

## گزارش درخت تصمیم

علیرضا طباطبائی – الیاس بصیری

ابتدا لازم به ذکر است که دستورات اصلی دارای کامنت و تمامی توابع دارای توضیحات میباشند . لذا به انجام قسمت الف و ب در این گزارش میپردازیم .

جهت اجرای کدها میبایست فایل Main\_Code را اجرا نمایید .

بقیه فایل ها ، توابع نوشته شده و مورد استفاده میباشند .

Accuracy_Fcn.m	12/9/2021 6:52 PM	M File	1 KB
adult.test.10k	9/17/2003 11:59 PM	Text Document	954 KB
adult.train.10k	9/17/2003 11:59 PM	Text Document	953 KB
Choose_Data_Randomly.m	12/9/2021 7:21 PM	M File	1 KB
Classifier.m	12/9/2021 6:53 PM	M File	5 KB
Correlator_Fcn.m	12/9/2021 6:51 PM	M File	2 KB
Data_Read.m	12/9/2021 6:52 PM	M File	1 KB
Entropy_calculator.m	12/9/2021 6:52 PM	M File	1 KB
Find_Features.m	12/9/2021 6:47 PM	M File	1 KB
Find_IG.m	12/9/2021 6:53 PM	M File	1 KB
Grouping_Fcn.m	12/9/2021 6:47 PM	M File	1 KB
Label_Fcn.m	12/9/2021 6:46 PM	M File	1 KB
Main_Code.asv	12/9/2021 7:48 PM	ASV File	6 KB
Main_Code.m	12/9/2021 7:57 PM	M File	6 KB
Max_IG.m	12/9/2021 6:53 PM	M File	1 KB
Number_of_a_Value_of_a_Feature.m	12/9/2021 6:53 PM	M File	1 KB
Pruning_Fcn.m	12/9/2021 6:43 PM	M File	3 KB
Report_9723052	12/9/2021 8:13 PM	Microsoft Word D...	16 KB
Root_Finder.m	12/9/2021 6:53 PM	M File	9 KB
Roots_Features_Fcn.m	12/9/2021 6:28 PM	M File	1 KB
Sarnes.m	12/9/2021 6:28 PM	M File	1 KB
Size_of_Tree.m	12/9/2021 6:27 PM	M File	1 KB
Ungrouped_Label_Fcn.m	12/9/2021 6:54 PM	M File	1 KB
UnSame.m	12/9/2021 6:23 PM	M File	1 KB
Which_Feature.m	12/9/2021 6:54 PM	M File	1 KB

قسمت الف :

تمامی داده های آموزش را استفاده کرده و درخت را ساختیم . سپس با استفاده از داده های تست ، آن را تست کردیم.

در اینجا به علت استفاده از روش اول عمق ، سایز درخت را برابر تعداد مسیر ها تعریف میکنیم.

a ( در این گام ، 45 درصد دادگان آموزش را بصورت تصادفی به درخت دادیم و صحت را بر روی کل دادگان تست امتحان کردیم :

تکرار 1 : ( سائز درخت یا همان تعداد مسیر ها : 1041 )

صحت بر روی 45 درصد دادگان آموزش : 88.60 درصد

صحت بر روی کل دادگان تست : 80.10 درصد

تکرار 2 : ( سائز درخت یا همان تعداد مسیر ها : 1015 )

صحت بر روی 45 درصد دادگان آموزش : 88.98 درصد

صحت بر روی کل دادگان تست : 80.33 درصد

تکرار 3 : ( سائز درخت یا همان تعداد مسیر ها : 1045 )

صحت بر روی 45 درصد دادگان آموزش : 88.24 درصد

صحت بر روی کل دادگان تست : 80.69 درصد

میانگین 3 تکرار : ( سائز درخت یا همان تعداد مسیر ها : 1033 )

صحت بر روی 45 درصد دادگان آموزش : 88.60 درصد

صحت بر روی کل دادگان تست : 80.37 درصد

b ( استفاده از 55 درصد دادگان آموزش :

تکرار 1 : ( سائز درخت یا همان تعداد مسیر ها : 1179 )

صحت بر روی 55 درصد دادگان آموزش : 88.73 درصد

صحت بر روی کل دادگان تست : 80.81 درصد

تکرار 2 : ( سائز درخت یا همان تعداد مسیر ها : 1184 )

صحت بر روی 55 درصد دادگان آموزش : 88.44 درصد

صحت بر روی کل دادگان تست : 80.31 درصد

تکرار 3 : ( سائز درخت یا همان تعداد مسیر ها : 1192 )

صحت بر روی 55 درصد دادگان آموزش : 88.11 درصد

صحت بر روی کل دادگان تست : 81.04 درصد

میانگین 3 تکرار : ( سائز درخت یا همان تعداد مسیر ها : 1185 )

صحت بر روی 55 درصد دادگان آموزش : 88.42 درصد

صحت بر روی کل دادگان تست : 80.72 درصد

b ( استفاده از 65 درصد دادگان آموزش :

تکرار 1 : ( سائز درخت یا همان تعداد مسیر ها : 1299 )

صحت بر روی 65 درصد دادگان آموزش : 88.12 درصد

صحت بر روی کل دادگان تست : 80.24 درصد

تکرار 2 : ( سائز درخت یا همان تعداد مسیر ها : 1375 )

صحت بر روی 65 درصد دادگان آموزش : 88.49 درصد

صحت بر روی کل دادگان تست : 80.56 درصد

تکرار 3 : ( سائز درخت یا همان تعداد مسیر ها : 1386 )

صحت بر روی 65 درصد دادگان آموزش : 88.12 درصد

صحت بر روی کل دادگان تست : 80.90 درصد

میانگین 3 تکرار : ( سائز درخت یا همان تعداد مسیر ها : 1353 )

صحت بر روی 65 درصد دادگان آموزش : 88.24 درصد

صحت بر روی کل دادگان تست : 80.56 درصد

b ( استفاده از 75 درصد دادگان آموزش :

تکرار 1 : ( سائز درخت یا همان تعداد مسیر ها : 1506 )

صحت بر روی 75 درصد دادگان آموزش : 88.05 درصد

صحت بر روی کل دادگان تست : 80.72 درصد

تکرار 2 : ( سائز درخت یا همان تعداد مسیر ها : 1511 )

صحت بر روی 75 درصد دادگان آموزش : 87.85 درصد

صحت بر روی کل دادگان تست : 80.96 درصد

تکرار 3 : ( سائز درخت یا همان تعداد مسیر ها : 1581 )

صحت بر روی 75 درصد دادگان آموزش : 88.01 درصد

صحت بر روی کل دادگان تست : 80.81 درصد

میانگین 3 تکرار : ( سائز درخت یا همان تعداد مسیر ها : 1532 )

صحت بر روی 75 درصد دادگان آموزش : 87.97 درصد

صحت بر روی کل دادگان تست : 80.83 درصد

b ( استفاده از 100 درصد دادگان آموزش :

تکرار 1 : ( سائز درخت یا همان تعداد مسیر ها : 1898 )

صحت بر روی 100 درصد دادگان آموزش : 87.54 درصد

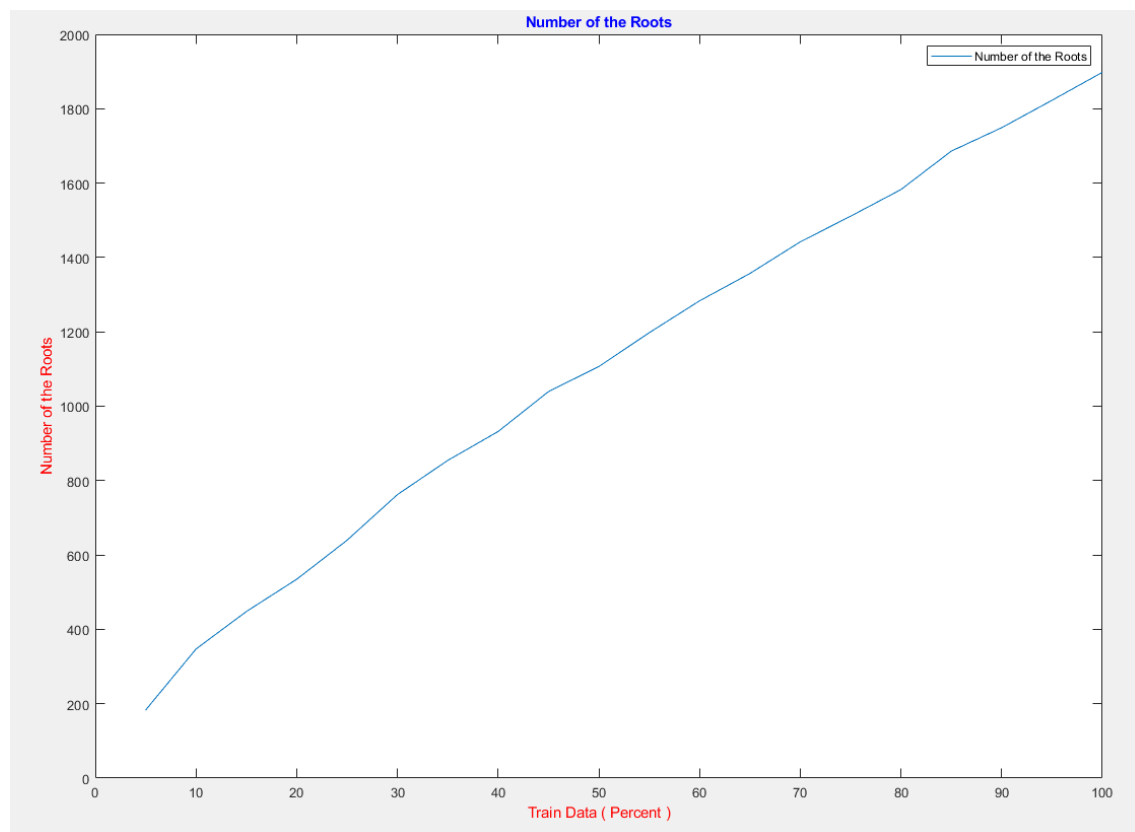
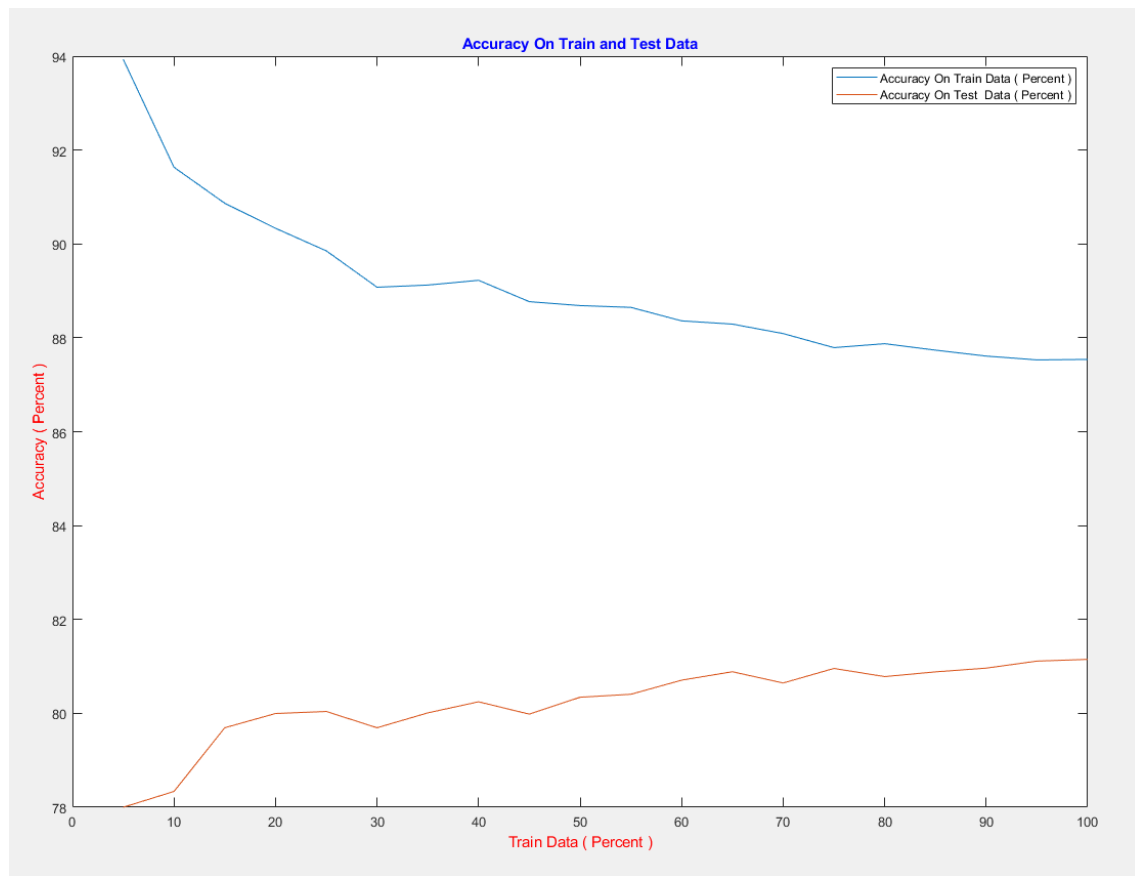
صحت بر روی کل دادگان تست : 81.15 درصد

بحث بر روی تعداد دادگان آموزش و تاثیر آن در صحت و سائز درخت :

هر چقدر تعداد دادگان آموزش بیشتر شود ، صحت بر روی مجموعه آموزش کمتر میشود زیرا هنگامی که داده ها کم باشد ، درخت به راحتی overfit میشود اما با دادگان زیاد ، دیگر نمیتواند overfit شود و در نتیجه صحت کمتری بر روی دادگان آموزش دارد .

همچنین با افزایش دادگان آموزش ، صحت بر روی مجموعه تست صعودی است .

افزایش دادگان آموزش به صورت کاملاً صعودی باعث افزایش طول درخت میشود .



## قسمت ب : Pruning

اجرای هر قسمت این سوال حدود 10 ساعت نیاز به Run شدن داشت که واقعا فرصت این کار را نداشتم ولی کد ها کامل نوشته شده است و بدون مشکل Run میشود:

( a

```
Main_Code.m x +
104
105 - for i = 1 : 50 % Prune the tree
106 -     [Roots , Test_Accuracy , Labels] = Pruning_Fcn(Data_Test,Roots,i,Features,Test_Accuracy,La
107 -     [Grouped_Data_Test , UnGrouped_Data_Test] = Grouping_Fcn(Roots,Data_Test_two,Features) ;
108 -     Roots_Features = Roots_Features_Fcn(Roots,Features) ;
109 -     [Number_of_the_Roots , ~] = size(Roots) ;
110
111 -     for k = 1 : Number_of_the_Roots
112 -         Grouped_Data_Test(k,2) = Labels(k,1) ;
113 -     end
114
115 -     [Row_UnGrouped_Data_Test,~] = size(UnGrouped_Data_Test) ;
116
117 -     for k = 1 : Row_UnGrouped_Data_Test
118 -         for j = 1 : 9
119 -             Data{1,j} = UnGrouped_Data_Test{k,j} ;
120 -         end
121 -         UnGrouped_Data_Test{k,Number_of_features+1} = Ungrouped_Label_Fcn(Data,Roots,Roots_Fea
122 -     end
123
124 -     Test_Accuracy_two = Accuracy_Fcn(Grouped_Data_Test,UnGrouped_Data_Test,My_Label_Column) ;
125
126 -     i
127 -     Test_Accuracy
128 -     Test_Accuracy_two
129 -     A(1,i) = Test_Accuracy ;|
130 -     A(2,i) = Test_Accuracy_two ;
131
132 - end
133
```



بحث در مورد Overfitting و تاثیر Pruning :

با هرس کردن درخت ، از overfit شدن آن جلوگیری میشود و باعث کاهش صحت بر روی مجموعه آموزش و افزایش صحت در داده های تست میشود که این امر بسیار مطلوب است .