

README

问题

安德森鸢尾花卉数据集(*Anderson's Iris dataset*), 其中包含150个样本, 对应数据集的每行数据。每行数据包含每个样本的4个特征和样本的1个类别信息。每个样本包含了花萼长度、花萼宽度、花瓣长度、花瓣宽度4个特征 (前4列)、1个品种信息, 即目标属性 (第5列, 也叫*target*)。

请将测试集与训练集按照1 : 4划分, 建立一个分类器, 分类器可以通过样本的4个特征来进行样本的分类, 判断样本属于山鸢尾、变色鸢尾还是维吉尼亚鸢尾 (三个品种名称 : 分别对应0、1、2)中的哪种。

数据特征内容

	数据特征
1	<i>sepal length(cm)</i> : 花萼长度
2	<i>sepal width(cm)</i> : 花萼宽度
3	<i>petal length(cm)</i> : 花瓣长度
4	<i>petal width(cm)</i> : 花瓣宽度
5	<i>target</i> : 品种信息 0, 1, 2

关键指标

通过一个混淆矩阵阐述指标:

混淆矩阵

实际 \ 预测	预测为好的样本数	预测为坏的样本数
实际为“好”的样本数	<i>TP</i>	<i>FN</i>
实际为“坏”的样本数	<i>FP</i>	<i>TN</i>

1、查准率: $P = \frac{TP}{TP+FP}$, 基于预测数据, 看实际“好”样本的占比

2、查全率: $R = \frac{TP}{TP+FN}$, 基于实际数据, 看实际“好”样本的占比

但我们能够发现, 查准率与查全率是一对矛盾的度量 (如: 查全率越高查准率就会越低)

3、想同时关注查准率和查全率, 可考虑对两个指标进行加权, 这里我们引入 F_1 值:

$$F_1 = \frac{2RP}{P + R}$$

此时 R, P 的权重各取 $\frac{1}{2}$, 将二者同等看待。

注：数据来源为Creator:R.A. Fisher, Donor:Michael Marshall (MARSHALL%PLU '@' <http://io.arc.nasa.gov>)