

حسین غلامی تمرین سری دو یادگیری ماشین - درخت تصمیم

در انجام این تمرین از آن جایی که فرمت اطلاعات به صورت SSV بودند به جهت ساده سازی و استفاده از متدهای pandas

نحوہ تبدیل:

نموده

جایگزین نموده و فایل را با فرمت CSV

ذخیره میکنیم

خروجی به شکل زیر حاصل شد از آن جایی که فایل CSV را در اکسل نیز میتوان مشاهده کرد، تصویر آن نیز قرار گرفت

[illegible]

پیاده سازی الگوریتم :

در سه مرحله انجام شد:

- 1- پیاده سازی اولیه
- 2- بررسی کتابخانه ها و عملکردشان و نوع کارکردن با آن ها
- 3- پیاده سازی نهایی

پیاده سازی اولیه

به این صورت که در ابتدا ریشه درخت مشخص میشود و در ادامه شاخه ها و نود ها پس از هم محاسبه میشوند. ولی این که پس از تمام شدن درخت چگونه از روی شاخه ها بازگشت ، قابل پیاده سازی نبود ، چرا که نام ویژگی های یکسانی وجود داشت و پیاده سازی امکان پذیر نبود

```
1 import pandas as pd
2 from math import log
3
4 def entr(x,y):# x:pos y:neg
5     return -((x/(x+y))*log((x/(x+y))))-((y/(x+y))*log((y/(x+y))))
6
7
8 path='F:/Data/noisy_train.csv'
9 df=pd.read_csv(path)
10
11 entropy=list()
12 gain=list()
13 tree=list()
14
15 index = df.columns.tolist()
16 index.remove('poisonous')
17 for feather in index:
18     for types ,z in df.groupby([feather]):
19         if(((0,types) in df.groupby(['poisonous',feather]).groups.keys())
20             &((1,types) in df.groupby(['poisonous',feather]).groups.keys())):
21             entropy.append(
22                 entr(df.groupby(['poisonous',feather]).size()[0][types]
23                     ,df.groupby(['poisonous',feather]).size()[1][types])
24                 *
25                 (df.groupby([feather]).size()[types]
26                 /df.groupby([feather]).size().sum() )
27             )
28         else:
29             entropy.append(0)
30     gain.append([sum(entropy),feather])
31     # print(entropy)
32     for types in range(0,len(entropy)):
33         entropy.pop()
34 gain.sort()
35 print(gain)
36 root=gain[0][1]
37 for types in range(0,len(gain)):
38     gain.pop()
```

```

39 #root has found
40
41 zanjir=list()
42 print('_____ \n')
43 stop=0
44 count =0
45 t=1
46 index.remove(root)
47
48 tree_lvl=-1
49
50 internal_counter=-1
51
52 node=root
53
54 while(stop==0):
55     #mohasebe zanjir ha
56     if( tree_lvl>0 ):
57         stop=1
58
59     if(count > 0):
60
61         if(internal_counter<0):
62             internal_counter=len(df.groupby([node]).size())
63             tree_lvl=tree_lvl+1
64
65
66
67         internal_counter=internal_counter-1
68         if(node in index ):
69             index.remove(node)
70
71
72         if(len(zanjir[tree_lvl])==6):
73             node=zanjir[internal_counter][3]
74             index.remove(node)
75         else:
76             node=zanjir[count][3]

```

```

81 count =count+1
82     ##calculate zanjir
83     for branch ,z1 in df.groupby([node]):
84         for feather in index:
85             for types ,z2 in df.groupby([feather]):
86                 if((1,types,branch) in df.groupby(['poisonous',feather,node]).groups.keys())
87                     &((0,types,branch) in df.groupby(['poisonous',feather,node]).groups.keys())):
88                     entropy.append(
89                         entr(df.groupby(['poisonous',feather,node]).size()[0][types][branch]
90                             ,df.groupby(['poisonous',feather,node]).size()[1][types][branch])
91                             *
92                             (df.groupby([feather]).size()[types]
93                             /df.groupby([feather]).size().sum() )
94                     )
95                 else:
96                     entropy.append(0)
97             GAIN=sum(entropy)
98             for types in range(0,len(entropy)):
99                 entropy.pop()
100             #label zadan
101             if(GAIN==0):
102                 t=1
103             else:
104                 gain.append([GAIN,feather])
105
106 gain.sort()
107 #print(branch,gain)
108 try:
109     if(t==0):
110         zanjir.append((tree_lvl,node,branch,gain[0][1]))
111     if(t==1):
112         zanjir.append((tree_lvl,node,branch,gain[0][1],'end'))
113     t=0
114
115     print(count,node,branch,gain[0][1])
116     for types in range(0,len(gain)):
117         gain.pop()

```

که به خروجی زیر رسیدم که شاخه های درخت را تشکیل میداد ولی برای محاسبه درخت مشکل داشت:

```

In [7]: zanjir
Out[7]:
[(-1, 'odor', 'a', 'stalkcolorbelowring', 'end'), (0, 'sporeprintcolor', 'h', 'capcolor'),
(-1, 'odor', 'c', 'gillsize'), (0, 'sporeprintcolor', 'k', 'stalksurfaceabovering'),
(-1, 'odor', 'f', 'sporeprintcolor'), (0, 'sporeprintcolor', 'n', 'stalksurfaceabovering'),
(-1, 'odor', 'l', 'stalkcolorbelowring'), (0, 'sporeprintcolor', 'n', 'ringnumber'),
(-1, 'odor', 'm', 'gillattachment'), (0, 'sporeprintcolor', 'u', 'gillspacing'),
(-1, 'odor', 'n', 'sporeprintcolor'), (0, 'sporeprintcolor', 'w', 'gillcolor'),
(-1, 'odor', 'p', 'gillsize'), (0, 'gillsize', 'b', 'stalksurfaceabovering'),
(0, 'gillsize', 'b', 'sporeprintcolor'), (0, 'gillsize', 'n', 'stalkcolorabovering'),
(0, 'gillsize', 'n', 'gillsize'), (0, 'sporeprintcolor', 'h', 'capcolor'),
(0, 'sporeprintcolor', 'h', 'sporeprintcolor'), (0, 'sporeprintcolor', 'k', 'stalksurfaceabovering'),
(0, 'sporeprintcolor', 'k', 'sporeprintcolor'), (0, 'sporeprintcolor', 'n', 'stalksurfaceabovering'),
(0, 'sporeprintcolor', 'n', 'sporeprintcolor'), (0, 'sporeprintcolor', 'n', 'ringnumber'),
(0, 'sporeprintcolor', 'r', 'sporeprintcolor'), (0, 'sporeprintcolor', 'u', 'gillspacing'),
(0, 'sporeprintcolor', 'r', 'sporeprintcolor'), (0, 'sporeprintcolor', 'w', 'gillcolor'),
(0, 'sporeprintcolor', 'u', 'sporeprintcolor'), (0, 'gillsize', 'b', 'stalksurfaceabovering'),
(0, 'sporeprintcolor', 'w', 'sporeprintcolor'), (0, 'gillsize', 'n', 'stalkcolorabovering'),
(0, 'stalkcolorbelowring', 'b', 'stalkcolorbelowring'), (1, 'sporeprintcolor', 'h', 'capcolor'),
(0, 'stalkcolorbelowring', 'c', 'gillattachment'), (1, 'sporeprintcolor', 'k', 'stalksurfaceabovering'),
(0, 'stalkcolorbelowring', 'g', 'stalkcolorbelowring'), (1, 'sporeprintcolor', 'n', 'stalksurfaceabovering'),
(0, 'stalkcolorbelowring', 'n', 'stalkcolorbelowring'), (1, 'sporeprintcolor', 'r', 'ringnumber'),
(0, 'stalkcolorbelowring', 'p', 'stalkcolorbelowring'), (1, 'sporeprintcolor', 'u', 'gillspacing'),
(0, 'stalkcolorbelowring', 'w', 'stalkcolorbelowring'), (1, 'sporeprintcolor', 'w', 'gillcolor'),
(0, 'gillattachment', 'a', 'gillattachment', 'end'), (1, 'sporeprintcolor', 'h', 'capcolor'),
(0, 'gillattachment', 'f', 'ringtype'), (1, 'sporeprintcolor', 'k', 'stalksurfaceabovering'),
(1, 'sporeprintcolor', 'n', 'stalksurfaceabovering'),
(1, 'sporeprintcolor', 'n', 'ringnumber')
.....

```

2-1) در مرحله بعد برای بررسی دقیق سایرین و فرایندی که در آینده استفاده خواهد شد ، به بررسی کتابخانه های موجود پرداختم که با جست و جویی ساده به دو کتابخانه Sklearn.tree و Id3 رسیدم ، که هر کدام را توضیح و مزایا و معایب آن ها را از دیدگاه خودم بررسی میکنم

Sklearn.tree کتابخانه ای کامل که علاوه بر روش entropy میتوان از روش gini نیز استفاده نمود. توانایی pruning

وجود ندارد. همچنین یکی از مشکلات اصلی بحث Passing categorical data to Sklearn Decision Tree است که

نمیتوان به این تابع مقدار کمی مانند 'a' نسبت داد و باید همه اطلاعات عدد باشند. راه حل آن این است که از

```
from sklearn import preprocessing
le = preprocessing.LabelEncoder()
le.fit(["paris", "paris", "tokyo", "amsterdam"])
le.transform(["tokyo", "tokyo", "paris"])
```

استفاده کرد و اطلاعات را به عدد تبدیل نمود و سپس از طریق

```
list(le.inverse_transform([2, 2, 1]))
```

به حالت اولیه برگرداند.

یک نمونه در ضمیمه قرار گرفت استفاده در ضمیمه قرار گرفت.

2-2) id3 کتابخانه و کلاس بعدی است که میتوان با نام decision-tree-id3 آن را نصب نمود. که بعد از اجرا خروجی dot.

میدهد ، همچنین میتوان خروجی dot. را با graphviz به فرمت png برده و نمایش داد.

از لینک زیر برای نصب graphviz می توان استفاده کرد :

<https://bobswift.atlassian.net/wiki/spaces/GVIZ/pages/20971549/How+to+install+Graphviz+software>

ware

باتوجه به نوع اطاللات موجود استفاده از id3 مناسب تر می باشد چرا که علاوه پذیرفتن و قبول کردن اطلاعات کتگورکال میتواند pruning را انجام دهد.

حال به بررسی اطلاعات و یافتن بهترین سطح هرس میپردازیم

```
10
17 feature_names = tr.columns.tolist()
18 feature_names.remove('poisonous')
19 #
20
21 Xtr=tr.values.tolist()
22 for i in range(0,len(tr.values)):
23     del Xtr[i][0]
24 ytr=list(tr['poisonous'])
25 #
26 Xva=va.values.tolist()
27 for i in range(0,len(va.values)):
28     del Xva[i][0]
29 yva=list(va['poisonous'])
30 #
31
32 Xtr_train, Xtr_test, ytr_train, ytr_test = train_test_split(Xtr, ytr
33                                                             ,test_size=0.0,shuffle=False)
34
35 Xva_train, Xva_test, yva_train, yva_test = train_test_split(Xva, yva
36                                                             ,test_size=0.0,shuffle=False)
37 # << start Learning >>
38
39 estimator = Id3Estimator(max_depth=None, min_samples_split=1, prune=False,
40                           gain_ratio=False, min_entropy_decrease=0.0, is_repeating=False)
41
42 estimator.fit(Xtr_train, ytr_train, check_input=False)
43 export_graphviz(estimator.tree_, 'tree.dot', feature_names,extensive=True)
44
45 y_predtr=estimator.predict(Xtr_train)
46
47
48 y_predva=estimator.predict(Xva_train)
49
50
51 print('Accuracy Score on train data: ',
52       accuracy_score(y_true = ytr_train , y_pred=y_predtr )
53       )
54 print('Accuracy Score on test data: ',
55       accuracy_score(y_true=yva_train,y_pred=y_predva)
56       )
57
```

بدون prune نتیجه به روی داده validation به شرح زیر است ،

```
Accuracy Score on train data:  1.0
Accuracy Score on test data:  0.8038984051978736
```

حال prune را فعال نموده ونتیجه :

```
Accuracy Score on train data:  0.9384138236597253
Accuracy Score on test data:  0.8411104548139398
```

حال `prune` را مجدد غیرفعال نموده و محدودیت را روی تعداد عمق درخت میبریم ، عمق درخت ما 11 است. پس به کم کردن آن میپردازیم تا جایی که میزان دقت اطلاعات `validation` حداکثر شود. اگر عمق درخت را 2 قرار دهیم به

```
Accuracy Score on train data: 0.898980948161276
Accuracy Score on test data: 0.8907265209686946
```

میرسیم که بیشترین میزان دقت در داده تست است.

حالت های مختلفی بررسی شد ، که بهترین نتیجه را بر روی داده `validation` از طریق کم کردن شاخه ها دریافت شد.

حال به رسم نتایج میپردازیم

فایل `.dot` را با کمک `graphviz` که پیشتر نصبش توضیح داده شد ، توسط `command`

`dot file.dot -Tpng -o image.png` در `cmd` به فایل `png` تبدیل میکنیم

درخت شماره 1



درخت شماره 2



درخت شماره 3



**** تصاویر در ضمیمه قرار گرفت ، اینجا نمیشد بزرگ کرد.**

حال اطلاعات `train` شده را برای بهترین حالت یادگیری (تغییر عمق شاخه بررسی می کنیم) که نتیجه زیر حاصل شد:

```
Accuracy Score on train data: 0.898980948161276
Accuracy Score on test data: 0.9037780401416765
```

روش سوم

با توجه به مرحله قبل دریافتیم که بهترین پاسخ را میتوان از تعیین عمق شاخه بدست آورد.

پس به پیاده سازی کامل بدون استفاده از کتابخانه id3 پرداخته شد. به صورتی که تابعی با ورودی های 1- دیتا فریم 2-لیست ویژگی ها 3- عمق نفوذ، نوشته شد، برای محاسبه بهترین عمق نفوذ، ابتدا نتیجه را روی داده validation بررسی کرده و درجایی که بهترین نتیجه را بگیریم، روی داده تست اجرا میکنیم:

برای بررسی بهترین عمق نفوذ چند مثال را بررسی میکنیم:

مثال 1:

```
tree = id3(df_shroom, attribute_names,depth_of_tree=None)
```

که خروجی آن:

Accuracy of train is 1.0

Accuracy of validation is 0.7672770230360307

مثال 2:

```
tree = id3(df_shroom, attribute_names,depth_of_tree=10)
```

که نتایج زیر حاصل شد:

Accuracy of train is 0.9091714665485158

Accuracy of validation is 0.8470171293561725

مثال 3:

```
tree = id3(df_shroom, attribute_names,depth_of_tree=5)
```

که نتایج زیر حاصل شد:

Accuracy of train is 0.9020824102791316

Accuracy of validation is 0.8865918487891317

مثال 4:

```
tree = id3(df_shroom, attribute_names,depth_of_tree=3)
```

که نتایج زیر حاصل شد:

Accuracy of train is 0.898980948161276

Accuracy of validation is 0.8907265209686946

مثال 5:

```
tree = id3(df_shroom, attribute_names, depth_of_tree=2)
```

که نتایج زیر حاصل شد:

Accuracy of train is 0.898980948161276

Accuracy of validation is 0.8907265209686946

مثال 5:

```
tree = id3(df_shroom, attribute_names, depth_of_tree=1)
```

که نتایج زیر حاصل شد:

Accuracy of train is 0.8932210899424015

Accuracy of validation is 0.8753691671588896

همانگونه که مشاهده میشود ، برای درختی دو یا سه سطحی ، بهترین نتیجه را داریم ، در زیر درخت دو سطحی نمایش داده میشود :

```
{'odor': {'a': {'sporeprintcolor': {'k': 0, 'n': 0, 'u': 0}},  
         'c': {'sporeprintcolor': {'k': 1, 'n': 1}},  
         'f': {'sporeprintcolor': {'h': 1}},  
         'l': {'sporeprintcolor': {'k': 0, 'n': 0, 'u': 0}},  
         'm': {'sporeprintcolor': {'w': 1}},  
         'n': {'sporeprintcolor': {'k': 0, 'n': 0, 'r': 1, 'w': 0}},  
         'p': {'sporeprintcolor': {'k': 1, 'n': 1}}}}
```

حال به بررسی اطلاعات مربوط به دادگان تست میپردازیم :

که نتیجه زیر حاصل شد :

Accuracy of train is 0.898980948161276

Accuracy of validation is 0.8907265209686946

Accuracy of test is 0.9037780401416765

که درصد ملاحظه شده ای نسبت به درخت بدون هرس ، یا عمق نفوذ بهبود یافته

همچنین کد بخش آخر در ادامه آمده است :

```

1 import pandas as pd
2 from collections import Counter
3 from math import log2
4 from pprint import pprint
5
6 path='F:/Data/noisy_train.csv'
7 df_shroom=pd.read_csv(path)
8
9 path1='F:/Data/noisy_valid.csv'
10 df_shroom_valid=pd.read_csv(path1)
11
12 path2='F:/Data/noisy_test.csv'
13 df_shroom_test=pd.read_csv(path2)
14
15
16 def entropy(prob):
17     return sum( [-p*log2(p) for p in prob] )
18
19 def entropy_of_list(a_list):
20     cnt = Counter(x for x in a_list)
21
22     num_instances = len(a_list)*1.0
23     #mohasebe ehtemal ha
24     probs = [x / num_instances for x in cnt.values()]
25
26     return entropy(probs)
27
28
29 def gain(feature_name):
30     # dar avardan category haye featue
31     df_category = df_shroom.groupby(feature_name)
32
33     # mohasebe entropy ha dar dataframe roye target_attribute_name
34     n = len(df_shroom.index) * 1.0
35     df_agg_ent=df_category.agg({'poisonous':
36         [entropy_of_list, lambda x: len(x)/n] })['poisonous']
37
38
39     df_agg_ent.columns = ['entropy', 'prob']
40
41     #jamh entropy ha va zarb dar ehtemaleshon
42     return sum( df_agg_ent['entropy'] * df_agg_ent['prob'] )
43
44 # _____ ta in ja be nahvi dar code aval neveshte shod _____

```

```

46 def id3(df , feature_names , depth_of_tree, last_lable=None):
47     #baraye moshakhas kardan depth tree
48     #aya dade ya na
49     if(depth_of_tree!=None):
50         dpt=depth_of_tree-1
51     else:
52         dpt=None
53     #bara hame lable ha :
54     cnt = Counter(x for x in df['poisonous'])
55
56     #sharayet tvaghof tree
57     ##1:agr hame label ha az yek model bashand
58     if len(cnt) == 1:
59         a=list(cnt.keys())
60         return a[0]
61     ##2:agr feature ha tamam shode bashand ya az hameye data estefade karde bashim
62     elif df.empty or (not feature_names):
63         return last_lable
64     ###hala bia va derakh ra beshkan
65     else:
66         # label aksariyat in ja ra baraye laye badi dashte bash
67         #voting
68         a=list(cnt.values())
69         index_of_max = a.index(max(cnt.values()))
70         a=list(cnt.keys())
71         last_lable = a[index_of_max]
72         #mohasebe gain
73         #bia gain ha ra hesab kon va kamtrin ra bardar
74         #choon az entropy ghabl kam nakrdi kamtarin gain bishtarin data ra darad
75         gainz = [gain(feature) for feature in feature_names]
76         index_of_max_gain = gainz.index(min(gainz))
77
78         best_feature = feature_names[index_of_max_gain]
79
80     ##3:agr depth baraye tree bashad va be on reside bashad
81     if(dpt!=None):
82         if(dpt<0):
83             return last_lable
84
85     # ezafe kardan shakhe khali daroon dict ghabli
86     tree = {best_feature:{}}
87     remaining_feature_names=[i for i in feature_names if i != best_feature]
88
89
90     #seda kardan recursively in algorithm ta yeki az sharayet tavaghof rokh dahad
91     for feature_val, data_subset in df.groupby(best_feature):
92
93         subtree = id3(data_subset,remaining_feature_names,dpt)
94
95         tree[best_feature][feature_val] = subtree
96     return tree
97

```

```

7
3 attribute_names = list(df_shroom.columns)
3 attribute_names.remove('poisonous')
3
1 tree = id3(df_shroom, attribute_names,depth_of_tree=3)
2
3 pprint(tree)
4

```

```

#func komaki baraye conv tree be list
def tree_to_list(d_pandas, tree):

    a = list(tree.keys())
    feature = a[0]
    if d_pandas[feature] in tree[feature].keys():
        result = tree[feature][d_pandas[feature]]
        if isinstance(result, dict):
            # tosh ye dic dg hast : derakh bayad shekaste shavad
            return tree_to_list(d_pandas, result)
        else:
            return result # this is a label
    else:
        return None

```

```

125 #__test of train__
126
127 train_data = df_shroom
128 #applay kardan func tree_to_list be tree va rikhtan ro predicted_train
129 train_data['predicted_train'] = train_data.apply(tree_to_list,axis='columns',args=(tree,) )
130
131 print ( 'Accuracy of train is ' ,
132         sum(train_data['poisonous']==train_data['predicted_train'])/ (len(train_data.index)*1.0)
133         )
134 #__test of valid__
135
136 valid_data = df_shroom_valid
137 valid_data['predicted_valid'] = valid_data.apply(tree_to_list,axis='columns',args=(tree,) )
138
139 print ( 'Accuracy of validation is ' ,
140         sum(valid_data['poisonous']==valid_data['predicted_valid'])/ (len(valid_data.index)*1.0)
141         )
142
143 # test of test
144 test_data = df_shroom_test
145 test_data['predicted_test'] = test_data.apply(tree_to_list,axis='columns',args=(tree,) )
146
147 print ( 'Accuracy of test is ' ,
148         sum(test_data['poisonous']==test_data['predicted_test'])/ (len(test_data.index))
149         )
150 #

```