

PyCharm

File Edit View Navigate Code Refactor Run Tools VCS Window Help

MLA2 - ID3.ipynb

100% Mon Oct 19 5:12 PM

MLA2 ID3.ipynb

Managed: http://localhost:8786 Python 3

Trusted

73

# builds and returns a dictionary decision tree

74

def build\_tree(examples, attributes):

75

pure = examples[target\_feature].value\_counts()

76

majority\_target\_feature\_value = pure.idxmax()

77

# handles all 3 base cases: if all target values are

78

if pure.size == 1 or len(attributes) == 0 or len(examples) == 0:

79

return majority\_target\_feature\_value

80

else:

81

best\_feature = d\_best(examples, attributes)

82

root = {best\_feature : {}} # dictionary of dict

83

# remove best\_feature from attributes to not test it again

84

attributes = [i for i in attributes if i != best\_feature]

85

unique\_values = unique\_attribute\_values[best\_feature]

86

for value in unique\_values:

87

# this selects all same value rows of an attribute

88

partition\_examples = examples.loc[examples[best\_feature] == value]

89

if partition\_examples.empty:

90

root[best\_feature][value] = majority\_target\_feature\_value

91

else: # make a subtree by passing same value

92

root[best\_feature][value] = build\_tree(partition\_examples, attributes)

93

return root

94

95

96

decision\_tree = build\_tree(example\_df, features)

97

decision\_tree

98

99

100

serialize()

162

{'Outlook': {'Sunny': {'Humidity': {'High': 'No', 'Normal': 'Yes'}}, 'Rain': {'Wind': {'Weak': 'Yes', 'Strong': 'No'}}, 'Overcast': 'Yes'}}

166

```
graph TD; Outlook([Outlook]) --> Sunny([Sunny]); Outlook --> Rain([Rain]); Outlook --> Overcast([Overcast]); Sunny --> Humidity([Humidity]); Humidity --> High([High]); High --> No([No]); Humidity --> Normal([Normal]); Normal --> Yes([Yes]); Rain --> Wind([Wind]); Wind --> Strong([Strong]); Strong --> No; Wind --> Weak([Weak]); Weak --> Yes; Overcast --> Yes;
```

6: Problems

Terminal

Git

TODO

Python Console

Event Log

PyCharm

File Edit View Navigate Code Refactor Run Tools VCS Window Help

MLA2 - ID3.ipynb

100% 5:13 PM

MLA2 ID3.ipynb

Managed: http://localhost:8786 Python 3

Trust

8

play\_tennis = pd.read\_csv("assets/playtennis.csv")

9

emails = pd.read\_csv("assets/emails.csv")

10

vegetation = pd.read\_csv("assets/vegetation.csv")

11

census = pd.read\_csv("assets/census\_training.csv")

12

13

# TYPE NAME OF CSV FILE TO TRAIN (census default)

14

example\_df = emails

15

pruning = 0

16

17

unique\_attribute\_values = {}

18

19

###

20

# removes columns with all distinct values to avoid over

21

def remove\_distinct\_attributes(data\_frame):

22

for column in data\_frame.columns:

23

if len(data\_frame[column].unique()) == len(data

24

del data\_frame[column]

25

return data\_frame

26

27

28

example\_df = remove\_distinct\_attributes(example\_df)

29

30

features = example\_df.columns.values.tolist()

31

target\_feature = features[-1]

32

features.pop() # remove target feature

33

# add feature to dictionary as a key and unique values

34

for feature in features:

35

unique\_attribute\_values[feature] = pd.value\_counts(e

172

{'SUSPICIOUS WORDS': {True: 'spam', False: 'ham'}}

176

```
graph TD; A([SUSPICIOUS WORDS]) --> B([True]); A --> C([False]); B --> D[spam]; C --> E[ham]
```

6 Problems

Terminal

Git

TODO

Python Console

Event Log

MLA2

ID3.ipynb

Python 3

Managed: http://localhost:8786

Python 3

Trust

Database

SciView

1: Project

2: Favorites

3: Problems

4: Terminal

5: Git

6: TODO

7: Python Console

8: Event Log

```
38 # returns uncertainty of an
39 def entropy(data_frame, attribute):
40     attribute_value_counts = data_fr
41     set_count = len(data_frame)
42     entropy = 0
43     for value_count in attribute_val
44         probability = value_count/se
45         entropy += -probability * np
46     return entropy
47
48 """
49 # returns information of gain of an
50 def information_gain(data_frame, att
51     value_counts = data_frame[attrib
52     target_feature_entropy = entropy
53     attribute_entropy = 0
54     for value in value_counts.index:
55         probability = value_counts[v
56         attribute_entropy += probabi
57     return target_feature_entropy -
58
59 """
60 # calculates best classifier to all a
61 def d_best(data_frame, attributes):
62     entropy = -1
63     best_gain = None
64     for attribute in attributes:
65         current_entropy = informatio
66         if current_entropy > entropy
```

```
184 { ELEVATION : { high : { SLOPE : { steep : chaparral ,
    'moderate': 'chaparral',
    'flat': 'conifer'}},
    'medium': { 'STREAM': {True: 'riparian', False: 'chaparral'}},
    'highest': 'conifer',
    'low': 'riparian'}}
```

```
186
```



