

数据集: 西瓜句

样本: 单个西瓜

特征向量:

属性.

得到模型: 分类 $\begin{cases} \text{二分类} \\ \text{多分类} \end{cases}$

回归: 价格预测:

聚类

进行预测:

奥卡剃刀: 这最简单的哪个.

模型评估与选择:

评估方法 $\begin{cases} \text{泛化能力} \\ \begin{cases} \text{训练集} \\ \text{测试集} \\ \text{验证集} \end{cases} \end{cases} \begin{cases} \text{留出法} \\ \text{交叉验证法} \\ \text{自助法} \end{cases}$

性能度量:

均方误差 $E(f; D) = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2$

引入概率密度 $p(x)$ $E(f; D) = \int_{x \in D} (f(x) - y)^2 p(x) dx$

错误率: $E(f; D) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(x_i) \neq y_i)$

精度: $acc(f; D) = 1 - E(f; D) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(x_i) = y_i)$

引入密度函数: $E(f; D) = \int_{x \in D} \mathbb{I}(f(x) \neq y) p(x) dx$

$acc(f; D) = \int_{x \in D} \mathbb{I}(f(x) = y) p(x) dx$

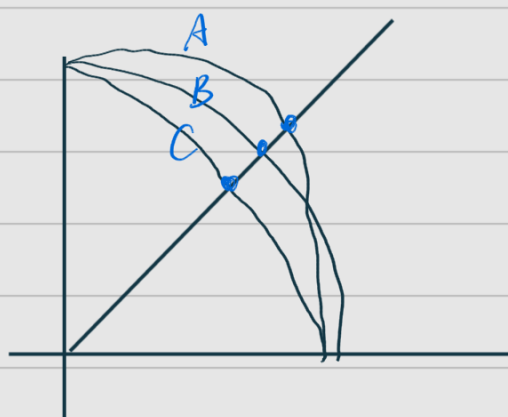
混淆矩阵:

真实	预测情况	
	正例.	反例.
正例	TP ①	FN ②
反例	FP ③	TN ④

查准率: $\frac{TP}{TP+FP}$

查全率: $\frac{TP}{TP+FN}$

P-R图.



$A > C$
 $B > C$
 比较 A-B.

F1 度量.

查准: $\frac{①}{①+③}$

查全: $\frac{①}{①+②}$

调和平均数.

$$2 \frac{①+③}{①} + 2 \frac{①+②}{①} = \frac{2①+②+③}{①}$$

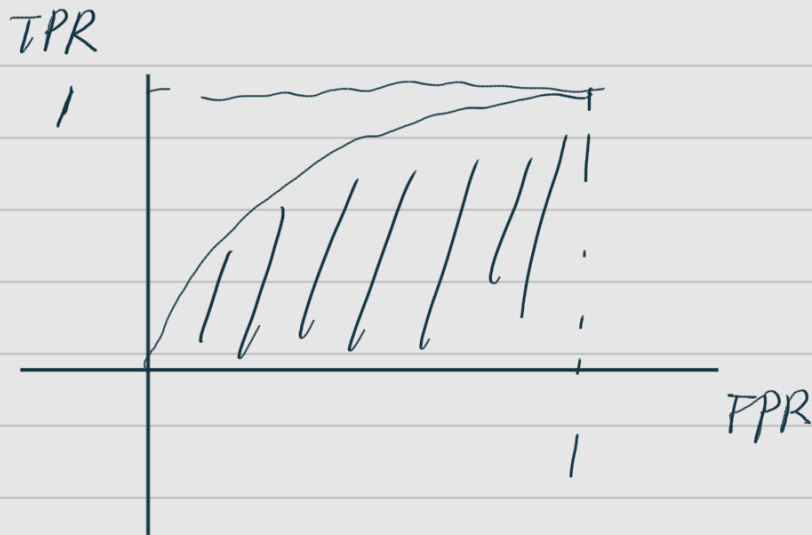
$$\frac{2①}{2①+②+③} = \frac{2 \times TP}{总 - TN + TP}$$

真实	预测情况	
	正例.	反例.
正例	TP ①	FN ②
反例	FP ③	TN ④

ROC & AUC

$$TPR = \frac{TP}{TP + FN} \quad (\text{True Positive Rate})$$

$$FPR = \frac{FP}{TP + FP} \quad (\text{False Positive Rate})$$



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

代价敏感错率, 代价曲线,
预测

真实	0	1
0	0	cost 1
1	cost 1	0

代价敏感错率:

$$E(f, D, w) = \frac{1}{m} \left(\sum_{x_i \in D} \mathbb{I}(f(x_i) \neq y_i) \cdot w(0) + \sum_{x_i \in D} \mathbb{I}(f(x_i) \neq y_i) \cdot w(1) \right)$$