

3.1 基本形式

$$x = (x_1; x_2; \dots; x_d).$$

$$f(x) = w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b.$$

$$f(x) = w^T x + b. \quad w = (w_1; w_2; \dots; w_d).$$

3.2 线性回归

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}, \text{ 其中 } x_i = (x_{i1}; x_{i2}; \dots; x_{id}), y_i \in \mathbb{R}.$$

输入属性为 x : $D = \{(x_i, y_i)\}_{i=1}^m, x_i \in \mathbb{R}^d$. 由 $f(x_i) = w x_i + b$ 使 $f(x_i) \approx y_i$. \leftarrow 回归目标.

均方误差: $(w^*, b^*) = \arg \min_{(w, b)} \sum_{i=1}^m (f(x_i) - y_i)^2 \leftarrow$ 欧氏距离.

基于均方误差最小化为最小二乘法.

$$E(w, b) = \sum_{i=1}^m (y_i - w x_i - b)^2 \text{ 对 } w, b \text{ 求导, 得 } \begin{cases} \frac{\partial E(w, b)}{\partial w} = 2(w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m y_i x_i) \\ \frac{\partial E(w, b)}{\partial b} = 2(mb - \sum_{i=1}^m y_i) \end{cases}$$

令导数为0, 得解.

多元线性回归: $f(x_i) = w^T x_i + b$, 使 $f(x_i) \approx y_i$. $\hat{w} = (w, b)$. D 变为 $m(d+1)$ 的矩阵 X

$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} & 1 \\ x_{21} & x_{22} & \dots & x_{2d} & 1 \\ \vdots & \vdots & \dots & \vdots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{md} & 1 \end{pmatrix} = \begin{pmatrix} x_1^T & 1 \\ x_2^T & 1 \\ \vdots & \vdots \\ x_m^T & 1 \end{pmatrix}$$

$$y = (y_1, y_2, \dots, y_m).$$

$$\hat{w}^* = \arg \min_{\hat{w}} (y - X \hat{w})^T (y - X \hat{w}).$$

令 $E \hat{w} = (y - X \hat{w})^T (y - X \hat{w})$ 对 \hat{w} 求导, 得 $\frac{\partial E \hat{w}}{\partial \hat{w}} = 2X^T(X \hat{w} - y)$. 矩阵微分公式: $\frac{\partial a^T x}{\partial x} = \frac{\partial x^T a}{\partial x} = a$

当 $X^T X$ 为满秩矩阵. $\hat{w}^* = (X^T X)^{-1} X^T y$ 令 $\hat{x}_i = (x_i, 1)$, 则 $f(\hat{x}_i) = \hat{x}_i^T (X^T X)^{-1} X^T y$

(列 \Rightarrow 行).

当 $X^T X$ 为非满秩矩阵. 正则化

广义线性回归: $y = g^{-1}(w^T x + b)$. $g(\cdot)$ 联系函数. $g(\cdot) = \ln(\cdot)$ $\ln y = w^T x + b$. 对数

3.3. 对数几率回归 (其实是分类).

二分类. $y \in \{0, 1\}$. $z = w^T x + b$.

单位阶跃函数. $y = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \end{cases} \rightarrow y = \frac{1}{1 + e^{-z}} \text{ sigmoid}$

$$y = \frac{1}{1 + e^{-(w^T x + b)}} \quad \ln \frac{y}{1-y} = w^T x + b. \quad \begin{matrix} \text{正例} \\ \ln \frac{y}{1-y} \text{ 对数几率} \\ \text{反例} \end{matrix}$$

↓ y 视为后验概率估计.

$$\ln \frac{p(y=1|x)}{p(y=0|x)} = w^T x + b.$$

$$\therefore p(y=1|x) = \frac{e^{w^T x + b}}{1 + e^{w^T x + b}}$$

$$p(y=0|x) = \frac{1}{1 + e^{w^T x + b}}$$

用极大似然 EMT 估计 w, b .

$$\text{给定 } \{(x_i, y_i)\}_{i=1}^m, \quad l(w, b) = \sum_{i=1}^m \ln p(y_i | x_i; w, b).$$

每个样本属于其真实标记概率越大越好.

$$\text{令 } \beta = (w, b), \quad \hat{x} = (x; 1). \quad \text{则 } w^T x + b = \beta^T \hat{x}.$$

$$\text{令 } p_1(\hat{x}; \beta) = p(y=1 | \hat{x}; \beta), \quad p_0(\hat{x}; \beta) = p(y=0 | \hat{x}; \beta) = 1 - p_1(\hat{x}; \beta).$$

$= [p_1(\hat{x}; \beta)]^{y_i} [p_0(\hat{x}; \beta)]^{1-y_i} \text{ 也可得证.}$

$$\text{则 } p(y_i | x_i; w, b) = y_i p_1(\hat{x}_i; \beta) + (1 - y_i) p_0(\hat{x}_i; \beta).$$

$$\therefore l(\beta) = \sum_{i=1}^m (-y_i \beta^T \hat{x}_i + \ln(1 + e^{\beta^T \hat{x}_i})).$$

$$\text{梯度下降. } \beta^* = \arg \min_{\beta} l(\beta)$$

3.4. 线性判别分析. LDA.

给定训练集, 同类在直线的投影点尽可能接近, 异类排列在投影区尽可能远.

$$D = \{(x_i, y_i)\}_{i=1}^m, \quad y_i \in \{0, 1\}, \quad x_i \text{ 集合. } \mu_i \text{ 均值. } \Sigma_i \text{ 协方差矩阵.}$$

将数据投影至直线 w . 则两类样本协方差为 $w^T \Sigma_0 w, w^T \Sigma_1 w$.

$$w^T \Sigma_0 w + w^T \Sigma_1 w \text{ 小, 而 } \|w^T \mu_0 - w^T \mu_1\|_2^2 \text{ 尽可能大.}$$

$$J = \frac{\|w^T \mu_0 - w^T \mu_1\|_2^2}{w^T \Sigma_0 w + w^T \Sigma_1 w} = \frac{w^T (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T w}{w^T (\Sigma_0 + \Sigma_1) w}.$$

$$\text{定义类间散度矩阵 } S_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T$$

$$\therefore J = \frac{w^T S_b w}{w^T S_w w}, \quad \text{即 } S_b, S_w \text{ 的广义瑞利商.}$$

$$\text{由于 } w \text{ 在分子分母都为二次项, 令 } w^T S_w w = 1. \quad \text{原式} = \max_w w^T S_b w \quad \text{s.t. } w^T S_w w = 1.$$

$$\text{由拉格朗日乘法, 上式等价于 } S_b w = \lambda S_w w. \quad \text{由于 } S_b w \text{ 方向为 } \mu_0 - \mu_1, \therefore S_b w = \lambda (\mu_0 - \mu_1)$$

$\therefore w = S_w^{-1}(\mu_0 - \mu_1)$. 对 S_w 奇异值分解, $S_w = UZV^T$. 再由 $S_w^{-1} = VZ^{-1}U^T$ 得 S_w^{-1} .

当两类数据同先验, 满足高斯分布且协方差相同. LDA 分类最优.

推广: 多分类, N 类, 类训练数 m_i . 则全局散度矩阵 $S_t = S_b + S_w = \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T$. $S_w = \sum_{i=1}^N S_{w_i}$. $S_{w_i} = \sum_{x \in X_i} (x - \mu_i)(x - \mu_i)^T$

$$S_b = S_t - S_w = \sum_{i=1}^N m_i (\mu_i - \mu)(\mu_i - \mu)^T.$$

采用目标 $\max_w \frac{tr(w^T S_b w)}{tr(w^T S_w w)}$. $w \in R^{d \times (N-1)}$. $S_b w = \lambda S_w w$.

\downarrow
投影矩阵. \longrightarrow 投影到 $N-1$ 维空间 (\ll 样本维数).

LDA 为经典监督降维.

3.5 多分类学习.

N 类, $C_1 \dots C_N$. 拆解法 \rightarrow 若干二分类.

- 对- $\rightarrow \frac{N(N-1)}{2}$ 个二分类. 投票产生最终结果.
- 对多 $\rightarrow N$ 个二分类. 一正 vs 其余反. 置信度最大为分类结果.

$D = \{(x_1, y_1), \dots, (x_m, y_m)\}$, $y_i \in \{C_1, \dots, C_N\}$.

多对多 \rightarrow 纠错输出码 EOC.
(有冗余冗余能力).

(\rightarrow DAG).

编码: 对 N 类划分 M 次. 每次二分类训练. 共 $M \times N$ 个模型.

编码矩阵

解码: 预测时记组成编码. 其中与各自编码比较距离最小的类别为最终结果.

3.6 类别不平衡问题.

m^+, m^- 训练集中正反例子数目. $\frac{y}{1-y} > \frac{m^+}{m^-}$ 预测为正例.

$$\text{再缩放: } \frac{y'}{1-y'} = \frac{y}{1-y} \times \frac{m^-}{m^+}$$

由于无偏采样经常不成立.

欠采样. 去除一些反例. EasyEnsemble. 将反例划为不同学习器. 减少信息丢失.
过采样. 增加一些正例. 不可简单重复采样. 会过拟合. better: SMOTE
阈值移动. 用 $\frac{y'}{1-y'}$ 替代 $\frac{y}{1-y}$. 代价敏感. 将 m^-/m^+ 换为 $\text{cost}^+/\text{cost}^-$.