

四、线性模型

目录

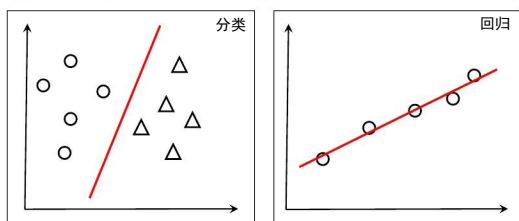
线性回归与逻辑回归

线性判别

多分类

非平衡分类

线性模型 (linear regression)



线性模型(linear model)试图学得一个通过属性的线性组合来进行预测的函数

$$f(x) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b$$

向量形式: $f(x) = w^T x + b$ 简单、基本、可理解性好

线性回归 (linear regression)

$$f(x) = wx_i + b \text{ 使得 } f(x_i) \simeq y_i$$

离散属性的处理: 若有“序”(order), 则连续化;
否则, 转化为 k 维向量

$$\begin{aligned} \text{令均方误差最小化, 有 } (w^*, b^*) &= \arg \min_{(w, b)} \sum_{i=1}^m (f(x_i) - y_i)^2 \\ &= \arg \min_{(w, b)} \sum_{i=1}^m (y_i - wx_i - b)^2 \end{aligned}$$

对 $E_{(w, b)} = \sum_{i=1}^m (y_i - wx_i - b)^2$ 进行最小二乘参数估计

线性回归

分别对 w 和 b 求导:

$$\begin{aligned} \frac{\partial E_{(w, b)}}{\partial w} &= 2 \left(w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b) x_i \right) \\ \frac{\partial E_{(w, b)}}{\partial b} &= 2 \left(mb - \sum_{i=1}^m (y_i - wx_i) \right) \end{aligned}$$

令导数为 0, 得到闭式(closed-form)解:

$$w = \frac{\sum_{i=1}^m y_i (x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{m} \left(\sum_{i=1}^m x_i \right)^2} \quad b = \frac{1}{m} \sum_{i=1}^m (y_i - wx_i)$$

多元(multi-variate)线性回归

$$f(x_i) = w^T x_i + b \text{ 使得 } f(x_i) \simeq y_i$$

$$x_i = (x_{i1}; x_{i2}; \dots; x_{id}) \quad y_i \in \mathbb{R}$$

把 w 和 b 吸收入向量形式 $\hat{w} = (w; b)$, 数据集表示为

$$X = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1d} & 1 \\ x_{21} & x_{22} & \cdots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{md} & 1 \end{pmatrix} = \begin{pmatrix} x_1^T & 1 \\ x_2^T & 1 \\ \vdots & \vdots \\ x_m^T & 1 \end{pmatrix} \quad y = (y_1; y_2; \dots; y_m)$$

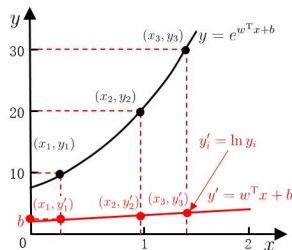
线性模型的变化

对于样例 (x, y) , $y \in \mathbb{R}$, 若希望线性模型的预测值逼近真实标记, 则得到线性回归模型 $y = \mathbf{w}^T \mathbf{x} + b$

令预测值逼近 y 的衍生物?

若令 $\ln y = \mathbf{w}^T \mathbf{x} + b$
则得到对数线性回归
(log-linear regression)

实际是在用 $e^{\mathbf{w}^T \mathbf{x} + b}$ 逼近 y



广义(generalized)线性模型-逻辑回归

一般形式: $y = g^{-1}(\mathbf{w}^T \mathbf{x} + b)$

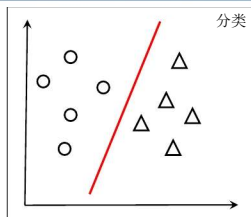
单调可微的 联系函数 (link function)

令 $g(\cdot) = \ln(\cdot)$ 则得到 对数线性回归

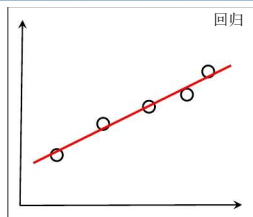
$$\ln y = \mathbf{w}^T \mathbf{x} + b$$

...

线性模型做“分类”



如何“直接”做分类?



例如, 对率回归

广义线性模型;
通过“联系函数”

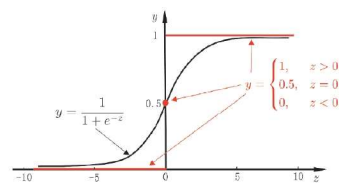
二分类任务

线性回归模型产生的实值输出 $z = \mathbf{w}^T \mathbf{x} + b$
期望输出 $y \in \{0, 1\}$

找 z 和 y 的
联系函数

理想的“单位阶跃函数”
(unit-step function)

$$y = \begin{cases} 0, & z < 0; \\ 0.5, & z = 0; \\ 1, & z > 0, \end{cases}$$



性质不好,
需找“替代函数”
(surrogate function)

单调可微、任意阶可导

$$y = \frac{1}{1 + e^{-z}}$$

对数几率函数
(logistic function)
简称“对率函数”

对率回归-逻辑回归

逻辑回
归

以对率函数为联系函数:

$$y = \frac{1}{1 + e^{-z}} \quad \text{变为} \quad y = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

即:

$$\ln\left(\frac{y}{1-y}\right) = \mathbf{w}^T \mathbf{x} + b$$

“对数几率”
(log odds, 亦称 logit) 几率(odds), 反映了 \mathbf{x} 作为正例的相对可能性

“对数几率回归”(logistic regression)
简称“对率回归”

- 无需事先假设数据分布
- 可得到“类别”的近似概率预测
- 可直接应用现有数值优化算法求取最优解

注意: 它是
分类学习算法!

求解思路

逻辑回
归

若将 y 看作类后验概率估计 $p(y = 1 | \mathbf{x})$, 则

$$\ln \frac{y}{1-y} = \mathbf{w}^T \mathbf{x} + b \quad \text{可写为} \quad \ln \frac{p(y = 1 | \mathbf{x})}{p(y = 0 | \mathbf{x})} = \mathbf{w}^T \mathbf{x} + b$$

$\theta^T \mathbf{x}$

$$h_{\theta}(\mathbf{x}) = g(\theta^T \mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$$

对于输入 \mathbf{x} 分类结果为类别1和类别0的概率分别为:

$$P(y = 1 | \mathbf{x}; \theta) = h_{\theta}(\mathbf{x})$$

$$P(y = 0 | \mathbf{x}; \theta) = 1 - h_{\theta}(\mathbf{x})$$

“极大似然法”

求解思路

逻辑回
归给定数据集 $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$

最大化“对数似然”(log-likelihood)函数:

$$l(\mathbf{w}, \mathbf{b}) = \sum_{i=1}^m \ln p(y_i | \mathbf{x}_i; \theta)$$

x被正确分类的概率:

$$P(y | \mathbf{x}; \theta) = (h_\theta(\mathbf{x}))^y (1 - h_\theta(\mathbf{x}))^{1-y}$$

求解思路

逻辑回
归

$$P(y | \mathbf{x}; \theta) = (h_\theta(\mathbf{x}))^y (1 - h_\theta(\mathbf{x}))^{1-y}$$

取似然函数为:

$$L(\theta) = \prod_{i=1}^m P(y^{(i)} | \mathbf{x}^{(i)}; \theta) = \prod_{i=1}^m (h_\theta(\mathbf{x}^{(i)}))^{y^{(i)}} (1 - h_\theta(\mathbf{x}^{(i)}))^{1-y^{(i)}}$$

对数似然函数为:

$$l(\theta) = \log L(\theta) = \sum_{i=1}^m (y^{(i)} \log h_\theta(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(\mathbf{x}^{(i)})))$$

对数似然函数为:

$$l(\theta) = \log L(\theta) = \sum_{i=1}^m (y^{(i)} \log h_\theta(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(\mathbf{x}^{(i)})))$$

逻辑回归的损失函数为:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h_\theta(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(\mathbf{x}^{(i)}))]$$

$$\text{Cost}(h_\theta(\mathbf{x}), y) = \begin{cases} -\log(h_\theta(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_\theta(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(\mathbf{x}^{(i)}), y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h_\theta(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(\mathbf{x}^{(i)}))]$$

$$J(\theta) = -\frac{1}{m} l(\theta)$$

$$h_\theta(\mathbf{x}) = g(\theta^T \mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$$

逻辑回
归

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(\mathbf{x}^{(i)}), y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h_\theta(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(\mathbf{x}^{(i)}))]$$

梯度下降法求的最小值: $\theta_j = \theta_j - \alpha \frac{\delta}{\delta \theta_j} J(\theta)$

$$\frac{\delta}{\delta \theta_j} J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \frac{1}{h_\theta(\mathbf{x}^{(i)})} \frac{\delta}{\delta \theta_j} h_\theta(\mathbf{x}^{(i)}) - (1 - y^{(i)}) \frac{1}{1 - h_\theta(\mathbf{x}^{(i)})} \frac{\delta}{\delta \theta_j} h_\theta(\mathbf{x}^{(i)}) \right)$$

$$= -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \frac{1}{g(\theta^T \mathbf{x}^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - g(\theta^T \mathbf{x}^{(i)})} \right) \frac{\delta}{\delta \theta_j} g(\theta^T \mathbf{x}^{(i)})$$

$$= -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \frac{1}{g(\theta^T \mathbf{x}^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - g(\theta^T \mathbf{x}^{(i)})} \right) g(\theta^T \mathbf{x}^{(i)}) (1 - g(\theta^T \mathbf{x}^{(i)})) \frac{\delta}{\delta \theta_j} \theta^T \mathbf{x}^{(i)}$$

$$\frac{\delta}{\delta \theta_j} J(\theta)$$

逻辑回
归

$$= -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \frac{1}{g(\theta^T \mathbf{x}^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - g(\theta^T \mathbf{x}^{(i)})} \right) g(\theta^T \mathbf{x}^{(i)}) (1 - g(\theta^T \mathbf{x}^{(i)})) \frac{\delta}{\delta \theta_j} \theta^T \mathbf{x}^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m (y^{(i)} (1 - g(\theta^T \mathbf{x}^{(i)})) - (1 - y^{(i)}) g(\theta^T \mathbf{x}^{(i)})) \mathbf{x}_j^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - g(\theta^T \mathbf{x}^{(i)})) \mathbf{x}_j^{(i)}$$

$$= \frac{1}{m} \sum_{i=1}^m (h_\theta(\mathbf{x}^{(i)}) - y^{(i)}) \mathbf{x}_j^{(i)}$$

求解思路

逻辑回
归

$$\theta_j := \theta_j - \alpha \frac{\delta}{\delta \theta_j} J(\theta)$$

$$\frac{\delta}{\delta \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(\mathbf{x}^{(i)}) - y^{(i)}) \mathbf{x}_j^{(i)}$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(\mathbf{x}^{(i)}) - y^{(i)}) \mathbf{x}_j^{(i)}$$

目录

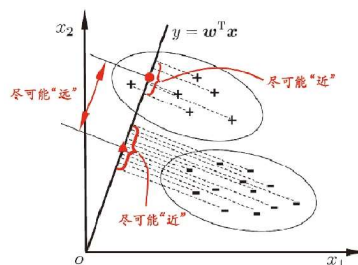
线性归与逻辑回归

Fisher线性判别

多分类

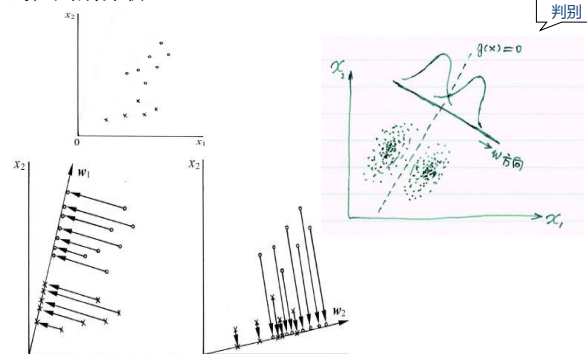
非平衡分类

线性判别分析 (Linear Discriminant Analysis)



由于将样例投影到一条直线（低维空间），因此也被视为一种“监督降维”技术 降维 → 第10章

线性判别分析



Fisher准则的基本原理：找到一个最合适的投影轴，使两类样本在该轴上投影之间的距离尽可能远，而每一类样本的投影尽可能紧凑，从而使分类效果为最佳。

21

d 维空间样本分布的描述量

- 各类样本均值向量 \mathbf{m}_i
$$\mathbf{m}_i = \frac{1}{N_i} \sum_{\mathbf{x} \in K_i} \mathbf{x} \quad i=1,2$$

- 样本类内离散度矩阵 \mathbf{S}_i 与总类内离散度矩阵 \mathbf{S}_w

$$\mathbf{S}_i = \sum_{\mathbf{x} \in K_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^T, \quad i=1,2 \quad \mathbf{S}_w = \mathbf{S}_1 + \mathbf{S}_2$$

- 样本类间离散度矩阵 \mathbf{S}_b :
$$\mathbf{S}_b = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$$

离散度矩阵在形式上与协方差矩阵很相似，但协方差矩阵是一种期望值，而离散矩阵只是表示有限个样本在空间分布的离散程度

22

一维 y 空间样本分布的描述量

- 各类样本均值
$$\tilde{m}_i = \frac{1}{N_i} \sum_{y \in V_i} y, \quad i=1,2$$

- 样本类内离散度和总类内离散度

$$\tilde{S}_i = \sum_{y \in V_i} (y - \tilde{m}_i)^2, \quad i=1,2 \quad \tilde{S}_w = \tilde{S}_1 + \tilde{S}_2$$

- 样本类间离散度

$$\tilde{S}_b = (\tilde{m}_1 - \tilde{m}_2)^2$$

以上定义描述 d 维空间样本点到一向量投影的分散情况，因此也就是对某向量 \mathbf{w} 的投影在 \mathbf{w} 上的分布。样本离散度的定义与随机变量方差相类似

23

样本与其投影统计量间的关系

- 样本 \mathbf{x} 与其投影 y 的统计量之间的关系：

$$\tilde{m}_i = \frac{1}{N_i} \sum_{y \in V_i} y = \frac{1}{N_i} \sum_{\mathbf{x} \in K_i} \mathbf{w}^T \mathbf{x} = \mathbf{w}^T \mathbf{m}_i, \quad i=1,2$$

$$\begin{aligned} \tilde{S}_b &= (\tilde{m}_1 - \tilde{m}_2)^2 = (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 \\ &= \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} = \mathbf{w}^T \mathbf{S}_b \mathbf{w} \end{aligned}$$

24

样本与其投影统计量间的关系

Fisher
判别

$$\begin{aligned}
 \tilde{S}_i &= \sum_{y \in W_i} (y - \tilde{m}_i)^2 \\
 &= \sum_{x \in K_i} (\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mathbf{m}_i)^2 \\
 &= \mathbf{w}^T \left[\sum_{x \in K_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^T \right] \mathbf{w} \\
 &= \mathbf{w}^T S_i \mathbf{w}
 \end{aligned}$$

$$\tilde{S}_1 + \tilde{S}_2 = \mathbf{w}^T (S_1 + S_2) \mathbf{w} = \mathbf{w}^T S_w \mathbf{w}$$

25

Fisher准则函数

Fisher
判别

- 评价投影方向 \mathbf{w} 的原则，使原样本向量在该方向上的投影能兼顾类间分布尽可能分开，类内尽可能密集的要求
- Fisher准则函数的定义：

$$J_F(\mathbf{w}) = \frac{\tilde{S}_b}{\tilde{S}_1 + \tilde{S}_2} = \frac{\mathbf{w}^T S_b \mathbf{w}}{\mathbf{w}^T S_w \mathbf{w}}$$

- ◆ Fisher最佳投影方向的求解

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} J_F(\mathbf{w})$$

26

Fisher最佳投影方向的求解

Fisher
判别

- 采用拉格朗日乘子算法解决

$$\mathbf{w}^* = S_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$

$\mathbf{m}_1, \mathbf{m}_2$ 是一向量，对与 $(\mathbf{m}_1 - \mathbf{m}_2)$ 平行的向量投影可使两均值点的距离最远。但是如从使类间分得较开，同时又使类内密集程度较高这样一个综合指标来看，则需根据两类样本的分布离散程度对投影方向作相应的调整，这就体现在对 $\mathbf{m}_1 - \mathbf{m}_2$ 向量按 S_w^{-1} 作一线性变换，从而使Fisher准则函数达到极值点

27

Fisher公式的推导

Fisher
判别

$$J_F(\mathbf{w}) = \frac{\tilde{S}_b}{\tilde{S}_1 + \tilde{S}_2} = \frac{\mathbf{w}^T S_b \mathbf{w}}{\mathbf{w}^T S_w \mathbf{w}}$$

$$\text{令 } \mathbf{w}^T S_w \mathbf{w} = c \neq 0$$

定义Lagrange函数: $L(\mathbf{w}, \lambda) = \mathbf{w}^T S_b \mathbf{w} - \lambda(\mathbf{w}^T S_w \mathbf{w} - c)$

$$\text{令: } \frac{\partial L(\mathbf{w}, \lambda)}{\partial \mathbf{w}} = S_b \mathbf{w} - \lambda S_w \mathbf{w} = 0$$

$$S_w^{-1} S_b \mathbf{w} = \lambda \mathbf{w}$$

$$\lambda \mathbf{w} = S_w^{-1} S_b \mathbf{w} = S_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} = S_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) R$$

$$\mathbf{w}^* = \frac{R}{\lambda} S_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) = S_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

28

决策规则: 若 $g(x) = \mathbf{w}^T \mathbf{x} + w_0 > 0$, 则 $x \in \omega_1$
 若 $g(x) = \mathbf{w}^T \mathbf{x} + w_0 < 0$, 则 $x \in \omega_2$

Fisher
判别

选阈值 w_0 :

(1) d 和 N 很大时, y 近似正态分布, 可在 Y 空间内用贝叶斯分类器。

(2) 经验, 如:

$$w_0 = -\frac{1}{2}(\tilde{m}_1 + \tilde{m}_2)$$

$$w_0 = -\tilde{m}$$

$$w_0 = -\frac{1}{2}(\tilde{m}_1 + \tilde{m}_2) - \frac{1}{N_1 + N_2 - 2} \ln \frac{P(\omega_1)}{P(\omega_2)}$$

2. 已知有两类数据, 分别为

Fisher
判别

$\omega_1: (1, 0), (2, 0), (1, 1)$

$\omega_2: (-1, 0), (0, 1), (-1, 1)$

试求: 该组数据的类内及类间离散矩阵 S_w 及 S_b 。

2. 解: 第一类的均值向量为

$$\mathbf{m}_1 = \left(\frac{4}{3}, \frac{1}{3} \right), \quad \mathbf{m}_2 = \left(-\frac{2}{3}, \frac{2}{3} \right)$$

$$\therefore S_1 = \frac{1}{9} \begin{pmatrix} 6 & -5 \\ -5 & 6 \end{pmatrix}, \quad S_2 = \frac{1}{3} \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$$

$$S_w = S_1 + S_2 = \frac{1}{9} \begin{pmatrix} 12 & -2 \\ -2 & 12 \end{pmatrix}$$

$$S_b = \begin{pmatrix} \frac{6}{3} \\ \frac{1}{3} \end{pmatrix} \begin{pmatrix} \frac{6}{3} & -\frac{1}{3} \end{pmatrix} = \frac{1}{9} \begin{pmatrix} 36 & -6 \\ -6 & 1 \end{pmatrix}$$

$$S_w^{-1} = \begin{bmatrix} 0.7714 & 0.1286 \\ 0.1286 & 0.7714 \end{bmatrix}$$

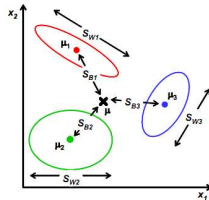
$$\mathbf{w}^* = S_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2) = [2.7407 \ 0.8889]^T$$

推广到多类

$$\mathbf{W} \in \mathbb{R}^{d \times (N-1)} \quad \mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{N-1}]$$

假定有 N 个类

- 全局散度矩阵 $\mathbf{S}_t = \mathbf{S}_b + \mathbf{S}_w = \sum_{i=1}^m (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$
- 类内散度矩阵 $\mathbf{S}_w = \sum_{i=1}^N \mathbf{S}_{w_i} \quad \mathbf{S}_{w_i} = \sum_{\mathbf{x} \in X_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T$
- 类间散度矩阵 $\mathbf{S}_b = \mathbf{S}_t - \mathbf{S}_w = \sum_{i=1}^N m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T$



推广到多类

$$\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{N-1}]$$

假定有 N 个类

- 全局散度矩阵 $\mathbf{S}_t = \mathbf{S}_b + \mathbf{S}_w = \sum_{i=1}^m (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$
 - 类内散度矩阵 $\mathbf{S}_w = \sum_{i=1}^N \mathbf{S}_{w_i} \quad \mathbf{S}_{w_i} = \sum_{\mathbf{x} \in X_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T$
 - 类间散度矩阵 $\mathbf{S}_b = \mathbf{S}_t - \mathbf{S}_w = \sum_{i=1}^N m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T$
- $$\mathbf{W} \in \mathbb{R}^{d \times (N-1)} \quad \mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{N-1}]$$

我们将样本点在这 $N-1$ 维向量投影后结果表示为: $\mathbf{y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{N-1}]$

$$\mathbf{y}_i = \mathbf{w}_i^T \mathbf{x}$$

$$\mathbf{y} = \mathbf{W}^T \mathbf{x}$$

推广到多类

假定有 N 个类

- 全局散度矩阵 $\mathbf{S}_t = \mathbf{S}_b + \mathbf{S}_w = \sum_{i=1}^m (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$
- 类内散度矩阵 $\mathbf{S}_w = \sum_{i=1}^N \mathbf{S}_{w_i} \quad \mathbf{S}_{w_i} = \sum_{\mathbf{x} \in X_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T$
- 类间散度矩阵 $\mathbf{S}_b = \mathbf{S}_t - \mathbf{S}_w = \sum_{i=1}^N m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T$

多分类LDA有多种实现方法: 采用 $\mathbf{S}_b, \mathbf{S}_w, \mathbf{S}_t$

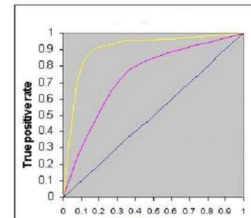
$$\text{例如, } \max_{\mathbf{W}} \frac{\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W})} \iff \mathbf{S}_b \mathbf{W} = \lambda \mathbf{S}_w \mathbf{W}$$

$$\mathbf{W} \in \mathbb{R}^{d \times (N-1)}$$

逻辑回归主要用来做回归吗?

逻辑回归中可以用以下哪种方法来训练数据?
A.最小二乘法 B.最大似然估计 C.杰卡德距离

10. 下图是3个逻辑回归模型的AUC-ROC曲线。不同的颜色表示不同的超参值, 哪一个会产生最佳的结果?



目录

线性归与逻辑回归

Fisher线性判别

多分类

非平衡分类

多类问题

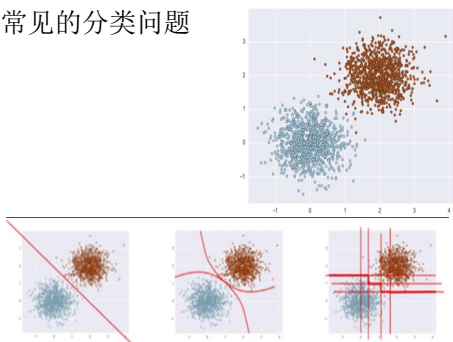
多分类

- 两类别问题可以推广到多类别问题
- 多个分类器
 - $\omega_i / \sim \omega_i$ 法: 将 C 类别问题化为 $(C-1)$ 个两类 (第 i 类与所有非 i 类) 问题, 按两类别问题确定其判别函数与决策面方程 (one-vs-rest)
 - ω_i / ω_j 法: 将 C 类中的每两类别单独设计其线性判别函数, 因此总共有 $C(C-1)/2$ 个线性判别函数方程 (one-vs-one)

非平衡分类问题

非平衡
分类

- 教课书中常见的分类问题



非平衡问题

非平衡
分类

- 研究不平衡类通常认为不平衡意味着少数类只占比10~20%。实际上，一些数据集远比这更不平衡。例如：
 - 每年大约有2%的信用卡账户被欺骗。（大多数欺诈检测领域严重不平衡。）
 - 状态医疗甄别通常在大量不存在此状态的人口中检测极少数有此状态的人（比如美国的HIV携带者仅占0.4%）。
 - 磁盘驱动器故障每年约1%。
 - 网络广告的转化率估计在 10^{-3} 到 10^{-6} 之间。

评价标准

非平衡
分类

- 混淆矩阵

状态与决策的可能关系		
决策 \ 状态	阳性	阴性
阳性	真阳性 (TP)	假阳性 (FP)
阴性	假阴性 (FN)	真阴性 (TN)

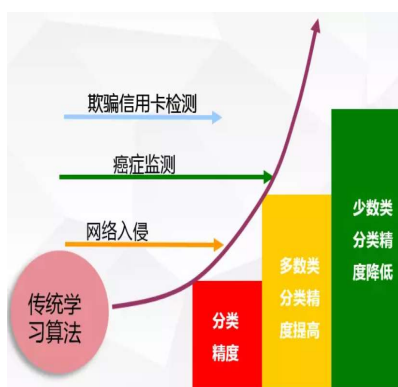
- 精度

$$\text{Precision} = \frac{TP}{TP + FP}$$

- 召回率

$$\text{Recall} = \frac{TP}{TP + FN}$$

- Kappa（衡量分类精度）和Roc曲线（曲线面积就是AUC）



类别不平衡 (class-imbalance)

非平衡
分类

不同类别的样本比例相差很大：“小类”往往更重要

基本思路：

若 $\frac{y}{1-y} > 1$ 则 预测为正例.

若 $\frac{y}{1-y} > \frac{m^-}{m^+}$ 则 预测为正例.

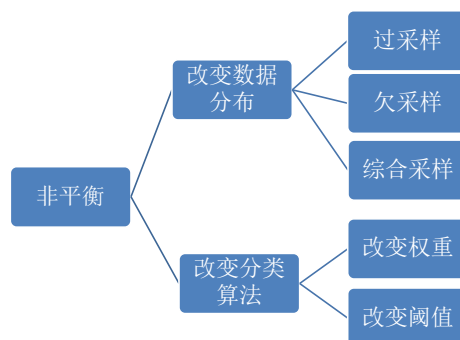
基本策略

——“再缩放” (rescaling):

$$\frac{y'}{1-y'} = \frac{y}{1-y} \times \frac{m^-}{m^+}$$

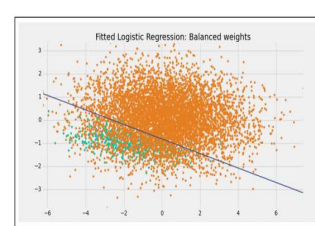
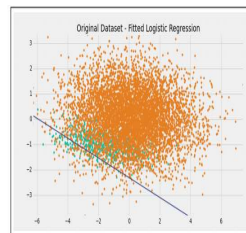
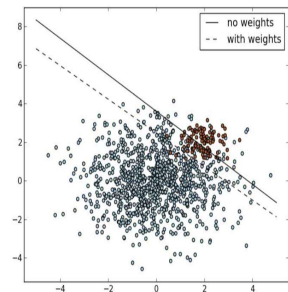
然而，精确估计 m^-/m^+ 通常很困难！

非平衡处理的方法

非平衡
分类

改变分类算法

- 一个处理非平衡数据常用的方法就是设置损失函数的权重，使得少数类别错误的损失大于多数类别错误的损失。在python的scikit-learn中我们可以使用class_weight参数来设置权重。
- 指定样本各类别的权重，主要是为了防止训练集某些类别的样本过多，导致训练的决策树过于偏向这些类别。这里可以自己指定各个样本的权重，或者用“balanced”，如果使用“balanced”，则算法会自己计算权重，样本量少的类别所对应的样本权重会高。当然，如果你的样本类别分布没有明显的偏倚，则可以不管这个参数，选择默认的“None”



The performance on the test set is as follows.

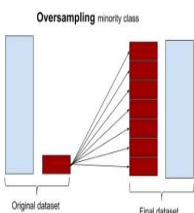
precision on L	recall on S
0.90	0.12

precision on L	recall on S
0.98	0.89

过采样

非平衡
分类

- 随机过采样
 - 采取简单复制样本的策略来增加少数类样本
 - 容易产生模型过拟合的问题
- SMOTE全称是Synthetic Minority Oversampling Technique即合成少数类过采样技术
 - 2002年Chawla提出了SMOTE算法
 - 对少数类样本进行分析并根据少数类样本人工合成新样本添加到数据集中



SMOTE算法流程

非平衡
分类

- SMOTE算法流程：
 - 对于少数类中每一个样本a，以欧氏距离为标准计算它到少数类样本集中所有样本的距离，得到其k近邻。
 - 根据样本不平衡比例设置一个采样比例以确定采样倍率N，对于每一个少数类样本a，从其k近邻中随机选择若干个样本，假设选择的近邻为b。
 - 对于每一个随机选出的近邻b，分别与原样本a按照如下的公式构建新的样本：

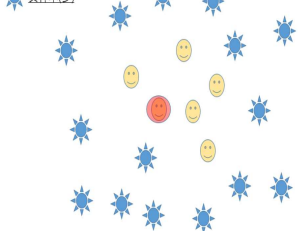
$$c = a + \text{rand}(0,1) * |a - b|$$

SMOTE 算法是建立在相距较近的少数类样本之间的样本仍然是少数类的假设基础上的。

步骤_1. 选一个正样本

😊 正样本(少)

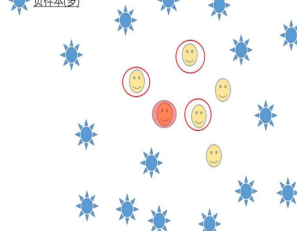
🌟 负样本(多)



步骤2: 找到该正样本的K个近邻 (假设K=3)

😊 正样本(少)

🌟 负样本(多)



步骤_3. 随机从k个近邻中选出一个样本



步骤_4. 在正样本和随机选出的这个近邻之间的连线上，随机找一点。这个点就是人工合成



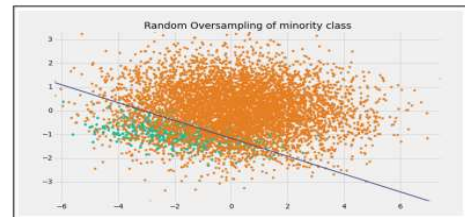
```

Algorithm SMOTE(T, N, k)
Input: Number of minority class samples T; Amount of SMOTE N%; Number of nearest
neighbors k
Output: (N/100) * T synthetic minority class samples
1. (* If N is less than 100%, randomize the minority class samples as only a random
percent of them will be SMOTEd *)
2. if N < 100
3. then Randomize the T minority class samples
4.   T = (N/100) * T
5.   N = 100
6. endif
7. N = (int)(N/100) (* The amount of SMOTE is assumed to be in integral multiples of
100. *)
8. k = Number of nearest neighbors
9. numattr = Number of attributes
10. Sample[ ] = array for original minority class samples
11. newindex: keeps a count of number of synthetic samples generated, initialized to 0
12. Synthetic[ ] = array for synthetic samples
(* Compute k nearest neighbors for each minority class sample only *)
13. for i = 1 to T
14.   Compute k nearest neighbors for i, and save the indices in the arrayray
15.   Populate(N, i, arrayray)
16. endwhile
Populate(N, i, arrayray) (* Function to generate the synthetic samples. *)
17. while N > 0
18.   Choose a random number between 1 and i, call it sn. This step chooses one of
the k nearest neighbors of i.
19.   for ally = 1 to numattr
20.     Compute: dif = Sample[arrayray[sn][ally]] - Sample[i][ally]
21.     Compute: gap = random number between 0 and 1
22.     Synthetic[newindex][ally] = Sample[i][ally] + gap * dif
23.   endwhile
24.   newindex++
25.   N = N - 1
26. endwhile
27. return (* End of Populate. *)
End of Pseudo-Code.

```

<http://blog.csdn.net/KeyCodey>

	L	S
Before resampling	6320	680
After resampling	6320	3160



precision on L	recall on S
0.97	0.76

Smote算法总结

非平衡
分类

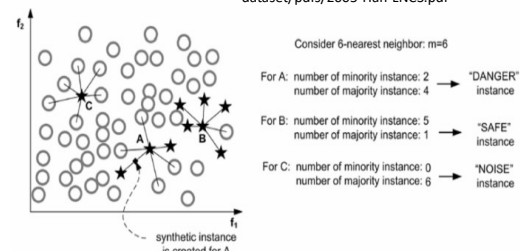
- 总结:
 - 对于少数类的每个样本寻找其同类样本中k个最近邻。其中, k通常是大于1的奇数。
 - 重复上述插值过程, 使得新生成的训练数据集数据达到均衡, 最后利用新的训练样本集进行训练
- 优点:
 - 有助于简单打破过抽样所产生的关系
 - 使得分类器的学习能力得到显著提高
- 缺陷:
 - 在近邻选择时, 存在一定的盲目性。
 - 边界模糊化

Borderline-SMOTE

非平衡
分类

- Han (2005) 等人在SMOTE算法基础上进行了改进, 提出了Borderline-SMOTE算法, 解决了生成样本重叠(Overlapping)的问题。算法在运行的过程中, 查找一个适当的区域, 该区域可以较好地反应数据集的性质, 然后在该区域内进行插值, 以使新增加的“人造”样本更有效。

<https://sci2s.ugr.es/keel/keel-dataset/pdfs/2005-Han-LNCS.pdf>



非平衡
分类

- Borderline-SMOTE

For each point p in S :

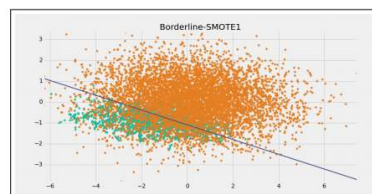
1. 计算点 p 在训练集 T 上的 m 个最近邻。我们称这个集合为 M_p ，然后设 $m' = |M_p \cap L|$ (表示点 p 的最近邻中属于 L 的数量)。
2. If $m' = m$, p 是一个噪声, 不做任何操作。
3. If $0 \leq m' \leq m/2$, 则说明 p 很安全, 不做任何操作。
4. If $m/2 \leq m' \leq m$, 那么点 p 就很危险了, 我们需要在这个点附近生成一些新的少数类点, 所以我们把它加入到DANGER中。

最后, 对于每个在DANGER中的点 d , 使用SMOTE算法生成新的样本。

非平衡
分类

- 应用Borderline-SMOTE的参数设置为 $k=5$

	$ L $	$ S $
Before resampling	6320	680
After resampling	6320	3160



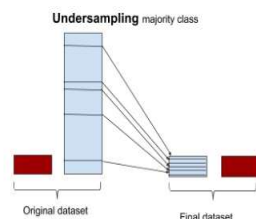
precision on L	recall on S
0.97	0.78

欠采样

非平衡
分类

- 随机欠采样

- 从多数类中随机抽取样本从而减少多数类样本的数量, 使数据达到平衡。

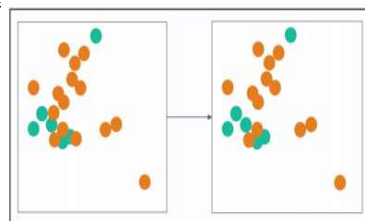


欠采样

非平衡
分类

- TomLink

- 如果有两个不同类别的样本, 它们的最近邻都是对方, 也就是A的最近邻是B, B的最近邻是A, 那么A,B就是 Tomek link。我们要做的就是讲所有 Tomek link 都删除掉。那么一个删除 Tomek link 的方法就是, 将组成 Tomek link 的两个样本, 如果有一个属于多数类样本, 就将该多数类样本删除掉。



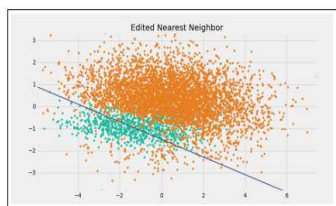
欠采样

非平衡
分类

- lomLink

- 针对minority的样本考察紧邻

	$ L $	$ S $
Before resampling	6320	680
After resampling	5120	680



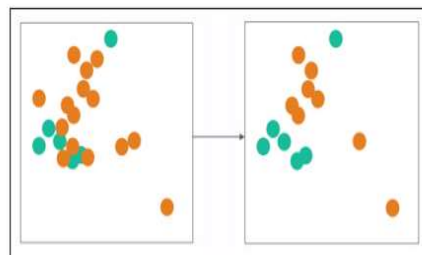
precision on L	recall on S
0.95	0.59

欠采样

非平衡
分类

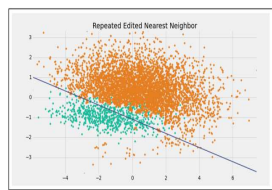
- ENN方法

- 对于那些多数类的样本, 如果他的大部分 k 近邻样本都跟他自己本身的类别不一样, 我们就将他删除。



非平衡
分类

- **Repeated Edited Nearest Neighbor**
- 这个方法就是不断的重复上述的删除过程，直到无法再删除为止。

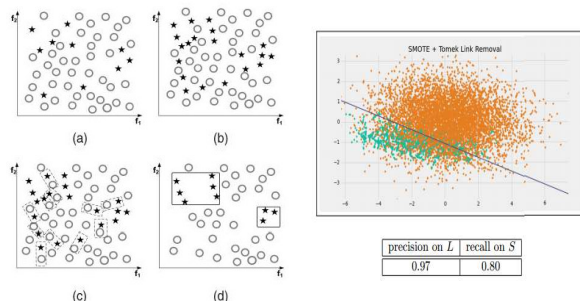


precision on L	recall on S
0.97	0.80

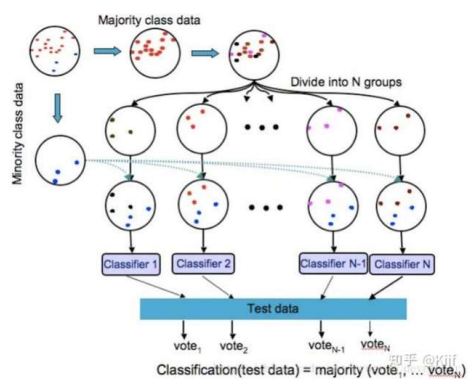
综合

非平衡
分类

- **SMOTE + Tomek**
 - 先合成SMOTE点，然后根据Tomek思想去除 Tomek links。



easyEnsemble

非平衡
分类

前往第五站.....

