

在现实任务中常会遇到这样的情况:不同类型的错误所造成的后果不同.例如在医疗诊断中，错误地把患者诊断为健康人与错误地把健康人诊断为患者，看起来都是犯了"一次错误"但后者的影响是增加了进一步检查的麻烦，前者的后果却可能是丧失了拯救生命的最佳时机;再如，门禁系统错误地把可通行人员拦在门外，将使得用户体验不佳，但错误地把陌生人放进门内，则会造成严重的安全事故.为权衡不同类型错误所造成的不同损失，可为错误赋予"非均等代价" (unequal cost). 以二分类任务为例，我们可根据任务的领域知识设定一个"代价矩阵" (cost matrix)，如下表，其中  $\text{cost}_{ij}$  表示将第  $i$  类样本预测为第  $j$  类样本的代价.一般来说， $\text{cost}_{ii}=0$ ;若将第0类判别为第1类所造成的损失更大，则  $\text{cost}_{01} > \text{cost}_{10}$ ; 损失程度相差越大， $\text{cost}_{01}$  与  $\text{cost}_{10}$  值的差别越大.

二分类代价矩阵		预测类别	
		第0类	第1类
真实类别	第0类	0	$\text{cost}_{01}$
	第1类	$\text{cost}_{10}$	0

回顾前面介绍的一些性能度量可看出，它们大都隐式地假设了均等代价，例如下式所定义的错误率是直接计算"错误次数"，并没有考虑不同错误会造成不同的后果.在非均等代价下，我们所希望的不再是简单地最小化错误次数，而是希望最小化"总体代价" (totalcost).若将上表中的第0类作为正类、第1类作为反类，令  $D^+$  与  $D^-$  分别代表样例集  $D$  的正例子集和反例子集，则"代价敏感"(cost-sensitive) 错误率为

$$E(f; D; \text{cost}) = \frac{1}{m} \left( \sum_{x_i \in D^+} \mathbf{1}(f(x_i) \neq y_i) \times \text{cost}_{01} + \sum_{z_i \in D^-} \mathbf{1}(f(z_i) \neq y_i) \times \text{cost}_{10} \right)$$

类似的，可给出基于分布定义的代价敏感错误率，以及其他一些性能度量如精度的代价敏感版本.若令  $\text{cost}_{ij}$  中的  $i, j$  取值不限于 0、1，则可定义出多分类任务的代价敏感性度量.

在非均等代价下，ROC曲线不能直接反映出学习器的期望总体代价，而"代价曲线" (cost curve) 则可达到该目的.代价曲线图的横轴是取值为  $[0, 1]$  的正例概率代价

$$P(+|\text{cost}) = \frac{p \times \text{cost}_{01}}{p \times \text{cost}_{01} + (1 - p) \times \text{cost}_{10}}$$

其中  $p$  是样例为正例的概率;纵轴是取值为  $[0, 1]$  的归一化代价

$$\text{cost}_{\text{norm}} = \frac{\text{FNR} \times p \times \text{cost}_{01} + \text{FPR} \times (1 - p) \times \text{cost}_{10}}{p \times \text{cost}_{01} + (1 - p) \times \text{cost}_{10}}$$

其中FPR(  $\text{FPR} = \frac{FP}{TN+FP}$  )是定义的假正例率，FNR=1-TPR(  $\text{TPR} = \frac{TP}{TP+FN}$  )是假反例率.代价曲线的绘制很简单:ROC由线上每...点对应了代价平面上的二条线段7设ROC曲线上点的坐标为(TPR,FPR)，则可相应计算出FNR，然后在代价平面上绘制一条从(O，FPR)到(1，FNR)的线段，线段下的面积即表示了该条件下的期望总体代价;如此将 ROC 曲线上的每个点转化为代价平面上的一条线段，然后取所有线段的下界，围成的自积即为在所有条件下学习器的期望总体代价，如下图.

