



南开大学
Nankai University

南 开 大 学

计 算 机 学 院

机器学习实验报告

实验四 决策树分类器

姓名：王泳鑫

学号：1911479

年级：2019级

专业：计算机科学与技术

指导教师：卫金茂

2021 年 11 月 28 日

摘要

关键字：层次聚类，Machine Learning，Deep Learning

目录

| | |
|--------------------------|---|
| 一、 实验描述与要求 | 1 |
| 二、 程序设计与代码实现 | 1 |
| (一) ID3算法决策树构建 | 1 |
| (二) CART决策树构建 | 4 |
| 三、 测试与结果演示 | 5 |
| (一) ID3决策树 | 5 |
| (二) CART决策树 | 5 |

一、 实验描述与要求

基本要求：绘制聚类前后样本分布情况

1. 基于 Watermelon-train1数据集（只有离散属性），构造ID3决策树；
2. 基于构造的 ID3 决策树，对数据集 Watermelon-test1进行预测，输出分类精度；

中级要求：

1. 对数据集Watermelon-train2，构造C4.5或者CART决策树，要求可以处理连续型属性；
2. 对测试集Watermelon-test2进行预测，输出分类精度；

二、 程序设计与代码实现

（一） ID3算法决策树构建

首先，ID3算法的核心思想应用信息增益准则作为标准,我们根据公式就可以实现计算信息熵，代码如下：

```
1 # 计算信息熵
2 def calculate_entropy(y):
3     log2 = math.log2
4     unique_labels = np.unique(y)
5     entropy = 0
6     for label in unique_labels:
7         count = len(y[y == label])
8         p = count / len(y)
9         entropy += -p * log2(p)
10    return entropy
```

接下来，我们需要定义树的节点，代码如下：

```
1 class DecisionNode():
2     def __init__(self, feature_i=None, threshold=None,
3                  value=None, true_branch=None, false_branch=None):
4         self.feature_i = feature_i
5         self.threshold = threshold
6         self.value = value
7         self.true_branch = true_branch
8         self.false_branch = false_branch
```

然后要对特征进行划分，代码如下：

```
1 def divide_on_feature(X, feature_i, threshold):
2     split_func = None
3     if isinstance(threshold, int) or isinstance(threshold, float):
4         split_func = lambda sample: sample[feature_i] >= threshold
5     else:
6         split_func = lambda sample: sample[feature_i] == threshold
7
8     X_1 = np.array([sample for sample in X if split_func(sample)])
9     X_2 = np.array([sample for sample in X if not split_func(sample)])
10
11    return np.array([X_1, X_2])
```

接着，就是我们决策树的内容，代码如下：

```

1  class DecisionTree(object):
2  def __init__(self, min_samples_split=2, min_impurity=1e-7,
3              max_depth=float("inf"), loss=None):
4      self.root = None #根节点
5      self.min_samples_split = min_samples_split
6      self.min_impurity = min_impurity
7      self.max_depth = max_depth
8      # 计算值如果是分类问题就是信息增益，回归问题就基尼指数
9      self._impurity_calculation = None
10     self._leaf_value_calculation = None #计算叶子
11     self.one_dim = None
12     self.loss = loss
13
14     def fit(self, X, y, loss=None):
15         self.one_dim = len(np.shape(y)) == 1
16         self.root = self._build_tree(X, y)
17         self.loss=None
18
19     def _build_tree(self, X, y, current_depth=0):
20         """递归求解树
21
22         """
23
24         largest_impurity = 0
25         best_criteria = None
26         best_sets = None
27
28         if len(np.shape(y)) == 1:
29             y = np.expand_dims(y, axis=1)
30
31         Xy = np.concatenate((X, y), axis=1)
32
33         n_samples, n_features = np.shape(X)
34
35         if n_samples >= self.min_samples_split and current_depth <= self.max_depth:
36             # 计算每一个特征的增益值
37             for feature_i in range(n_features):
38                 feature_values = np.expand_dims(X[:, feature_i], axis=1)
39                 unique_values = np.unique(feature_values)
40
41                 for threshold in unique_values:
42                     Xy1, Xy2 = divide_on_feature(Xy, feature_i, threshold)
43
44                     if len(Xy1) > 0 and len(Xy2) > 0:
45                         y1 = Xy1[:, n_features:]
46                         y2 = Xy2[:, n_features:]
47
48                         # 计算增益值
49                         impurity = self._impurity_calculation(y, y1, y2)
50
51                         if impurity > largest_impurity:
52                             largest_impurity = impurity
53                             best_criteria = {"feature_i": feature_i, "threshold":
                                     threshold}

```

```

54         best_sets = {
55             "leftX": Xy1[:, :n_features],
56             "lefty": Xy1[:, n_features:],
57             "rightX": Xy2[:, :n_features],
58             "righty": Xy2[:, n_features:]
59         }
60
61     if largest_impurity > self.min_impurity:
62         true_branch = self._build_tree(best_sets["leftX"], best_sets["lefty"],
63                                         current_depth + 1)
64         false_branch = self._build_tree(best_sets["rightX"], best_sets["righty"],
65                                         current_depth + 1)
66
67     return DecisionNode(feature_i=best_criteria["feature_i"], threshold=
68                         best_criteria[
69                             "threshold"], true_branch=true_branch, false_branch=false_branch)
70
71     # 计算节点的目标值
72     leaf_value = self._leaf_value_calculation(y)
73
74     return DecisionNode(value=leaf_value)
75
76 def predict_value(self, x, tree=None):
77     """预测
78
79     """
80
81     if tree is None:
82         tree = self.root
83
84     if tree.value is not None:
85         return tree.value
86
87     feature_value = x[tree.feature_i]
88
89     branch = tree.false_branch
90     if isinstance(feature_value, int) or isinstance(feature_value, float):
91         if feature_value >= tree.threshold:
92             branch = tree.true_branch
93     elif feature_value == tree.threshold:
94         branch = tree.true_branch
95
96     return self.predict_value(x, branch)
97
98 def predict(self, X):
99     y_pred = []
100     for x in X:
101         y_pred.append(self.predict_value(x))
102     return y_pred

```

最后是分类器的实现，代码如下：

```

1 class ClassificationTree(DecisionTree):
2     def _calculate_information_gain(self, y, y1, y2):
3         # 计算信息增益
4         p = len(y1) / len(y)

```

```

5         entropy = calculate_entropy(y)
6         info_gain = entropy - p * calculate_entropy(y1) - (1 - p) * calculate_entropy(
            y2)
7
8         return info_gain
9
10    def _majority_vote(self, y):
11        most_common = None
12        max_count = 0
13        for label in np.unique(y):
14            # 投票决定当前的节点为哪一个类
15            count = len(y[y == label])
16            if count > max_count:
17                most_common = label
18                max_count = count
19
20        return most_common
21
22    def fit(self, X, y):
23        self._impurity_calculation = self._calculate_information_gain
24        self._leaf_value_calculation = self._majority_vote
25        super(ClassificationTree, self).fit(X, y)

```

(二) CART决策树构建

对于CART决策树，我们只需要把ID3算法中计算信息熵函数修改为计算基尼系数的函数，代码如下：

```

1    def calculate_variance(X):
2        """ Return the variance of the features in dataset X """
3        mean = np.ones(np.shape(X)) * X.mean(0)
4        n_samples = np.shape(X)[0]
5        variance = (1 / n_samples) * np.diag((X - mean).T.dot(X - mean))
6
7        return variance

```

同时把分类器中计算函数修改为：

```

1    def _calculate_variance_reduction(self, y, y1, y2):
2        var_tot = calculate_variance(y)
3        var_1 = calculate_variance(y1)
4        var_2 = calculate_variance(y2)
5        frac_1 = len(y1) / len(y)
6        frac_2 = len(y2) / len(y)
7
8        # 使用方差缩减
9        variance_reduction = var_tot - (frac_1 * var_1 + frac_2 * var_2)
10
11    return sum(variance_reduction)

```

三、 测试与结果演示

(一) ID3决策树

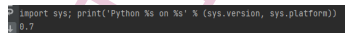
通过测试数据集来测试，代码如下：

```

1 from sklearn.metrics import accuracy_score
2
3 x_train = np.loadtxt('D:\\Data\\lab6\\Watermelon-train1.csv', dtype='str', delimiter=',',
4                       , skiprows=1, usecols=(1, 2, 3, 4))
5
6 y_train = np.loadtxt('D:\\Data\\lab6\\Watermelon-train1.csv', dtype='str', delimiter=',',
7                       , skiprows=1, usecols=(5))
8
9 clf = ClassificationTree()
10 clf.fit(x_train, y_train)
11
12
13 y_pred = clf.predict(x_test)
14
15 accuracy = accuracy_score(y_pred, y_test)
16
17 print(accuracy)

```

结果如图2所示



```

P import sys; print("Python %s on %s" % (sys.version, sys.platform))
0.7

```

图 1: ID3决策树结果

(二) CART决策树

这里需要注意的是我们需要把密度从字符串形式转换为float类型，代码如下：

```

1 from sklearn.metrics import accuracy_score
2
3 x_train = np.loadtxt('D:\\Data\\lab6\\Watermelon-train2.csv', dtype='str', delimiter=',',
4                       , skiprows=1, usecols=(1, 2, 3, 4, 5))
5
6 y_train = np.loadtxt('D:\\Data\\lab6\\Watermelon-train2.csv', dtype='str', delimiter=',',
7                       , skiprows=1, usecols=(6))
8
9
10 for i in range(17):
11     x_train[i][4] = float(x_train[i][4])
12
13 x_test = np.loadtxt('D:\\Data\\lab6\\Watermelon-test2.csv', dtype='str', delimiter=',',
14                     , skiprows=1, usecols=(1, 2, 3, 4, 5))
15
16 y_test = np.loadtxt('D:\\Data\\lab6\\Watermelon-test2.csv', dtype='str', delimiter=',',
17                     , skiprows=1, usecols=(6))
18
19
20 for i in range(5):
21     x_test[i][4] = float(x_train[i][4])
22
23

```

```
15 clf = RegressionTree()  
16 clf.fit(x_train,y_train)  
17  
18  
19 y_pred = clf.predict(x_test)  
20  
21 accuracy = accuracy_score(y_pred,y_test)  
22  
23 print (accuracy)
```

结果如图2所示



```
return np.array([4, 1, 4, 2])  
0.7  
1.0
```

图 2: CART决策树结果