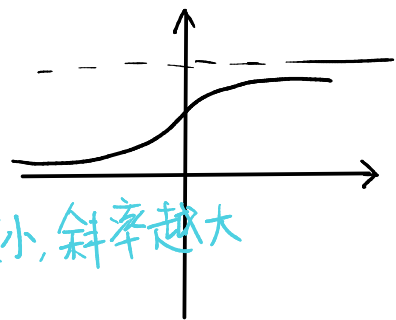


Logistic Regression 逻辑回归

1. Logistic Distribution

$$F(x) = \frac{1}{1 + \exp(-(x-\mu)/\sigma)}$$



• Output Y (class label)

Binary response $Y \in \{0, 1\}$

$$P(Y=1|X) = \frac{\exp(x^T \beta)}{1 + \exp(x^T \beta)}$$

$\Phi(x^T \beta)$ 正态分布 cdf

$$P(Y=0|X) = \frac{1}{1 + \exp(x^T \beta)}$$

$$\text{if } x^T \beta \rightarrow +\infty \quad P(Y=1|X) \rightarrow 1$$

$$\rightarrow -\infty \quad P(Y=1|X) \rightarrow 0$$

2. Model Estimation

$$L(\beta) = \prod_{i=1}^N p(x_i, \beta)^{y_i} \{1 - p(x_i, \beta)\}^{1-y_i}, \quad p(x_i, \beta) = \frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)}$$

$$l(\beta) = \sum_{i=1}^N (y_i x_i^T \beta - \log(1 + \exp(x_i^T \beta)))$$

$$\begin{aligned} \frac{\partial L}{\partial \beta} &= \sum_{i=1}^N y_i x_i - \frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)} \cdot x_i & E\left\{\frac{\partial L}{\partial \beta} \mid \beta = \beta_0\right\} &= \sum_{i=1}^n x_i E\{y_i - p(x_i, \beta)\} = 0 \\ &= \sum_{i=1}^N x_i (y_i - p(x_i, \beta)) \end{aligned}$$

• Optimization Newton-Raphson algorithm

$$\beta^{\text{new}} = \beta^{\text{old}} - \underbrace{\left(\frac{\partial^2 L}{\partial \beta \partial \beta^T}\right)^{-1}}_{\text{= Fisher information}} \frac{\partial L(\beta)}{\partial \beta} \quad \frac{\partial^2 L}{\partial \beta \partial \beta^T} = - \sum x_i x_i^T p(x_i, \beta) \{1 - p(x_i, \beta)\}$$

$$\text{Let } \sum x_i x_i^T = X^T X, \quad W = \text{diag}\{W_1, \dots, W_n\} \quad W_i = p(x_i, \beta) (1 - p(x_i, \beta))$$

$$\Rightarrow \frac{\partial^2 L}{\partial \beta \partial \beta^T} = X^T W X$$

3. Multi-nomial Logistic Regression model

$$Y = \{1, 2, \dots, K\} \quad P(Y=k | X=x) = \frac{\exp(\beta_k^T x)}{1 + \sum_{i=1}^K \exp(\beta_i^T x)}$$

4. Model Evaluation

1. 总体衡量

(1) Accuracy (精度) $\frac{TP+TN}{P+N}$

(2) Error Rate (错分率) $\frac{FP+FN}{P+N}$

		Predicted	
		Yes	No
Actual	Yes	TP	FN
	No	FP	TN

2. 查准率与查全率

(1) Precision (查准率) $= \frac{TP}{TP+FP}$

(2) Recall (查全率) $= \frac{TP}{TP+FN}$

(3) F_1 度量 (precision 和 recall 的调和平均) $= \frac{2 \times P \times R}{P+R}$

(4) F_β 度量: $F_\beta = \frac{(1+\beta^2) \times P \times R}{\beta^2 P + R}$

3. ROC 曲线 & AUC

评估方法不依赖于阈值

原理: 如果一个分类器能尽量多地将正例排在负例之前, 则认为它有较好的分类能力

ROC (Receiver Operation Characteristic 曲线)

横轴 假正例率 False positive rate $FPR = \frac{FP}{TN+FP}$

纵轴 真正例率 True positive rate $TPR = \frac{TP}{TP+FN}$

AUC (Area Under Curve) ROC 曲线下的面积

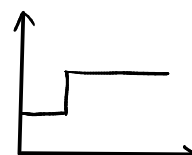
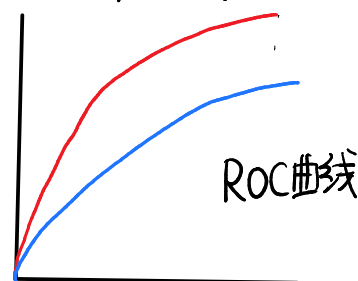
有限样本下绘制的 ROC 曲线

(1) 将阈值设成最大, 此时 $TPR=FPR=0$, 此时无预测的正例

(2) 将样本的预测值从高到低排序, 将阈值依次设成预测值, 即将每个样本标成一个点

(3) 设前一个点 (x, y) $\begin{cases} \text{当前预测真正例 } (x, y + \frac{1}{n}) \\ \text{当前预测假正例 } (x + m, y) \end{cases}$

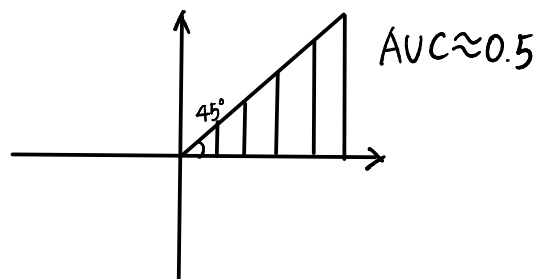
TPR 高 FPR 也高



(4) 假如一个分类器随机预测正例负例

AUC的取值只与排序有关

排序“损失”(rank loss)



假设有 m^+ 正例和 m^- 负例, 令 D^+ 与 D^- 分别代表正例与反例的集合

$$l_{\text{rank}} = \frac{1}{m^+ m^-} \sum_{x^+ \in D^+} \sum_{x^- \in D^-} [I(\underbrace{f(x^+) < f(x^-)}_{\text{模型预测值}}) + \frac{1}{2} I(f(x^+) = f(x^-))]$$

模型预测值

若正例预测值小于负例, 则记一个罚分, 若相等, 则记 0.5 罚分

$$AUC = 1 - l_{\text{rank}}$$

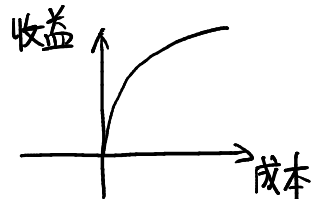
4. 成本-收益线

成本的质量(覆盖率)

$$\frac{TP + FP}{P + N} \quad (\text{预测的正值比例})$$

收益的质量(捕获率)

Recall (查全率)



5. 多次采样

由于五折交叉验证存在随机性, 一般重复 K 次实验取度量的平均值

6. 类别不均衡

(1) 设置域值 $\frac{p_i}{1-p_i} > \frac{m^+}{m^-}$ (不再是 $\frac{p_i}{1-p_i} \geq 1$)

(2) 过采样 (Over-sampling) SMOTE 算法

(3) 欠采样 (Under-sampling) Easy-Ensemble 算法

· 广义线性模型

(1) 指数分布族: $f(y|\theta, \psi) = \exp\left\{\frac{yb(\theta) - c(\theta)}{a(\psi)} + d(y, \psi)\right\}$

θ 典型参数: 与 y 的均值 μ 有关

ψ 刻度参数: 与 y 的方差 σ^2 有关