

《数据仓库与数据挖掘》实验报告

姓名:	朱志儒	学号:	SA20225085	日期:	2020/12/14
上机题目:	决策树 ID3 算法				
操作环境: OS: Window 10 CPU: AMD Ryzen 5 3600X 6-Core Processor 4.25GHz GPU: GeForce RTX 2070 super					
<p>一、基础知识:</p> <p>决策树</p> <p>决策树是一种树形结构,可以是二叉树或非二叉树,树中每个非叶节点表示一个特征属性上的测试,每个分支代表这个特征属性在某个值域上的输出,每个叶节点存放一个类别。使用决策树进行分类的过程就是从根节点开始,测试待分类项中相应的特征属性,按照其值选择输出分支,直到到达叶子节点,将叶子节点存放的类别作为分类结果。</p> <p>ID3</p> <p>ID3 模型是以信息熵和信息增益作为衡量标准的分类模型。</p> <p>熵是指信息的混乱程度,熵值越大,变量的不确定性也就越大,计算信息熵的公式:</p> $\text{Entropy}(S) = - \sum_{i=1}^m p(u_i) \log_2(p(u_i))$ <p>其中, $p(u_i) = \frac{ u_i }{ S }$, $p(u_i)$ 为类别 u_i 在样本 S 中出现的概率。</p> <p>条件熵是指在已知第二个随机变量 X 的值的条件下,随机变量 Y 的信息熵。计算特征 A 对数据集 S 的条件熵的公式:</p> $H(S A) = \sum_{V \in \text{Value}(A)} \frac{ S_V }{ S } \text{Entropy}(S_V)$ <p>其中, A 表示样本特征, $\text{Value}(A)$ 是特征 A 所有的取值集合, V 是 A 中一个特征值, S_V 是 S 中 A 的值为 V 的样例集合。</p> <p>信息增益是指在某个条件下,信息复杂度,即不确定性,减少的程度。计算信息增益的公式:</p> $\text{infoGain}(S,A) = \text{Entropy}(S) - H(S A)$ <p>其中, A 表示样本特征。</p>					

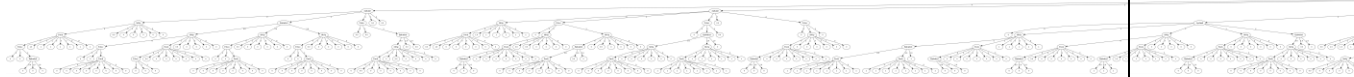
在构建决策树时，选择信息增益最大的特征作为决策点。

二、实验过程：

根据上次实验的方法处理训练集数据，依据这些数据构建决策树，由于整个数据的特征项并不是很多，所以构建好决策树后没有进行剪枝操作。最后，根据构建好的决策树对测试集数据进行分类，分类的过程就是从根节点进行搜索直到叶节点，叶节点的标签就是该数据的分类标签。

三、结果分析：

构建决策树，如图所示（原图见附件）



对测试集数据进行分类，准确率如下：

ID3准确率：0.7583732057416268

附录

算法源代码（C/C++/JAVA 描述）：

```
1. def read_data_set():
2.     '''处理数据，提取特征'''
3.     trainData = pd.read_csv("train.csv")
4.     testData = pd.read_csv("test.csv")
5.     # 将训练集和测试集整合
6.     data = pd.concat([trainData, testData], axis=0).reset_index(drop=True)
7.     # male: 0, female: 1
8.     data['Sex'].replace(['male', 'female'], [0, 1], inplace=True)
9.     # S: 0, C: 1, Q: 2
10.    data['Embarked'].replace(['S', 'C', 'Q'], [0, 1, 2], inplace=True)
11.
12.    # print(data[data['Fare'].isnull()])
13.    # Pclass: 3, Embarked: 0, Sex: 0
14.    # 填补 Fare 为 NaN 的数据
15.    data['Fare'] = data['Fare'].fillna(
16.        np.mean(data[((data['Pclass'] == 3) & (data['Embarked'] == 0) & (data['Sex'] == 0))['Fare']]))
17.    # print(data[data['Fare'].isnull()])
18.    # Empty DataFrame
19.
20.    # data['FareLimit'] = pd.qcut(data['Fare'], 4)
21.    # print(data.groupby(['FareLimit'])['Survived'].mean())
22.    # 使用 FareLimit 替代 Fare
```

```

23.     data['FareLimit'] = 0
24.     data.loc[data['Fare'] <= 8.662, 'FareLimit'] = 0
25.     data.loc[(data['Fare'] > 8.662) & (data['Fare'] <= 14.4
54), 'FareLimit'] = 1
26.     data.loc[(data['Fare'] > 14.454) & (data['Fare'] <= 53.
1), 'FareLimit'] = 2
27.     data.loc[data['Fare'] > 53.1, 'FareLimit'] = 3
28.
29.     # print(data[data['Embarked'].isnull()])
30.     # Pclass: 1, Sex: 1
31.     # 填补 Embarked 为 NaN 的数据
32.     data['Embarked'] = data['Embarked'].fillna(
33.         stats.mode(data[((data['Pclass'] == 1) & (data['Sex
'] == 1))]['Embarked'])[0][0])
34.     # print(data[data['Embarked'].isnull()])
35.     # Empty DataFrame
36.
37.     # 添加新特征: 家人 Family
38.     data['Family'] = data['SibSp'] + data['Parch']
39.
40.     # 填补 Age 为 NaN 的数据
41.     dataAgeNaNIndex = data[data['Age'].isnull()].index
42.     for i in dataAgeNaNIndex:
43.         # 取 Pclass、Family 相同的数据的平均值
44.         meanAge = data['Age'][
45.             (data['Pclass'] == data.iloc[i]['Pclass']) & (d
ata['Family'] == data.iloc[i]['Family'])].mean()
46.         data['Age'].iloc[i] = meanAge
47.
48.     # data['AgeLimit'] = pd.cut(data['Age'], 5)
49.     # print(data.groupby(['AgeLimit'])['Survived'].mean())

50.     # 使用 AgeLimit 替代 Age
51.     data['AgeLimit'] = 0
52.     data.loc[data['Age'] <= 16, 'AgeLimit'] = 0
53.     data.loc[(data['Age'] > 16) & (data['Age'] <= 32), 'Age
Limit'] = 1
54.     data.loc[(data['Age'] > 32) & (data['Age'] <= 48), 'Age
Limit'] = 2
55.     data.loc[(data['Age'] > 48) & (data['Age'] <= 60), 'Age
Limit'] = 3
56.     data.loc[data['Age'] > 60, 'AgeLimit'] = 4
57.
58.     # 删除无用列

```

```

59.     data.drop(labels=["Age", "Fare", "Ticket", "Cabin", "Name", "PassengerId", 'Family'], axis=1, inplace=True)
60.
61.     data = data[['Pclass', 'Sex', 'SibSp', 'Parch', 'Embarked', 'FareLimit', 'AgeLimit', 'Survived']]
62.
63.     trainX = data[:len(trainData)].values
64.     testY = data[len(trainData):]['Survived'].values.tolist()
65.     testX = data[len(trainData):].drop(labels='Survived', axis=1).values
66.
67.     return trainX, testX, testY
68.
69.
70. def empirical_entropy(train_set, D):
71.     '''根据数据集 train_set 计算类别 D 的信息熵'''
72.     dict_of_kinds = {}
73.     for i in range(len(train_set)):
74.         label = train_set[i][D]
75.         if label not in dict_of_kinds.keys():
76.             dict_of_kinds[label] = 1
77.         else:
78.             dict_of_kinds[label] += 1
79.     summ = len(train_set)
80.     for key in dict_of_kinds.keys():
81.         pd = dict_of_kinds[key] / summ
82.         dict_of_kinds[key] = pd * math.log(pd)
83.     return -sum(list(dict_of_kinds.values()))
84.
85.
86. def condition_entropy(train_set, A, D):
87.     '''根据数据集 train_set, 在已知条件 A 的前提下, 计算 D 的条件熵'''
88.     dict_of_kinds = {}
89.     for i in range(len(train_set)):
90.         label = train_set[i][A]
91.         if label not in dict_of_kinds.keys():
92.             dict_of_kinds[label] = [i]
93.         else:
94.             dict_of_kinds[label].append(i)
95.     for key in dict_of_kinds.keys():
96.         dict_of_acct = {}
97.         for j in dict_of_kinds[key]:

```

```

98.         label = train_set[j][D]
99.         if label not in dict_of_acct.keys():
100.             dict_of_acct[label] = 1
101.         else:
102.             dict_of_acct[label] += 1
103.         summ = len(dict_of_kinds[key])
104.         for keyy in dict_of_acct.keys():
105.             pd = dict_of_acct[keyy] / summ
106.             dict_of_acct[keyy] = pd * math.log(pd)
107.             dict_of_kinds[key] = len(dict_of_kinds[key]) / len(
                train_set) * (-sum(list(dict_of_acct.values()))))
108.         return sum(list(dict_of_kinds.values()))
109.
110.
111. def informatin_gain(train_set, D, A):
112.     '''计算信息增益'''
113.     return empirical_entropy(train_set, D) - condition_entropy(
        train_set, A, D)
114.
115.
116. def get_child_set(items):
117.     '''求 items 的非空真子集'''
118.     result = [[]]
119.     for x in items:
120.         result.extend([subset + [x] for subset in result]
            )
121.     return result[1:len(result) - 1]
122.
123.
124. class decision_node:
125.     '''定义节点类'''
126.
127.     def __init__(self, col, value=None, child_node=None):
128.
129.         self.col = col
130.         self.value = value
131.         self.child_node = child_node
132.
133. class Decision_tree:
134.     '''定义决策树类'''
135.
136.     def __init__(self, train_set, D, dict_of_labels, function, dict_of_col):

```

```

137.         self.labels = [i for i in range(len(dict_of_labels.keys()))]
138.         self.train_set = train_set
139.         self.D = D
140.         self.dict_of_labels = dict_of_labels
141.         self.function = function
142.         self.dict_of_col = dict_of_col
143.         self.decision_tree = self.build_tree(train_set, [
            ], function)
144.
145.         def is_same(self, remain_set):
146.             '''数据集 D 中的样本属于同一类别 C，则将当前结点标记为
                C 类叶结点'''
147.             acct = remain_set[0][self.D]
148.             for row in remain_set:
149.                 if acct != row[self.D]:
150.                     return False
151.             return True
152.
153.         def classify(self, test_set):
154.             '''根据已生成的决策树将 test_set 中的数据分类'''
155.             results = []
156.             for row in test_set:
157.                 head = self.decision_tree
158.                 while head.col != -1:
159.                     for key in head.child_node.keys():
160.                         if key == row[head.col]:
161.                             head = head.child_node[key]
162.                             break
163.                 results.append(head.value)
164.             return results
165.
166.         def find_mode(self, remain_set):
167.             '''找出 remain_set 中出现次数最多的类别'''
168.             result = [0, 0]
169.             for line in remain_set:
170.                 result[int(line[self.D])] += 1
171.             if result[0] > result[1]:
172.                 return 0
173.             else:
174.                 return 1
175.
176.         def build_tree(self, remain_set, used_col, function):

```

```

177.         '''构建决策树'''
178.         if self.is_same(remain_set):
179.             '''数据集 D 中的样本属于同一类别 C, 则将当前结点标
                记为 C 类叶结点'''
180.             return decision_node(-1, value=remain_set[0][
                self.D])
181.         if len(used_col) == len(self.dict_of_labels.keys(
                )):
182.             '''特征集 A 为空集, 或数据集 D 中所有样本在 A 中所
                有特征上取值相同, 此时无法划分。将当前结点标记为叶结点, 类别为 D
                中出现最多的类'''
183.             return decision_node(-1, value=self.find_mode
                (remain_set))
184.
185.         '''构建 ID3 模型决策树'''
186.         entropies = []
187.         for i in self.labels:
188.             if i not in used_col:
189.                 entropies.append(function(remain_set, sel
                f.D, i))
190.             else:
191.                 entropies.append(-1)
192.             '''选择信息增益最大的特征作为决策点'''
193.             choose = entropies.index(max(entropies))
194.             new_used_col = used_col + [choose]
195.             child_node = {}
196.             for label in self.dict_of_labels[choose]:
197.                 new_remain_set = []
198.                 for row in remain_set:
199.                     if row[choose] == label:
200.                         new_remain_set.append(row)
201.                 if len(new_remain_set) == 0:
202.                     '''数据集 D 为空集, 则将当前结点标记为叶结点,
                        类别为父结点中出现最多的类'''
203.                     child_node[label] = decision_node(-1, val
                        ue=self.find_mode(remain_set))
204.                 else:
205.                     child_node[label] = self.build_tree(new_r
                        emain_set, new_used_col, function)
206.             return decision_node(choose, child_node=child_nod
                e)
207.
208.         def visit_node(self, graph, father_col, node, strings
                ):

```

```

209.         # 绘制决策树时访问节点函数
210.         if node.col != -1:
211.             new_col = father_col + self.dict_of_col[node.col] + str(strings)
212.             graph.node(new_col, self.dict_of_col[node.col])
213.             graph.edge(father_col, new_col, str(strings))
214.         for key in node.child_node.keys():
215.             self.visit_node(graph, new_col, node.child_node[key], key)
216.         else:
217.             new_col = father_col + str(node.value) + str(strings)
218.             graph.node(new_col, str(node.value))
219.             graph.edge(father_col, new_col, str(strings))
220.
221.     def draw_tree(self):
222.         # 绘制决策树的图像
223.         tree_name = 'ID3_Decision_Tree.gv'
224.         graph = Digraph(tree_name, format='png')
225.         node = self.decision_tree
226.         graph.node(str(node.col), self.dict_of_col[node.col])
227.         for key in node.child_node.keys():
228.             self.visit_node(graph, str(node.col), node.child_node[key], key)
229.         graph.render(tree_name, view=True)
230.
231.     def validation(testY, result, function_name):
232.         count = 0
233.         for i in range(len(result)):
234.             if int(result[i]) == int(testY[i]):
235.                 count += 1
236.         print(function_name + '准确率:', count / len(result))
237.
238.
239. if __name__ == '__main__':
240.     dict_of_col = {0: 'Pclass', 1: 'Sex', 2: 'SibSp', 3: 'Parch', 4: 'Embarked', 5: 'FareLimit', 6: 'AgeLimit'}
241.     dict_of_train_labels = {0: [1, 2, 3], 1: [0, 1], 2: [0, 1, 2, 3, 4, 5, 8], 3: [0, 1, 2, 3, 4, 5, 6, 9],

```


	<pre>242. 4: [0.0, 1.0, 2.0], 5: [0, 1, 2, 3], 6: [0, 1, 2, 3, 4]} 243. train_set, test_set, testY = read_data_set() 244. dt_ig = Decision_tree(train_set, 7, dict_of_train_labels, informatin_gain, dict_of_col) 245. dt_ig.draw_tree() 246. re1 = dt_ig.classify(test_set) 247. validation(testY, re1, 'ID3')</pre>
--	---