

重 庆 大 学

学 生 实 验 报 告

实验课程名称 人工智能导论

开课实验室 DS1502

学 院 软件学院 年级 2021 专业班 软工 X 班

学 生 姓 名 XXX 学 号 2021XXXX

开 课 时 间 2023 至 2024 学年第 1 学期

总 成 绩	
教师签名	

大数据与软件学院制

《人工智能导论》实验报告

开课实验室：DS1502

2023 年 12 月 9 日

学院	大数据与软件学院	年级、专业、班	21 软件工程 X 班	姓名	XXX	成绩	
课程名称	人工智能导论	实验项目名称	基于决策树的企鹅分类	指导教师		XX	
教师评语	<div>教师签名：2023 年 月 日</div>						
<div>一、实验目的</div> <p>实验目的是通过决策树对企鹅进行分类，并评估模型的性能，以了解模型在解决企鹅分类问题上的有效性。</p> <div>二、实验内容</div> <div>① 数据准备：使用了一个包含企鹅相关特征的数据集，例如岛屿、喙长度、喙深度、鳍长度、体重、性别和年龄。</div> <div>② 决策树构建：利用数据集和标签，以及标签属性（是分类属性还是连续属性）创建了一个决策树。</div> <div>③ 数据集划分：将数据集分为训练集和测试集，其中训练集用于构建决策树，测试集用于评估模型的性能。</div> <div>④ 模型训练：使用训练集对决策树进行训练。</div> <div>⑤ 模型测试：使用测试集进行分类预测，评估模型在新数据上的性能。</div> <div>⑥ 性能评估：计算分类准确率，并生成混淆矩阵以更详细地了解模型的性能。</div> <div>⑦ 可视化：绘制混淆矩阵的热力图，以及决策树的结构图，以便更直观地理解模型的工作方式。</div>							

三、使用仪器、材料

1. 操作系统: Windows 11
2. 开发设备: Lenovo Legion R9000P2021H
3. 开发平台: PyCharm 2023.1

四、实验过程原始记录(数据、图表、计算等):

(一) 源代码

```
1  from math import log
2  import pandas as pd
3  import numpy as np
4  import operator
5  import seaborn as sns
6  from sklearn.metrics import confusion_matrix
7  from sklearn.model_selection import train_test_split
8  from sklearn.metrics import accuracy_score
9  import matplotlib.pyplot as plt
10 import copy
11
12 # 获取决策树叶子节点数量
13 3 用法
14 def get_num_leafs(tree):
15     num_leafs = 0
16     first_str = next(iter(tree))
17     second_dict = tree[first_str]
18     for key in second_dict.keys():
19         if type(second_dict[key]) is dict:
20             num_leafs += get_num_leafs(second_dict[key])
21         else:
22             num_leafs += 1
23     return num_leafs
24
25 # 获取决策树深度
26 2 用法
27 def get_tree_depth(tree):
28     max_depth = 0
29     first_str = next(iter(tree))
30     second_dict = tree[first_str]
31     for key in second_dict.keys():
32         if type(second_dict[key]) is dict:
33             this_depth = 1 + get_tree_depth(second_dict[key])
34         else:
35             this_depth = 1
36         if this_depth > max_depth:
37             max_depth = this_depth
38     return max_depth
39
40 # 绘制决策树节点
41 2 用法
42 def plot_node(node_txt, center_pt, parent_pt, node_type):
43     arrow_args = dict(arrowstyle="<-")
44     create_plot.ax1.annotate(
45         node_txt,
46         xy=parent_pt,
47         xycoords='axes fraction',
48         xytext=center_pt,
49         textcoords='axes fraction',
50         va="center",
51         ha="center",
52         bbox=node_type,
53         arrowprops=arrow_args
54     )
```

```

53 # 在父节点和子节点之间绘制文本
54 2 用法
55 def plot_mid_text(cntr_pt, parent_pt, txt_string):
56     x_mid = (parent_pt[0] - cntr_pt[0]) / 2.0 + cntr_pt[0]
57     y_mid = (parent_pt[1] - cntr_pt[1]) / 2.0 + cntr_pt[1]
58     create_plot.ax1.text(x_mid, y_mid, txt_string, va="center", ha="center", rotation=30)
59
60 # 绘制决策树
61 23 用法
62 def plot_tree(tree, parent_pt, node_txt):
63     decision_node = dict(boxstyle="sawtooth", fc="0.8")
64     leaf_node = dict(boxstyle="round4", fc="0.8")
65     num_leafs = get_num_leafs(tree)
66     first_str = next(iter(tree))
67     cntr_pt = (plot_tree.x_off + (1.0 + float(num_leafs)) / 2.0 / plot_tree.total_w,
68               plot_tree.y_off)
69     plot_mid_text(cntr_pt, parent_pt, node_txt)
70     plot_node(first_str, cntr_pt, parent_pt, decision_node)
71     second_dict = tree[first_str]
72     plot_tree.y_off = plot_tree.y_off - 1.0 / plot_tree.total_d
73     for key in second_dict.keys():
74         if type(second_dict[key]).__name__ == 'dict':
75             plot_tree(second_dict[key], cntr_pt, str(key))
76         else:
77             plot_tree.x_off = plot_tree.x_off + 1.0 / plot_tree.total_w
78             plot_node(second_dict[key], (plot_tree.x_off, plot_tree.y_off), cntr_pt, leaf_node)
79             plot_mid_text((plot_tree.x_off, plot_tree.y_off), cntr_pt, str(key))
80     plot_tree.y_off = plot_tree.y_off + 1.0 / plot_tree.total_d
81
82 # 创建决策树绘图
83 4 用法
84 def create_plot(tree):
85     fig = plt.figure(1, facecolor='white')
86     fig.clf()
87     axprops = dict(xticks=[], yticks=[])
88     create_plot.ax1 = plt.subplot(111, frameon=False, **axprops)
89     plot_tree.total_w = float(get_num_leafs(tree))
90     plot_tree.total_d = float(get_tree_depth(tree))
91     plot_tree.x_off = -0.5 / plot_tree.total_w
92     plot_tree.y_off = 1.0
93     plot_tree(tree, (0.5, 1.0), '')
94     plt.show()

```

```

93 # 将字符串标签转换为数值
94 3 用法
95 def transition(x):
96     if x == data['Species'].unique()[0]:
97         return 0
98     if x == data['Species'].unique()[1]:
99         return 1
100     if x == data['Species'].unique()[2]:
101         return 2
102     if x == data['Island'].unique()[0]:
103         return 0
104     if x == data['Island'].unique()[1]:
105         return 1
106     if x == data['Island'].unique()[2]:
107         return 2
108     if x == data['Sex'].unique()[0]:
109         return 0
110     if x == data['Sex'].unique()[1]:
111         return 1
112     if x == data['Sex'].unique()[2]:
113         return -1.0

```

```

114 # 创建数据集
115 # 1 个用法
116 def create_data_set():
117     global data
118     data = data[[
119         'Island',
120         'Culmen Length (mm)',
121         'Culmen Depth (mm)',
122         'Flipper Length (mm)',
123         'Body Mass (g)',
124         'Sex',
125         'Age',
126         'Species',
127     ]]
128     data = data.fillna(-1)
129     data['Species'] = data['Species'].apply(transition)
130     data['Island'] = data['Island'].apply(transition)
131     data['Sex'] = data['Sex'].apply(transition)
132     data_set = []
133     for i in range(344):
134         data_set.append(list(data.iloc[i, :]))
135     labels = [
136         'Island', 'Culmen Length (mm)', 'Culmen Depth (mm)',
137         'Flipper Length (mm)', 'Body Mass (g)', 'Sex', 'Age'
138     ]
139     return data_set, labels
140
141 # 计算信息熵
142 # 4 用法
143 def calc_ent(data_set):
144     num_entries = len(data_set)
145     label_counts = {}
146     for feat_vec in data_set:
147         current_label = feat_vec[-1]
148         label_counts[current_label] = label_counts.get(current_label, 0) + 1
149     info_ent = 0.0
150     for key in label_counts:
151         prob = float(label_counts[key]) / num_entries
152         info_ent -= prob * log(prob, 2)
153     return info_ent
154
155 # 根据特征和特征值划分数据集
156 # 2 用法
157 def split_data_set(data_set, axis, value):
158     ret_data_set = []
159     for feat_vec in data_set:
160         if feat_vec[axis] == value:
161             reduced_feat_vec = feat_vec[:axis]
162             reduced_feat_vec.extend(feat_vec[axis + 1:])
163             ret_data_set.append(reduced_feat_vec)
164     return ret_data_set

```

```

163 # 根据数值型特征和划分值划分数据集
164 # 4 用法
165 def split_data_set_c(data_set, axis, value, lor_r='L'):
166     ret_data_set = []
167     if lor_r == 'L':
168         for feat_vec in data_set:
169             if float(feat_vec[axis]) < value:
170                 ret_data_set.append(feat_vec)
171     else:
172         for feat_vec in data_set:
173             if float(feat_vec[axis]) > value:
174                 ret_data_set.append(feat_vec)
175     return ret_data_set

```

```

176 # 选择最佳划分特征
177 1 个用法
178 def choose_best_feature_to_split(data_set, label_property):
179     num_features = len(label_property)
180     base_entropy = calc_ent(data_set)
181     best_info_gain = 0.0
182     best_feature = -1
183     best_part_value = None
184     for i in range(num_features):
185         feat_list = [example[i] for example in data_set]
186         unique_vals = set(feat_list)
187         new_entropy = 0.0
188         best_part_value_i = None
189         if label_property[i] == 0:
190             for value in unique_vals:
191                 sub_data_set = split_data_set(data_set, i, value)
192                 prob = len(sub_data_set) / float(len(data_set))
193                 new_entropy += prob * calc_ent(sub_data_set)
194         else:
195             sorted_unique_vals = list(unique_vals)
196             sorted_unique_vals.sort()
197             min_entropy = float('inf')
198             for j in range(len(sorted_unique_vals) - 1):
199                 part_value = (float(sorted_unique_vals[j]) +
200                             float(sorted_unique_vals[j + 1])) / 2
201                 data_set_left = split_data_set_c(data_set, i, part_value, 'L')
202                 data_set_right = split_data_set_c(data_set, i, part_value, 'R')
203                 prob_left = len(data_set_left) / float(len(data_set))
204                 prob_right = len(data_set_right) / float(len(data_set))
205                 entropy = prob_left * calc_ent(data_set_left) + prob_right * calc_ent(data_set_right)
206                 if entropy < min_entropy:
207                     min_entropy = entropy
208                     best_part_value_i = part_value
209             new_entropy = min_entropy
210             info_gain = base_entropy - new_entropy
211             if info_gain > best_info_gain:
212                 best_info_gain = info_gain
213                 best_feature = i
214                 best_part_value = best_part_value_i
215     return best_feature, best_part_value
216
217 # 统计类别标签出现频率，返回出现频率最高的标签
218 2 用法
219 def majority_count(class_list):
220     class_count = {}
221     for vote in class_list:
222         class_count[vote] = class_count.get(vote, 0) + 1
223     sorted_class_count = sorted(class_count.items(),
224                                key=operator.itemgetter(1),
225                                reverse=True)
226     return sorted_class_count[0][0]

```



```

226 # 递归构建决策树
227 4 用法
228 def create_tree(data_set, labels, label_property):
229     class_list = [example[-1] for example in data_set]
230     if class_list.count(class_list[0]) == len(class_list):
231         return class_list[0]
232     if len(data_set[0]) == 1:
233         return majority_count(class_list)
234     best_feat, best_part_value = choose_best_feature_to_split(
235         data_set, label_property)
236     if best_feat == -1:
237         return majority_count(class_list)
238     if label_property[best_feat] == 0:
239         best_feat_label = labels[best_feat]
240         my_tree = {best_feat_label: {}}
241         labels_new = copy.copy(labels)
242         label_property_new = copy.copy(label_property)
243         del labels_new[best_feat]
244         del label_property_new[best_feat]
245         feat_values = [example[best_feat] for example in data_set]
246         unique_value = set(feat_values)
247         for value in unique_value:
248             sub_labels = labels_new[:]
249             sub_label_property = label_property_new[:]
250             my_tree[best_feat_label][value] = create_tree(
251                 split_data_set(data_set, best_feat, value), sub_labels,
252                 sub_label_property)
253     else:
254         best_feat_label = labels[best_feat] + '<' + str(best_part_value)
255         my_tree = {best_feat_label: {}}
256         sub_labels = labels[:]
257         sub_label_property = label_property[:]
258         value_left = 'Yes'
259         my_tree[best_feat_label][value_left] = create_tree(
260             split_data_set_c(data_set, best_feat, best_part_value, 'L'), sub_labels,
261             sub_label_property)
262         value_right = 'No'
263         my_tree[best_feat_label][value_right] = create_tree(
264             split_data_set_c(data_set, best_feat, best_part_value, 'R'), sub_labels,
265             sub_label_property)
266     return my_tree
267
268
269

```

```

270 # 对测试样本进行分类
271 4 用法
272 def classify(input_tree, feat_labels, feat_label_properties, test_vec):
273     first_str = list(input_tree.keys())[0]
274     first_label = first_str
275     less_index = first_str.find('<')
276     if less_index > -1:
277         first_label = first_str[:less_index]
278     second_dict = input_tree[first_label]
279     feat_index = feat_labels.index(first_label)
280     class_label = None
281     for key in second_dict.keys():
282         if feat_label_properties[feat_index] == 0:
283             if test_vec[feat_index] == key:
284                 if type(second_dict[key]).__name__ == 'dict':
285                     class_label = classify(second_dict[key], feat_labels,
286                                             feat_label_properties, test_vec)
287                 else:
288                     class_label = second_dict[key]
289             else:
290                 part_value = float(str(first_str)[less_index + 1:])
291                 if test_vec[feat_index] < part_value:
292                     if type(second_dict['Yes']).__name__ == 'dict':
293                         class_label = classify(second_dict['Yes'], feat_labels,
294                                                 feat_label_properties, test_vec)
295                     else:
296                         class_label = second_dict['Yes']
297                 else:
298                     if type(second_dict['No']).__name__ == 'dict':
299                         class_label = classify(second_dict['No'], feat_labels,
300                                                 feat_label_properties, test_vec)
301                     else:
302                         class_label = second_dict['No']
303     return class_label
304

```

```

304 # 读取数据集
305 data = pd.read_csv('penguins_data.csv')
306
307 # 特征名称
308 feature_name = [
309     'Island', 'Culmen Length (mm)', 'Culmen Depth (mm)', 'Flipper Length (mm)', 'Body Mass (g)', 'Sex', 'Age'
310 ]
311
312 # 创建数据集
313 data_set, labels = create_data_set()
314
315 # 定义一个标签属性列表，表示每个属性是分类属性 (1) 还是连续属性 (0)
316 label_properties = [0, 1, 1, 1, 1, 0, 1]
317
318 # 使用数据集、标签和标签属性创建决策树
319 my_tree = create_tree(data_set, labels, label_properties)
320
321 # 从数据集中选择特征列
322 feature = data[[
323     'Island', 'Culmen Length (mm)', 'Culmen Depth (mm)', 'Flipper Length (mm)', 'Body Mass (g)', 'Sex', 'Age'
324 ]]
325
326 # 设置目标变量为物种类别
327 goal = data['Species']
328
329 # 划分数据集为训练集和测试集
330 x_train, x_test, y_train, y_test = train_test_split(feature, goal, test_size=0.2)
331
332 # 将测试集的特征转换为NumPy数组，并再次转换为列表
333 x_test = np.array(x_test)
334 x_test = x_test.tolist()
335
336 # 用训练好的决策树进行测试集的分类预测
337 test_pre = []
338 for i in range(len(x_test)):
339     result = classify(my_tree, labels, label_properties, x_test[i])
340     test_pre.append(result)
341
342 # 打印分类准确率
343 print("预测准确率为:", accuracy_score(y_test, test_pre))
344
345 # 生成混淆矩阵并用热力图进行可视化
346 confu_matrix = confusion_matrix(test_pre, y_test)
347 plt.figure(figsize=(8, 6))
348 sns.heatmap(confu_matrix, annot=True, cmap='Blues')
349 plt.xlabel('Predicted labels')
350 plt.ylabel('True labels')
351 plt.show()
352
353 # 绘制决策树的图形
354 create_plot(my_tree)
355
356 # 打印决策树结构
357 print(my_tree)

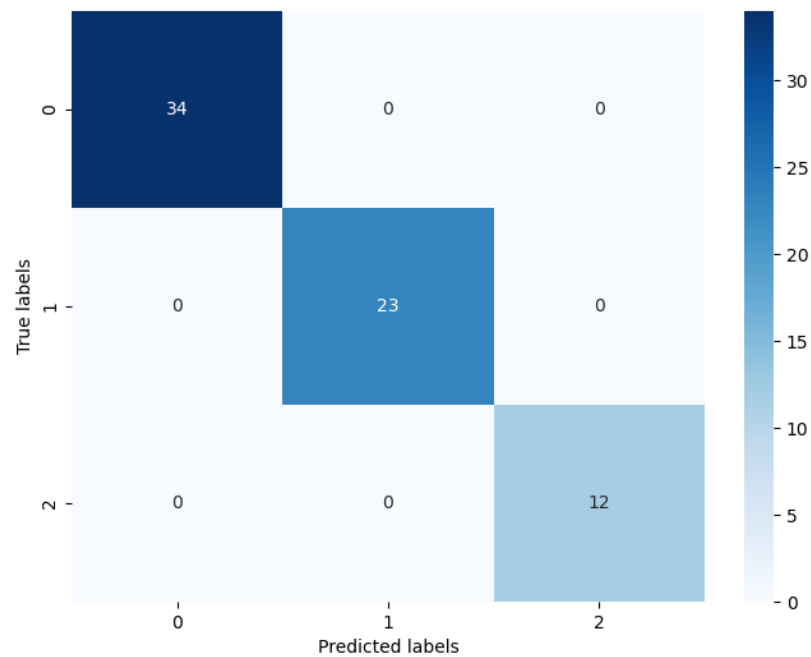
```


(二) 实现效果

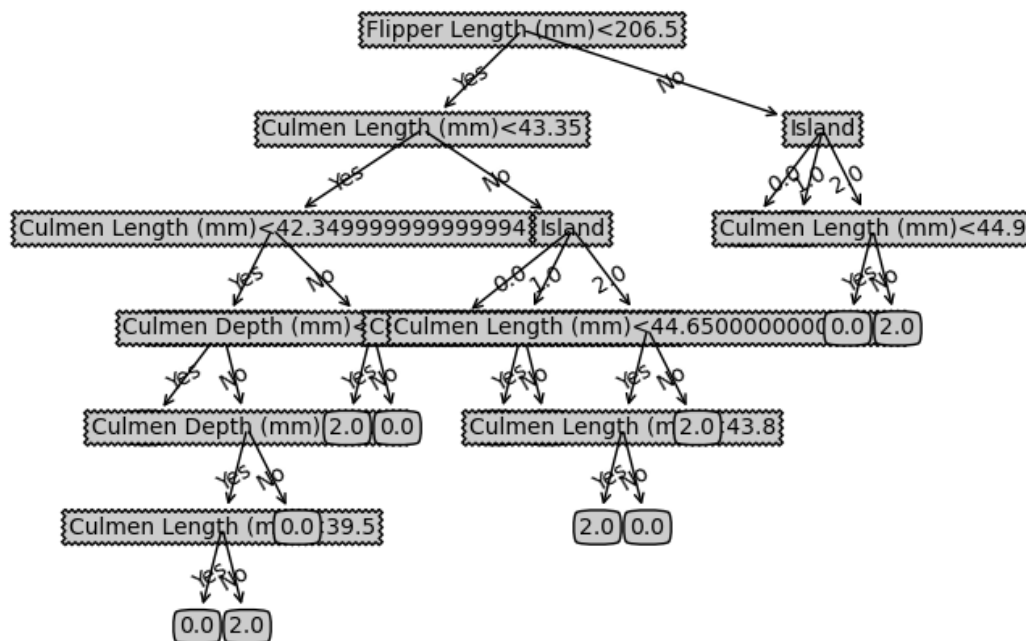
1. 预测准确率

预测准确率为: 1.0

2. 混淆矩阵



3. 决策树



4. 总结

本实验实现了一个基于决策树的企鹅分类模型，主要包括以下关键步骤：

（1）决策树构建和可视化：

① 利用信息熵和信息增益选择最佳特征进行数据集划分，支持处理连续型特征。

② 通过递归方式绘制决策树的节点，使用矩形框表示决策节点和叶子节点，并通过箭头表示决策流向。

（2）数据处理： 将原始数据集中的字符串标签转换为数值，处理缺失值。

（3）模型训练和测试： 将数据集划分为训练集和测试集，使用训练集训练决策树模型。利用测试集评估模型的分类准确率。

（4）混淆矩阵可视化：生成混淆矩阵，通过热力图进行可视化，直观展示模型在不同类别上的性能表现。

（5）决策树结构展示：打印构建好的决策树结构，以更详细地了解模型的构建。

（6）代码组织：通过函数进行模块化设计，提高了代码的可读性和可维护性。

总体而言，这个实验通过决策树对企鹅进行分类，并通过可视化和混淆矩阵的方式对模型进行了评估，展示了决策树在分类问题上的应用和性能评估过程。