其中, $d_{ij} = \operatorname{dist}(x_i, x_j)$, d_c 为截断距离

$$b_i \text{ 是第 } i \text{ 数据点与其他数据点的最小距离} \quad b_i = \min_{j: j \neq i} (d_{ij})$$

10.7

(不同点)

	原型	数据点	训练数据	训练数据	过程
分类	将数据点分为已知类别	已知类别	✓	寻找超平面 (决策边界)	$\hat{f} = \underset{f \in H}{\operatorname{argmin}} p(f)$
聚类	寻找内在数据点为簇	未知类别	✗	硬性/软性	迭代直到收敛结束

相同点: 分类和聚类问题均数据点进行划分, 得到划分结果

10.8 排名

① 定义: 基于某种准则将数据点进行排序.

② 前提: 设输入空间: X .

未知数据: $D = x \times x$

输出空间: $Y = \{x \times x, \geq\}$ 排序集 $\forall (x_i, x_j) \in D: x_i \geq x_j$

映射函数: $f: x \times x \rightarrow Y$

训练集 S : $S = \{(x_i, x_j, \geq) | (x_i, x_j) \in D, i \neq j \text{ and } i, j \in [1, m]\}$

$$\hat{f} = \underset{h \in H}{\operatorname{argmin}} f(h) = \underset{h \in H}{\operatorname{argmin}} P_{r(x_i, y_i)} [f(x_i, y_i)] \approx \underset{h \in H}{\operatorname{argmin}} [f(x_i, y_i) \neq 0 \wedge f(x_i, y_i) < h(x_i) - h(x_i)] \leq 0]$$

③ 类型

1° 基于距离 (欧式距离)

2° 基于相似 (余弦值)

3° 基于相关性:

样本数足够:

$$r_{x,y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$\bar{x} = \frac{1}{n} \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \bar{y} = \frac{1}{n} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}$$

群内 pearson.

$$\rho_{X,Y} = \frac{\operatorname{cov}(X,Y)}{\sigma_X \sigma_Y}$$

$$r_s = \rho_{r_{gx}, r_{gy}} = \frac{\operatorname{cov}(r_{gx}, r_{gy})}{\sigma_{r_{gx}} \sigma_{r_{gy}}}$$

10.9 降维 ① 将高维数据映射到低维空间, 并保留数据的主要特征

② 降维: 输入空间 $X^k \subseteq \mathbb{R}^N$ ($N \geq 3$)

输出空间 $Y^L \subseteq \mathbb{R}^N$

高维数据集 $S = \{x_1, \dots, x_m\} \subseteq X^{k \times m}$

映射集 $\tilde{S}: X^k \rightarrow Y^L$

映射后的数据集: $T = \{y_i | i \in [1, m], y_i \in \phi(x_i)\} \in Y^{L \times m}$

③ 线性与非线性:

1° { (i) PLA, LCA, CCA, LDA 寻找投影矩阵 (线性降维)
(ii) 非线性降维

2° 非线性

(i) 马尔可夫

(ii) 核函数

(iii) 流形学习