

数据挖掘作业二：分类数据集 letter-recognition

利用 2000 个训练数据，构造决策树分类器，预测 2000 个测试数据的类别，并给出测试集正确率的估计值及依据。

Relevant Information:

The objective is to identify each of a large number of black-and-white rectangular pixel displays as one of the 26 capital letters in the English alphabet. The character images were based on 20 different fonts and each letter within these 20 fonts was randomly distorted to produce a file of 20,000 unique stimuli. Each stimulus was converted into 16 primitive numerical attributes (statistical moments and edge counts) which were then scaled to fit into a range of integer values from 0 through 15.



Number of Attributes: 17 (Letter category and 16 numeric features)

No. training data: 4000

No. test data: 16000

Attribute Information:

1. lettr capital letter (26 values from A to Z)
2. x-box horizontal position of box (integer)
3. y-box vertical position of box (integer)
4. width width of box (integer)
5. high height of box (integer)
6. onpix total # on pixels (integer)
7. x-bar mean x of on pixels in box (integer)
8. y-bar mean y of on pixels in box (integer)
9. x2bar mean x variance (integer)
10. y2bar mean y variance (integer)
11. xybar mean x y correlation (integer)
12. x2ybr mean of $x * x * y$ (integer)
13. xy2br mean of $x * y * y$ (integer)
14. x-egge mean edge count left to right (integer)
15. xegvy correlation of x-egge with y (integer)

- 16. y-ege mean edge count bottom to top (integer)
- 17. yegvx correlation of y-ege with x (integer)