

python 高级应用 决策树

2024 年 5 月 14 日

实验目标

以树为模型的分类器

编程语言：Python3

使用库：Numpy

分类任务：对个体的收入进行判断 $>50K$ or $\leq 50K$

数据集：Adult (UCI Machine Learning Repository: Adult Data Set)

数据预处理

设计思路

使用 `np.genfromtxt` 函数读入训练数据，缺失项用众数填补。

分出 `X,y`

对 `y` 进行二值编码 `y = y == y[0]`

对 `X` 的非数值列进行数值编码。

测试数据同样按照上述方式操作，编码使用到的映射与训练数据一致。

读入数据

```
1 # 训练集
2 data = np.genfromtxt("adult/adult.data", delimiter
    = ",", dtype=str)
3 data = data[~np.any(data == '?', axis=1)]
4 X = data[:, :14]
5 y = data[:, 14]
```

编码

```
1 # 01 编码
2 y = y == y[0]
3
4 # 数值编码
5 for i in [1, 3, 5, 6, 7, 8, 9, 13]:
6     values = np.unique(X[:, i])
7     encoding = {value: idx for idx, value in
8                 enumerate(values)}
9     for j in range(X.shape[0]):
10         X[j][i] = encoding[X[j][i]]
11
12 X = X.astype(int)
13 y = y.astype(int)
```

分割数据集

```
1 X_train = X[:X.shape[0] // 3 * 2, :]  
2 y_train = y[:X.shape[0] // 3 * 2]  
3  
4 X_test = X[X.shape[0] // 3 * 2:, :]  
5 y_test = y[X.shape[0] // 3 * 2:]
```

结点设置

```
1 class TreeNode:
2     def __init__(
3         self, feature_index=None, threshold=
4             None,
5             value=None, left=None, right=None
6     ):
7         self.feature_index = feature_index
8         self.threshold = threshold
9         self.value = value
10        self.left = left
11        self.right = right
```

算法思路

采用递归建树的思路

```
1 function generate(X: array, y: array, depth: int)
   → Node:
2     找到X的最佳划分D
3     按照D划分为 $X_l$ ,  $X_r$ ,  $y_l$ ,  $y_r$ 
4     把划分数据交给左右子树, 递归建树
5
6     终止条件: y只剩下一类或达到最大递归深度
7     返回值: 记录有划分标准D与左右孩子L,R的结点
```


参考代码

```
1 def generateTree(self, X, y, depth=1):
2     if depth == self.max_depth or len(set(y)) == 1:
3         node = TreeNode(value=np.sort(y)[y.shape[0] >> 1])
4     else:
5         best_feature_index, best_threshold = find_best_split(X, y)
6         left_indices = X[:, best_feature_index] <= best_threshold
7         right_indices = ~left_indices
8
9         left = self.generateTree(X[left_indices], y[left_indices], depth + 1)
10        right = self.generateTree(X[right_indices], y[right_indices], depth + 1)
11
12        node = TreeNode(
13            feature_index=best_feature_index, threshold=best_threshold,
14            left=left, right=right
15        )
16    return node
```

最佳划分

```
1 def split_data(X, y):
2     best_feature_index = None
3     best_threshold = None
4     best_entropy = inf
5
6     for feature in X.features:
7         D = set(X[feature])
8         for v in D:
9             left = i where X[i][feature] <= v
10            right = i where X[i][feature] > v
11
12            entropy = left / y.shape[0] * entropy(y[left]) + right / y.shape[0] * entropy(y[right])
13
14            if entropy < best_entropy:
15                best_entropy = entropy
16                best_feature = feature
17                best_threshold = v
18
19     return best_feature, best_threshold
```