

به نام خدا

ایمان علیپور

۹۸۱۰۲۰۲۴

تمرین پیاده سازی ۲ هوش مصنوعی

استاد: دکتر هجران دوست

هدف:

پیاده سازی decision tree.

چالش ها:

اصلی ترین چالشی که با آن روبرو شدم این بود که زمانی که کدی که برای رستوران زده بودم را به دیتاست دوم که برای دیابت بود انتقال دادم، الگوریتمم اجرایش تمام نمیشد و تلاش میکرد تا میشود عمق را زیاد کند، برای همین برای آن مجبور شدم محدودیت عمق قرار دهم که خروجی تولید شود.

چالش دیگر من چیزی بود که نفهمیدم چگونه ایجاد میشود، بعضی وقت ها که تابع محاسبه گر دقت را روی داده های train صدا میزد، اروری میگرفتم که علت آن را نمیدانم اما احتمالاً بخاطر overflow بوده و آنرا پیدا نکردم. این تابع برای داده های test هیچوقت خطا نداد!

نحوه اجرا:

فایل های csv را در کنار فایل های اصلی کد قرار دهید و فایل کد را اجرا کنید. توجه کنید برای هر دو درخت من یک تابع نوشتم که entropy را نیز چاپ کند، اما چون درخت ناخوانا میشد در نهایت از آن استفاده نکردم.

خروجی:

خروجی اینگونه است که هریار در درخت پایین میرویم یک indentation ایجاد میشود. برای هر دو کد یک تابع نوشتم که درخت را با مقدار آنتروپی تست کند اما آنرا کامنت کردم چون خروجی را زشت میکرد!

نحوه پیاده سازی:

بسیاری از توابع برای بخش رستوران را مشابه اسلاید ها و با توجه به سود و کد های اسلاید نوشتم، از جمله توابع آنتروپی و gain و remainder و تابع training برای همین توضیحی برای آنها اینجا نمی آورم.

تابع انتخاب attribute:

```
33 def select_an_attribute(data_file, attributes, res_name):
34     attributes_importance = {}
35     entropy = {}
36     remains = {}
37     for tmp_attributes in attributes:
38         attributes_importance[tmp_attributes], entropy[tmp_attributes], remains[tmp_attributes] = gain(data_file, tmp_attributes, res_name)
39     answer = max(attributes_importance.items(), key=operator.itemgetter(1))[0]
40     return answer, attributes_importance[answer], entropy[answer], remains[answer]
41
```

این تابع با بررسی attribute هایی که هنوز استفاده نشده اند، با توجه به میزان gain هر کدام، بهترین را انتخاب کرده و آنرا برمیگرداند.

```
16
17 def all_similar(data_file: pd.DataFrame):
18     a = data_file.to_numpy()
19     return (a[0] == a[1:]).all()
20
```

این تابع چک میکند آیا اعضای یک پنجره داده یا dataframe یکی هستند یا نه یعنی آیا تقسیم بندی ای داریم که نتایج آن همگی برابر باشند و یا نه.

```

117 def decision_tree_training(examples: pd.DataFrame, attributes: list, parent_examples, outcome_name, curr_depth
118                             ):
119     if examples.empty:
120         return value_node_of_plural_attributes(parent_examples, outcome_name)
121     if all_similar(examples[outcome_name]):
122         return decision_tree_leaf(examples[outcome_name].iloc[0], entropy=0)
123     if not attributes:
124         return value_node_of_plural_attributes(examples, outcome_name)
125     columns = list(examples.columns)
126     columns.remove(outcome_name)
127     attributes, gain_, entropy_, remains_ = select_an_attribute(examples, columns, outcome_name)
128     all_values = examples[attributes].unique()
129     tree = decision_tree_multi_way(attributes, value_node_of_plural_attributes(examples, outcome_name), gain_=gain_, entropy=entropy_,
130                                   remainder=remains_)
131     for vk in all_values:
132         new_columns = attributes
133         if (attributes in new_columns):
134             new_columns.remove(attributes)
135         subtree = decision_tree_training([examples[examples[attributes] == vk], new_columns, examples, outcome_name, curr_depth + 1
136                                         ])
137         tree.add(vk, subtree)
138     return tree
139

```

این تابع قلب اصلی کد است و عملیات یادگیری را انجام میدهد، همچنین آنرا مشابه اسلاید ها پیاده سازی کردم.

حال کلاس ها را توضیح میدهم:

```

97 class decision_tree_leaf:
98     def __init__(self, result, entropy=None):
99         self.result = result
100        self.entropy = entropy
101        def __call__(self, example):
102            return self.result
103        def __str__(self):
104            string = ' ' + str(self.result) + '\n'
105            return string
106        def display(self, scope=0):
107            print('RESULT =', self.result)
108

```

کلاس برگ درخت تصمیم که برای آن توابع روتین call و str را نوشته تا اگر خواستم آنرا روی خود نود صدا بزنم ممکن باشد، همچنین بتوان آنرا به string تبدیل کرد و چاپ کرد. هر نود برگ مقدار entropy و result خود را نگه

میدارد.

```
43 class decicion_tree_multi_way:
44     def __init__(self, attributes, default_child=None, branches=None, entropy=None, gain=None, remainder=None):
45         self.attributes = attributes
46         self.default_child = default_child
47         self.branches = branches or {}
48         self.entropy = entropy
49         self.gain_ = gain_
50         self.remainder_ = remainder_
51
52     def __call__(self, example):
53
54         attributes_val = example[self.attributes]
55         if attributes_val in self.branches:
56             return self.branches[attributes_val](example)
57         else:
58
59             return self.default_child(example)
60
61     def add(self, val, subtree):
62
63         self.branches[val] = subtree
64
65     def print_tree(self, scope=0):
66         name = self.attributes
67         print('Test ' + name + '?')
68         for (val, subtree) in self.branches.items():
69             if name == 'Patrons':
70                 if val == 0:
71                     val = 'None'
72                 elif val == 1:
73                     val = 'Some'
74                 else:
75                     val = 'Full'
76             elif name == 'Type':
77                 if val == 0:
78                     val = 'French'
79                 elif val == 1:
80                     val = 'Thai'
81                 elif val == 2:
82                     val = 'Burger'
83                 elif val == 3:
84                     val = 'Italian'
85             else:
86                 if val == 0:
87                     val = 'No'
88                 else:
89                     val = 'Yes'
90             print(' ' * 4 * scope, 'entropy = ', str(self.entropy) + ', ', name, '=', val, '---->', end=' ')
91             subtree.print_tree(scope + 1)
92
```

کلاس decision tree که علت multiway ایجاد قابلیت چند راهی در آن است، مجددا تابع call را برای آن نوشتم تا بتوان آنرا روی node صدا زد و بازگشتی آنرا پیمایش کرد. همچنین بجای str برای آن تابع print_tree را نوشتم تا بتوان درخت را چاپ کرد. در انتها نتیجه خروجی کد برای دیتای رستوران ها که فایل آنرا دستی ساختم اینگونه است:

The image shows a VS Code editor window with a Python script named `tree.py` and its output in the terminal. The script implements a decision tree algorithm for classifying restaurant patrons based on various features like `Alternate`, `Bar`, `Friday`, `Hungry`, `Patrons`, `Price`, `Rain`, `Reservation`, `Type`, `WaitEstimation`, and `Outcome`.

Script Content:

```
tree.py > decision_tree_multi_way > display
def calculate_entropy(q):
    if q == 0 or q == 1:
        return 0
    return calculate_entropy(q) + calculate_entropy(1 - q)

##### Driver Code #####
data_file = pd.read_csv('restaurants.csv')
columns = list(data_file.columns)
columns.remove('Outcome')

generated_tree = decision_tree_training(data_file, columns, None, 'Outcome', 0, 6)
generated_tree.display()
```

Terminal Output:

```
(base) InanAlipour@Inan-MacBook-Pro: IMP2_98102024 % /Users/InanAlipour/opt/anaconda3/bin/python /Users/InanAlipour/Documents/Programming/Python/AI/IMP2_98102024/tree.py
Test Patrons:
entropy = 1.0, Patrons = Some -> RESULT = 1
entropy = 1.0, Patrons = Full -> Test Hungry?
entropy = 0.9182958346544896, Hungry = Yes -> Test Type?
entropy = 1.0, Type = Thai -> Test Friday?
entropy = 1.0, Friday = No -> RESULT = 0
entropy = 1.0, Friday = Yes -> RESULT = 1
entropy = 1.0, Type = Italian -> RESULT = 0
entropy = 1.0, Type = Burger -> RESULT = 1
entropy = 0.9182958346544896, Hungry = No -> RESULT = 0
entropy = 1.0, Patrons = None -> RESULT = 0
(base) InanAlipour@Inan-MacBook-Pro: IMP2_98102024 %
```

restaurants.csv Data:

Alternate	Bar	Friday	Hungry	Patrons	Price	Rain	Reservation	Type	WaitEstimation	Outcome
1	0	0	1	1	2	0	1	0	0	1
3	1	0	0	1	2	0	0	0	1	2
4	0	1	0	0	1	0	0	2	0	1
5	1	0	1	2	0	1	0	1	1	1
6	1	0	1	0	2	0	1	0	3	0
7	0	1	0	1	1	1	1	3	0	1
8	0	1	0	0	0	1	0	2	0	0
9	0	0	0	1	1	1	1	1	0	1
10	0	1	0	2	0	1	0	2	3	0
11	1	1	1	2	2	0	1	3	1	0
12	0	0	0	0	0	0	0	1	0	0
13	1	1	1	2	0	0	0	2	2	1

درخت تولید شده با درخت اسلاید فرق میکند اما درست است، فرمت خروجی هم اینگونه است که هر لایه فرورفتگی نشانه پایین رفتن یک لایه در درخت است، اگر نیاز به تست جدید باشد ابتدا سوال مطرح میشود و زیر درخت چاپ میشود.

کد بخش دیابت:

حال با استفاده از کد قسمت قبل و ایجاد اندکی تغییر در آن، برای دیتاست دیابت کد را توضیح میدهم.

اولا کد نیاز به کمی تغییر داشت، همانطور که در چالش ها گفتیم، الگوریتم من با این دیتاست هی تلاش میکرد عمق را زیاد کند پس من با اضافه کردن یک `limit` برای عمق این چالش را حل کردم که نمیدانم ایده خوبی بود یا نه. بقیه توابع اضافه شده اینها هستند:

```
179 def accuracy_calculator(test_data, diabetes_decision_tree):
180     number_of_correct_guesses = 0
181     for i in range(len(test_data)):
182         test_res = diabetes_decision_tree(test_data.iloc[i])
183         if test_res == test_data.iloc[i].Outcome:
184             number_of_correct_guesses += 1
185
186     return (number_of_correct_guesses / len(test_data) * 100)
187
```

که دقت را برای داده های داده شده محاسبه میکند.

```
93 def discrete_find_clusters(data_file: pd.DataFrame, column_name, number_of_clusters):
94     max_i_value = data_file[column_name].max()
95     min_i_value = data_file[column_name].min()
96     difference = (max_i_value - min_i_value) / number_of_clusters
97     clusters = []
98     clusters.append(round(min_i_value - difference, 2))
99     for i in range(number_of_clusters):
100         clusters.append(round(min_i_value + difference * i, 2))
101     clusters.append(round(max_i_value, 2))
102     clusters.append(round(max_i_value + difference, 2))
103     return clusters
```

که داده هارا کلاستر بندی میکند و همانطور که در متن تمرین آمده آنها را گسسته میکند، درواقع مرز دسته هارا با `min` و `max` مشخص میکند و یک دسته برای مقادیر بیشتر و یا کمتر هم درنظر میگیرد(همانطور که در متن تمرین

آمده بود)

```
34 def discretify_data(training_data: pd.DataFrame, test_data: pd.DataFrame):
35     column_name = list(training_data.columns)
36
37     Pregnancies_clusters = discrete_find_clusters(training_data, 'Pregnancies', 5)
38     training_data = make_discrete(training_data, 'Pregnancies', Pregnancies_clusters)
39     test_data = make_discrete(test_data, 'Pregnancies', Pregnancies_clusters)
40
41     Glucose_clusters = discrete_find_clusters(training_data, 'Glucose', 5)
42     training_data = make_discrete(training_data, 'Glucose', Glucose_clusters)
43     test_data = make_discrete(test_data, 'Glucose', Glucose_clusters)
44
45     BloodPressure_clusters = discrete_find_clusters(training_data, 'BloodPressure', 5)
46     training_data = make_discrete(training_data, 'BloodPressure', BloodPressure_clusters)
47     test_data = make_discrete(test_data, 'BloodPressure', BloodPressure_clusters)
48
49     SkinThickness_clusters = discrete_find_clusters(training_data, 'SkinThickness', 5)
50     training_data = make_discrete(training_data, 'SkinThickness', SkinThickness_clusters)
51     test_data = make_discrete(test_data, 'SkinThickness', SkinThickness_clusters)
52
53     Insulin_clusters = discrete_find_clusters(training_data, 'Insulin', 5)
54     training_data = make_discrete(training_data, 'Insulin', Insulin_clusters)
55     test_data = make_discrete(test_data, 'Insulin', Insulin_clusters)
56
57     BMI_clusters = discrete_find_clusters(training_data, 'BMI', 5)
58     training_data = make_discrete(training_data, 'BMI', BMI_clusters)
59     test_data = make_discrete(test_data, 'BMI', BMI_clusters)
60
61     DiabetesPedigreeFunction_clusters = discrete_find_clusters(data_file, 'DiabetesPedigreeFunction', 5)
62     training_data = make_discrete(training_data, 'DiabetesPedigreeFunction', DiabetesPedigreeFunction_clusters)
63     test_data = make_discrete(test_data, 'DiabetesPedigreeFunction', DiabetesPedigreeFunction_clusters)
64
65     Age_clusters = discrete_find_clusters(data_file, 'Age', 5)
66     training_data = make_discrete(training_data, 'Age', Age_clusters)
67     test_data = make_discrete(test_data, 'Age', Age_clusters)
68
```

این تابع با خروجی تابع قبل ستون هارا گسسته سازی میکند، همچنین من مقادیر مختلفی را برای تعداد کلاستر ها امتحان کردم که ۵ و ۶ بهتر بودند و ۵ را نگه داشتم چون نسبت به ۶ عملکرد بهتری داشت.

```
198 data_file = pd.read_csv('diabetes.csv')
199 training_data, test_data = train_test_split(data_file, test_size=0.2)
200 training_data, test_data = discretify_data(training_data, test_data)
```

اینگونه داده هارا از فایل خواندم و با تابع آماده sklearn دیتاست را به یک بخش تست و یک بخش train تقسیم کردم که بخش تست ۰.۲ کل داده ها باشد، سپس با تابع discretify data داده هارا گسسته کردم، حال کافیت الگوریتم را اجرا کنیم

خروجی اینگونه است:

```
BloodPressure cluster = 5 ----> RESULT = 0
BloodPressure cluster = 0 ----> RESULT = 0
Age cluster = 4 ----> RESULT = 1
Age cluster = 0 ----> RESULT = 0
Age cluster = 3 ----> RESULT = 0
Pregnancies cluster = 2 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 1 ----> Test BMI cluster
BMI cluster = 3 ----> Test Age cluster
Age cluster = 0 ----> RESULT = 0
Age cluster = 1 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> RESULT = 0
SkinThickness cluster = 1 ----> RESULT = 0
SkinThickness cluster = 0 ----> RESULT = 0
Age cluster = 2 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> RESULT = 0
SkinThickness cluster = 1 ----> RESULT = 0
SkinThickness cluster = 0 ----> RESULT = 0
Age cluster = 3 ----> Test Insulin cluster
Insulin cluster = 0 ----> RESULT = 1
Insulin cluster = 1 ----> RESULT = 0
Age cluster = 4 ----> RESULT = 0
BMI cluster = 2 ----> RESULT = 0
BMI cluster = 4 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> RESULT = 0
SkinThickness cluster = 0 ----> Test BloodPressure cluster
BloodPressure cluster = 4 ----> RESULT = 0
BloodPressure cluster = 3 ----> RESULT = 1
SkinThickness cluster = 3 ----> RESULT = 0
BMI cluster = 5 ----> RESULT = 1
DiabetesPedigreeFunction cluster = 3 ----> RESULT = 1
DiabetesPedigreeFunction cluster = 2 ----> Test SkinThickness cluster
SkinThickness cluster = 1 ----> Test Pregnancies cluster
Pregnancies cluster = 2 ----> Test Pregnancies cluster
Pregnancies cluster = 2 ----> RESULT = 0
SkinThickness cluster = 0 ----> RESULT = 0
SkinThickness cluster = 3 ----> RESULT = 1
SkinThickness cluster = 2 ----> Test BloodPressure cluster
BloodPressure cluster = 2 ----> RESULT = 0
BloodPressure cluster = 4 ----> Test Age cluster
Age cluster = 2 ----> RESULT = 1
Age cluster = 1 ----> RESULT = 0
BloodPressure cluster = 3 ----> RESULT = 0
Pregnancies cluster = 3 ----> Test SkinThickness cluster
SkinThickness cluster = 0 ----> Test Age cluster
Age cluster = 2 ----> Test BloodPressure cluster
BloodPressure cluster = 3 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 1 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 2 ----> RESULT = 1
BloodPressure cluster = 0 ----> RESULT = 0
BloodPressure cluster = 4 ----> RESULT = 0
BloodPressure cluster = 5 ----> RESULT = 0
Age cluster = 1 ----> Test BMI cluster
BMI cluster = 3 ----> Test BloodPressure cluster
BloodPressure cluster = 0 ----> RESULT = 0
BloodPressure cluster = 3 ----> RESULT = 0
BloodPressure cluster = 4 ----> RESULT = 1
BMI cluster = 2 ----> RESULT = 0
BMI cluster = 0 ----> RESULT = 1
Age cluster = 4 ----> RESULT = 0
Age cluster = 3 ----> RESULT = 0
SkinThickness cluster = 3 ----> RESULT = 0
SkinThickness cluster = 2 ----> Test Age cluster
Age cluster = 2 ----> Test Insulin cluster
Insulin cluster = 2 ----> RESULT = 1
Insulin cluster = 0 ----> Test BloodPressure cluster
BloodPressure cluster = 4 ----> RESULT = 0
BloodPressure cluster = 3 ----> RESULT = 0
BloodPressure cluster = 5 ----> RESULT = 1
Insulin cluster = 1 ----> RESULT = 0
Age cluster = 3 ----> RESULT = 1
Age cluster = 1 ----> RESULT = 1
```



```
SkinThickness cluster = 1 ----> RESULT = 0
Pregnancies cluster = 4 ----> Test BloodPressure cluster
BloodPressure cluster = 3 ----> Test SkinThickness cluster
SkinThickness cluster = 3 ----> RESULT = 0
SkinThickness cluster = 0 ----> Test Age cluster
Age cluster = 2 ----> RESULT = 1
Age cluster = 3 ----> RESULT = 0
BloodPressure cluster = 4 ----> RESULT = 0
Pregnancies cluster = 5 ----> RESULT = 1
Glucose cluster = 5 ----> Test BMI cluster
BMI cluster = 3 ----> Test Insulin cluster
Insulin cluster = 0 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 1 ----> Test Pregnancies cluster
Pregnancies cluster = 2 ----> Test BloodPressure cluster
BloodPressure cluster = 3 ----> RESULT = 1
BloodPressure cluster = 4 ----> RESULT = 0
BloodPressure cluster = 0 ----> RESULT = 1
Pregnancies cluster = 3 ----> Test Age cluster
Age cluster = 3 ----> RESULT = 1
Age cluster = 4 ----> RESULT = 0
Age cluster = 2 ----> RESULT = 0
Pregnancies cluster = 0 ----> RESULT = 1
Pregnancies cluster = 5 ----> RESULT = 1
DiabetesPedigreeFunction cluster = 3 ----> RESULT = 1
DiabetesPedigreeFunction cluster = 2 ----> Test Pregnancies cluster
Pregnancies cluster = 3 ----> RESULT = 0
Pregnancies cluster = 1 ----> RESULT = 1
Pregnancies cluster = 2 ----> RESULT = 0
Insulin cluster = 1 ----> Test BloodPressure cluster
BloodPressure cluster = 3 ----> Test Pregnancies cluster
Pregnancies cluster = 2 ----> RESULT = 1
Pregnancies cluster = 1 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> RESULT = 0
SkinThickness cluster = 1 ----> RESULT = 1
Pregnancies cluster = 3 ----> RESULT = 1
BloodPressure cluster = 4 ----> RESULT = 0
Insulin cluster = 2 ----> RESULT = 1
Insulin cluster = 3 ----> Test Pregnancies cluster
Pregnancies cluster = 3 ----> RESULT = 1
Pregnancies cluster = 0 ----> RESULT = 1
Pregnancies cluster = 1 ----> RESULT = 0
Insulin cluster = 4 ----> RESULT = 1
Insulin cluster = 5 ----> RESULT = 1
BMI cluster = 4 ----> Test Age cluster
Age cluster = 1 ----> RESULT = 1
Age cluster = 2 ----> Test Pregnancies cluster
Pregnancies cluster = 3 ----> RESULT = 1
Pregnancies cluster = 1 ----> RESULT = 1
Pregnancies cluster = 0 ----> RESULT = 1
Pregnancies cluster = 2 ----> RESULT = 0
Age cluster = 3 ----> RESULT = 1
Age cluster = 4 ----> RESULT = 0
BMI cluster = 2 ----> RESULT = 1
BMI cluster = 5 ----> Test Insulin cluster
Insulin cluster = 1 ----> RESULT = 1
Insulin cluster = 2 ----> RESULT = 0
Insulin cluster = 5 ----> RESULT = 0
Glucose cluster = 4 ----> Test Age cluster
Age cluster = 4 ----> Test BMI cluster
BMI cluster = 3 ----> Test Pregnancies cluster
Pregnancies cluster = 2 ----> Test Insulin cluster
Insulin cluster = 0 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> RESULT = 0
SkinThickness cluster = 0 ----> RESULT = 0
Insulin cluster = 2 ----> RESULT = 1
Pregnancies cluster = 3 ----> RESULT = 0
Pregnancies cluster = 4 ----> RESULT = 0
Pregnancies cluster = 1 ----> Test SkinThickness cluster
SkinThickness cluster = 1 ----> RESULT = 0
SkinThickness cluster = 0 ----> RESULT = 1
Pregnancies cluster = 0 ----> RESULT = 0
```

```

BMI cluster = 2 ----> RESULT = 0
BMI cluster = 4 ----> RESULT = 1
BMI cluster = 0 ----> RESULT = 0
Age cluster = 1 ----> Test BloodPressure cluster
BloodPressure cluster = 4 ----> Test BMI cluster
BMI cluster = 4 ----> Test Pregnancies cluster
Pregnancies cluster = 1 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 2 ----> RESULT = 0.0
DiabetesPedigreeFunction cluster = 1 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 3 ----> RESULT = 0
Pregnancies cluster = 2 ----> Test Insulin cluster
Insulin cluster = 0 ----> RESULT = 1
Insulin cluster = 1 ----> RESULT = 0
Insulin cluster = 2 ----> RESULT = 1
Pregnancies cluster = 0 ----> RESULT = 0
BMI cluster = 3 ----> Test Pregnancies cluster
Pregnancies cluster = 2 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 1 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 2 ----> RESULT = 1
Pregnancies cluster = 1 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 2 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 1 ----> RESULT = 0
Pregnancies cluster = 0 ----> Test SkinThickness cluster
SkinThickness cluster = 1 ----> RESULT = 0
SkinThickness cluster = 0 ----> RESULT = 0
SkinThickness cluster = 2 ----> RESULT = 0
Pregnancies cluster = 3 ----> RESULT = 1
BMI cluster = 5 ----> RESULT = 1
BloodPressure cluster = 3 ----> Test Pregnancies cluster
Pregnancies cluster = 1 ----> Test Insulin cluster
Insulin cluster = 0 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> RESULT = 0
SkinThickness cluster = 0 ----> RESULT = 0
SkinThickness cluster = 3 ----> RESULT = 0
Insulin cluster = 4 ----> RESULT = 1
Insulin cluster = 1 ----> Test BMI cluster
BMI cluster = 3 ----> RESULT = 0
BMI cluster = 4 ----> RESULT = 0
Insulin cluster = 2 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 1 ----> RESULT = 1
DiabetesPedigreeFunction cluster = 4 ----> RESULT = 0
Insulin cluster = 3 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> RESULT = 1
SkinThickness cluster = 1 ----> RESULT = 0
Pregnancies cluster = 0 ----> RESULT = 1
Pregnancies cluster = 2 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 1 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> RESULT = 0
SkinThickness cluster = 0 ----> RESULT = 0
SkinThickness cluster = 1 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 3 ----> RESULT = 1
DiabetesPedigreeFunction cluster = 2 ----> RESULT = 0
Pregnancies cluster = 3 ----> RESULT = 1
BloodPressure cluster = 0 ----> RESULT = 1
BloodPressure cluster = 5 ----> Test SkinThickness cluster
SkinThickness cluster = 3 ----> RESULT = 1
SkinThickness cluster = 2 ----> RESULT = 0
BloodPressure cluster = 2 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> RESULT = 0
SkinThickness cluster = 3 ----> RESULT = 1
SkinThickness cluster = 1 ----> RESULT = 0
Age cluster = 3 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> Test BMI cluster
BMI cluster = 3 ----> RESULT = 1
BMI cluster = 4 ----> Test Pregnancies cluster
Pregnancies cluster = 4 ----> RESULT = 1
Pregnancies cluster = 2 ----> RESULT = 1
Pregnancies cluster = 3 ----> RESULT = 0
BMI cluster = 2 ----> RESULT = 0
SkinThickness cluster = 0 ----> Test BMI cluster
BMI cluster = 3 ----> Test Pregnancies cluster

```

```
Pregnancies cluster = 3 --> Test DiabetesPedigreeFunction cluster
  DiabetesPedigreeFunction cluster = 1 --> RESULT = 0
  DiabetesPedigreeFunction cluster = 3 --> RESULT = 0
Pregnancies cluster = 1 --> RESULT = 1
Pregnancies cluster = 2 --> Test BloodPressure cluster
  BloodPressure cluster = 4 --> RESULT = 1
  BloodPressure cluster = 3 --> RESULT = 0
BMI cluster = 0 --> RESULT = 1
BMI cluster = 4 --> Test Pregnancies cluster
  Pregnancies cluster = 4 --> RESULT = 0
  Pregnancies cluster = 3 --> RESULT = 0
  Pregnancies cluster = 2 --> RESULT = 0
  Pregnancies cluster = 0 --> RESULT = 1
BMI cluster = 2 --> RESULT = 0
SkinThickness cluster = 3 --> RESULT = 1
SkinThickness cluster = 1 --> RESULT = 0
Age cluster = 2 --> Test Pregnancies cluster
  Pregnancies cluster = 3 --> Test SkinThickness cluster
    SkinThickness cluster = 0 --> Test BloodPressure cluster
      BloodPressure cluster = 3 --> Test DiabetesPedigreeFunction cluster
        DiabetesPedigreeFunction cluster = 1 --> RESULT = 0
        DiabetesPedigreeFunction cluster = 2 --> RESULT = 1
      BloodPressure cluster = 4 --> Test BMI cluster
        BMI cluster = 3 --> RESULT = 1
        BMI cluster = 4 --> RESULT = 0
      BloodPressure cluster = 0 --> RESULT = 1
    SkinThickness cluster = 2 --> Test Insulin cluster
      Insulin cluster = 3 --> Test BloodPressure cluster
        BloodPressure cluster = 3 --> RESULT = 0
        BloodPressure cluster = 4 --> RESULT = 1
      Insulin cluster = 0 --> RESULT = 1
      Insulin cluster = 1 --> Test BloodPressure cluster
        BloodPressure cluster = 3 --> RESULT = 1
        BloodPressure cluster = 4 --> RESULT = 0
      Insulin cluster = 2 --> Test BloodPressure cluster
        BloodPressure cluster = 3 --> RESULT = 0
        BloodPressure cluster = 4 --> RESULT = 1
    SkinThickness cluster = 1 --> RESULT = 0
    SkinThickness cluster = 3 --> Test Insulin cluster
      Insulin cluster = 0 --> RESULT = 0
      Insulin cluster = 1 --> RESULT = 1
      Insulin cluster = 3 --> RESULT = 0
  Pregnancies cluster = 4 --> Test BloodPressure cluster
    BloodPressure cluster = 0 --> RESULT = 1
    BloodPressure cluster = 3 --> RESULT = 1
    BloodPressure cluster = 4 --> Test Insulin cluster
      Insulin cluster = 1 --> RESULT = 1
      Insulin cluster = 0 --> Test DiabetesPedigreeFunction cluster
        DiabetesPedigreeFunction cluster = 1 --> RESULT = 0
        DiabetesPedigreeFunction cluster = 2 --> RESULT = 1
    BloodPressure cluster = 5 --> RESULT = 1
  Pregnancies cluster = 0 --> Test BloodPressure cluster
    BloodPressure cluster = 4 --> RESULT = 0
    BloodPressure cluster = 3 --> RESULT = 1
  Pregnancies cluster = 5 --> RESULT = 1
  Pregnancies cluster = 2 --> Test Insulin cluster
    Insulin cluster = 2 --> RESULT = 0
    Insulin cluster = 1 --> Test SkinThickness cluster
      SkinThickness cluster = 1 --> RESULT = 1
      SkinThickness cluster = 2 --> RESULT = 0
      SkinThickness cluster = 3 --> RESULT = 1
    Insulin cluster = 0 --> Test BMI cluster
      BMI cluster = 3 --> Test BloodPressure cluster
        BloodPressure cluster = 4 --> RESULT = 0
        BloodPressure cluster = 3 --> RESULT = 0
      BMI cluster = 4 --> RESULT = 1
    Insulin cluster = 4 --> RESULT = 0
  Pregnancies cluster = 1 --> Test SkinThickness cluster
    SkinThickness cluster = 2 --> RESULT = 1
    SkinThickness cluster = 3 --> RESULT = 0
    SkinThickness cluster = 0 --> RESULT = 0
```

```
[ 'Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome' ]
Test Glucose cluster
Glucose cluster = 3 ----> Test Pregnancies cluster
Pregnancies cluster = 1 ----> Test BloodPressure cluster
BloodPressure cluster = 3 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 1 ----> Test BMI cluster
BMI cluster = 4 ----> RESULT = 0
BMI cluster = 3 ----> RESULT = 0
BMI cluster = 2 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 2 ----> Test BMI cluster
BMI cluster = 3 ----> RESULT = 0
BMI cluster = 4 ----> RESULT = 1
DiabetesPedigreeFunction cluster = 3 ----> RESULT = 0
SkinThickness cluster = 0 ----> RESULT = 0
SkinThickness cluster = 1 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 1 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 2 ----> Test BMI cluster
BMI cluster = 2 ----> RESULT = 0
BMI cluster = 3 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 3 ----> RESULT = 0
SkinThickness cluster = 3 ----> Test Insulin cluster
Insulin cluster = 1 ----> RESULT = 0
Insulin cluster = 0 ----> RESULT = 1
BloodPressure cluster = 4 ----> Test SkinThickness cluster
SkinThickness cluster = 2 ----> Test Insulin cluster
Insulin cluster = 1 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 3 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 1 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 2 ----> RESULT = 0
Insulin cluster = 2 ----> RESULT = 1
Insulin cluster = 0 ----> RESULT = 0
SkinThickness cluster = 1 ----> RESULT = 0
SkinThickness cluster = 0 ----> Test Age cluster
Age cluster = 3 ----> Test BMI cluster
BMI cluster = 2 ----> RESULT = 0
BMI cluster = 4 ----> RESULT = 0
BMI cluster = 3 ----> RESULT = 1
Age cluster = 1 ----> RESULT = 0
Age cluster = 0 ----> Test Pregnancies cluster
Pregnancies cluster = 1 ----> RESULT = 0
Age cluster = 2 ----> Test BMI cluster
BMI cluster = 3 ----> RESULT = 0
BMI cluster = 4 ----> RESULT = 1
SkinThickness cluster = 3 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 1 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 2 ----> RESULT = 1
BloodPressure cluster = 5 ----> RESULT = 0
BloodPressure cluster = 0 ----> RESULT = 0
BloodPressure cluster = 2 ----> RESULT = 0
Pregnancies cluster = 0 ----> Test Age cluster
Age cluster = 1 ----> Test SkinThickness cluster
SkinThickness cluster = 1 ----> RESULT = 0
SkinThickness cluster = 0 ----> Test BMI cluster
BMI cluster = 3 ----> Test BloodPressure cluster
BloodPressure cluster = 4 ----> RESULT = 0
BloodPressure cluster = 0 ----> RESULT = 0
BloodPressure cluster = 3 ----> RESULT = 0
BMI cluster = 2 ----> RESULT = 0
BMI cluster = 4 ----> RESULT = 0
BMI cluster = 0 ----> RESULT = 0
SkinThickness cluster = 2 ----> Test Insulin cluster
Insulin cluster = 1 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 1 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 2 ----> RESULT = 0
Insulin cluster = 0 ----> RESULT = 0
Insulin cluster = 2 ----> RESULT = 0
SkinThickness cluster = 3 ----> RESULT = 0
SkinThickness cluster = 4 ----> RESULT = 0
Age cluster = 2 ----> Test BloodPressure cluster
BloodPressure cluster = 4 ----> RESULT = 1
SkinThickness cluster = 0 ----> RESULT = 0
Age cluster = 0 ----> RESULT = 0
Age cluster = 5 ----> Test Pregnancies cluster
Pregnancies cluster = 3 ----> RESULT = 0
Pregnancies cluster = 2 ----> RESULT = 1
Glucose cluster = 2 ----> Test Pregnancies cluster
Pregnancies cluster = 0 ----> RESULT = 0
Pregnancies cluster = 4 ----> RESULT = 0
Pregnancies cluster = 3 ----> RESULT = 0
Pregnancies cluster = 1 ----> Test Insulin cluster
Insulin cluster = 1 ----> Test BloodPressure cluster
BloodPressure cluster = 4 ----> RESULT = 0
BloodPressure cluster = 3 ----> Test DiabetesPedigreeFunction cluster
DiabetesPedigreeFunction cluster = 1 ----> Test Pregnancies cluster
Pregnancies cluster = 1 ----> RESULT = 0
DiabetesPedigreeFunction cluster = 3 ----> RESULT = 0
BloodPressure cluster = 2 ----> RESULT = 0
Insulin cluster = 0 ----> RESULT = 0
Pregnancies cluster = 2 ----> RESULT = 0
Glucose cluster = 0 ----> Test Pregnancies cluster
Pregnancies cluster = 1 ----> RESULT = 0
Pregnancies cluster = 2 ----> RESULT = 1
training data accuracy = 92.18241042345277
(base) ImanAlipour@Imans-MacBook-Pro IMP2_98102024 %
```

همانطور که دیده میشود دقت روی داده های train ۹۲ درصد است و روی داده های تست 65 درصد، با کاهش داده های test و افزایش داده های تست به ۰.۵ و ۰.۵ دقت زیاد شد که خلاف انتظار من بود، من انتظار داشتم با کاهش داده train عملکرد در داده های test پایین بیاید که این اتفاق رخ نداد و دقت برای train به ۹۵ رسید و برای تست به 70! که نشان از overfitting است. همچنین با افزایش محدودیت عمق overfitting کمتر رخ میداد اما زمان اجرای کد طولانی تر میشد.

آموخته های من از این تمرین:

اولا با پیاده سازی شبه کد های آنتروپی و محاسبه B و remainder کمی نسبت به آنها احساس ترس کمتری پیدا کردم، متأسفانه خیلی باگ های زیادی خوردم و بیشتر آن هم به این دلیل بود که سعی داشتم کدم را کوتاه بنویسم پ با خیلی از روش هایی که استفاده کردم دوستی زیادی نداشتم و این منجر به طولانی شدن فرایند کد زدن شد.

در کل نسبت به سختی این مسئله و دلیل اینکه بجای این روش از روش های دیگر استفاده میشود و عدم کارایی زیاد آن آگاهی پیدا کردم.

برای افزایش دقت حس میکنم افزایش یا برداشتن محدودیت ارتفاع درخت ایده خوبی باشد، همچنین روش های هوشمندانه تر برای گسسته سازی داده قطعا راهگشا خواهند بود. در این [لینک](#) اینطور ایده ها آمده است که من آنها البته پیاده سازی نکردم اما برای بسیار جذاب بودند.

حس میکنم با کمی fine tuning میشد دقت را بیشتر کرد اما خب این به تست های زیاد نیاز دارد. در کل با انجام تمرین حس میکنم فهم بهتری نسبت به گسسته سازی داده ها و محاسبه آنتروپی و ... پیدا کردم. در انتها باید بگویم که بسیاری از ایده های پیاده سازی را از این [لینک](#) گرفتم، در واقع در ابتدا خودم بخشی از کد را زدم اما چون نمیدانستم مسیر را چگونه باید ادامه دهم کمی گشتم و یک مثال دیدم و سپس آنرا سعی کردم خودم پیاده سازی کنم که به باگ های بسیار بسیار زیادی برخورددم، خواهشا این را در نظر بگیرید کدی کپی نکردم و فقط ایده های آنها را دیدم تا بفهمم باید چگونه کلاستر بندی و ... را انجام دهم.

همچنین از sklearn فقط برای تقسیم کردن دیتاست استفاده کردم و از آن استفاده دیگری نکردم!