

Q5 - Naive Bayes and Decision Trees

Import the libraries

```
1 from sklearn.naive_bayes import MultinomialNB
2 from IPython.display import display
3 from IPython.display import SVG
4 from id3 import export_graphviz
5 from id3 import Id3Estimator
6 from sklearn import tree
7
8 import matplotlib.pyplot as plt
9 import numpy as np
10 import graphviz
11 import os
12 %matplotlib inline
```

Load the data

```
1 def data_and_headers(filename):
2     data = None
3     with open(filename) as fp:
4         data = [x.strip().split(',') for x in fp.readlines()]
5         headers = data[0]
6         headers = np.asarray(headers)
7         class_field = len(headers) - 1
8         data_x = [[x[i] for i in range(class_field)] for x in data[1:]]
9         data_x = np.asarray(data_x)
10        data_y = [[x[i] for i in range(class_field, class_field + 1)] for x in
11                    data[1:]]
12        data_y = np.asarray(data_y)
13        return headers, data_x, data_y
```

```
1 headers, X, Y = data_and_headers('Data' + os.sep + 'hw2q5.csv')
2 indexes=[int(x) for x in list(X[:,0])]
3 X=X[:,1:]
```

(A) K-Fold Splits

```

1 def createfolds(indexes):
2     folds = {i:{'train':[], 'test':[]}} for i in range(1,6)}
3     for i in range(len(indexes)):
4         for j in range(1,6):
5             if indexes[i] % 5 == j-1:
6                 folds[j]['test'].append(i)
7             else:
8                 folds[j]['train'].append(i)
9     return folds

```

```

1 folds = createfolds(indexes)

```

Naive Bayes

```

1 X = X.tolist()
2 Y = np.ravel(Y).tolist()
3 d1={'presbyopic':2, 'pre-presbyopic':1, 'young':0, 'myope':0,
4     'hypermetropia':1, 'no':0, 'yes':1, 'normal':1, 'reduced':0}
5 d2={'Yes':1, 'No':0}
6 X = [[d1[X[i][j]] for j in range(len(X[0]))] for i in range(len(X))]
7 X = np.asarray(X)
8 Y = [d2[Y[i]] for i in range(len(Y))]
9 Y = np.asarray(Y)

```

```

1 cnt = 0
2 subheaders = headers[1:-1]
3 for i in sorted(folds.keys()):
4     print('Fold '+str(i))
5     nb = MultinomialNB(alpha=1)
6     nb=nb.fit(X[folds[i]['train']],Y[folds[i]['train']])
7     ypred=nb.predict(X[folds[i]['test'],:])
8     print('\tTest IID -\t' + ', '.join([str(x) for x in
9     np.asarray(indexes)[folds[i]['test']]))
10    print('\tActual -\t'+', '.join(['Yes' if x==1 else 'No' for x in
11    Y[folds[i]['test']]))
12    print('\tPredict -\t' + ', '.join(['Yes' if x==1 else 'No' for x in
13    ypred]))
14    for j in range(len(ypred)):
15        if ypred[j]!=Y[folds[i]['test']][j]:
16            cnt+=1
17    print('\tProbabilities - ')
18    dt={0:'Yes', 1:'No'}
19    for i in range(len(nb.feature_log_prob_)):

```

```

17         for j in range(len(subheaders)):
18             print('\t\tP({}|Class={}) = {:.3f}'.format(subheaders[j],
dt[i], np.exp(nb.feature_log_prob_)[i][j]))
19 print('\nNaive-Bayes 5-fold CV accuracy - ' + str((24-cnt)*100/24) + '%')

```

```

1 Fold 1
2     Test IID - 5, 10, 15, 20
3     Actual - No, Yes, No, Yes
4     Predict - No, No, No, No
5     Probabilities -
6         P(patient age|Class=Yes) = 0.472
7         P(spectacle prescription|Class=Yes) = 0.194
8         P(astigmatic|Class=Yes) = 0.222
9         P(tear production rate|Class=Yes) = 0.111
10        P(patient age|Class=No) = 0.227
11        P(spectacle prescription|Class=No) = 0.227
12        P(astigmatic|Class=No) = 0.182
13        P(tear production rate|Class=No) = 0.364
14 Fold 2
15     Test IID - 1, 6, 11, 16, 21
16     Actual - No, Yes, No, No, No
17     Predict - No, Yes, No, Yes, No
18     Probabilities -
19         P(patient age|Class=Yes) = 0.452
20         P(spectacle prescription|Class=Yes) = 0.226
21         P(astigmatic|Class=Yes) = 0.226
22         P(tear production rate|Class=Yes) = 0.097
23         P(patient age|Class=No) = 0.308
24         P(spectacle prescription|Class=No) = 0.154
25         P(astigmatic|Class=No) = 0.192
26         P(tear production rate|Class=No) = 0.346
27 Fold 3
28     Test IID - 2, 7, 12, 17, 22
29     Actual - Yes, No, Yes, No, Yes
30     Predict - Yes, No, No, No, No
31     Probabilities -
32         P(patient age|Class=Yes) = 0.444
33         P(spectacle prescription|Class=Yes) = 0.222
34         P(astigmatic|Class=Yes) = 0.222
35         P(tear production rate|Class=Yes) = 0.111
36         P(patient age|Class=No) = 0.250
37         P(spectacle prescription|Class=No) = 0.200
38         P(astigmatic|Class=No) = 0.200
39         P(tear production rate|Class=No) = 0.350
40 Fold 4

```

```

41     Test IID -   3, 8, 13, 18, 23
42     Actual -    No, Yes, No, No, No
43     Predict -   No, Yes, No, Yes, No
44     Probabilities -
45         P(patient age|Class=Yes) = 0.433
46         P(spectacle prescription|Class=Yes) = 0.233
47         P(astigmatic|Class=Yes) = 0.233
48         P(tear production rate|Class=Yes) = 0.100
49         P(patient age|Class=No) = 0.320
50         P(spectacle prescription|Class=No) = 0.160
51         P(astigmatic|Class=No) = 0.160
52         P(tear production rate|Class=No) = 0.360
53     Fold 5
54     Test IID -   4, 9, 14, 19, 24
55     Actual -    Yes, No, Yes, No, No
56     Predict -   Yes, No, Yes, No, No
57     Probabilities -
58         P(patient age|Class=Yes) = 0.419
59         P(spectacle prescription|Class=Yes) = 0.258
60         P(astigmatic|Class=Yes) = 0.226
61         P(tear production rate|Class=Yes) = 0.097
62         P(patient age|Class=No) = 0.304
63         P(spectacle prescription|Class=No) = 0.174
64         P(astigmatic|Class=No) = 0.174
65         P(tear production rate|Class=No) = 0.348
66
67     Naive-Bayes 5-fold CV accuracy - 75.0%

```

Decision Tree

```

1 headers, X, Y = data_and_headers('Data' + os.sep + 'hw2q5.csv')
2 indexes=[int(x) for x in list(X[:,0])]
3 X=X[:,1:]

```

```

1 cnt = 0
2 subheaders = headers[1:-1]
3 for i in sorted(folds.keys()):
4     print('Fold '+str(i))
5     #dt = tree.DecisionTreeClassifier(criterion='entropy',
6     splitter='best')
7     dt = Id3Estimator(gain_ratio=True)
8     dt = dt.fit(X[folds[i]['train'],:],Y[folds[i]['train']])
9     ypred=dt.predict(X[folds[i]['test'],:])

```

```

9     print('\tTest IID -\t' + ', '.join([str(x) for x in
np.asarray(indexes)[folds[i]['test']]))
10    print('\tActual -\t'+', '.join(np.ravel(Y[folds[i]['test']]))
11    print('\tPredict -\t' + ', '.join(ypred))
12    #     print('\tActual -\t'+', '.join(['Yes' if x==1 else 'No' for x in
Y[folds[i]['test']]))
13    #     print('\tPredict -\t' + ', '.join(['Yes' if x==1 else 'No' for x in
ypred]))
14    for j in range(len(ypred)):
15        if ypred[j]!=Y[folds[i]['test']][j]:
16            cnt+=1
17    dot_data = export_graphviz(dt.tree_,
'fold'+str(i)+'.dot',feature_names = subheaders)
18    print('\nDecision Tree 5-fold CV accuracy - '+ str((24-cnt)*100/24) + '%')

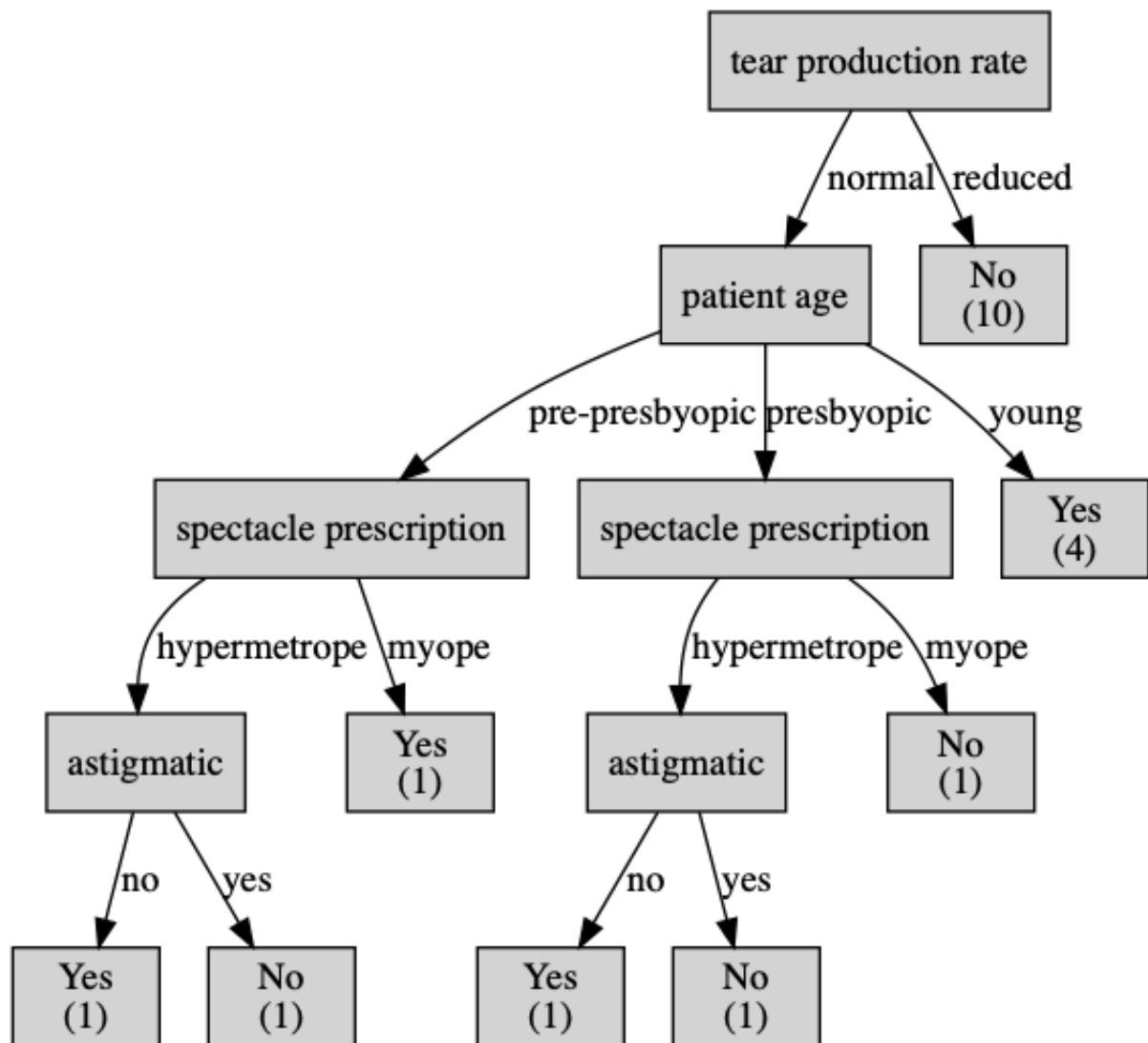
```

```

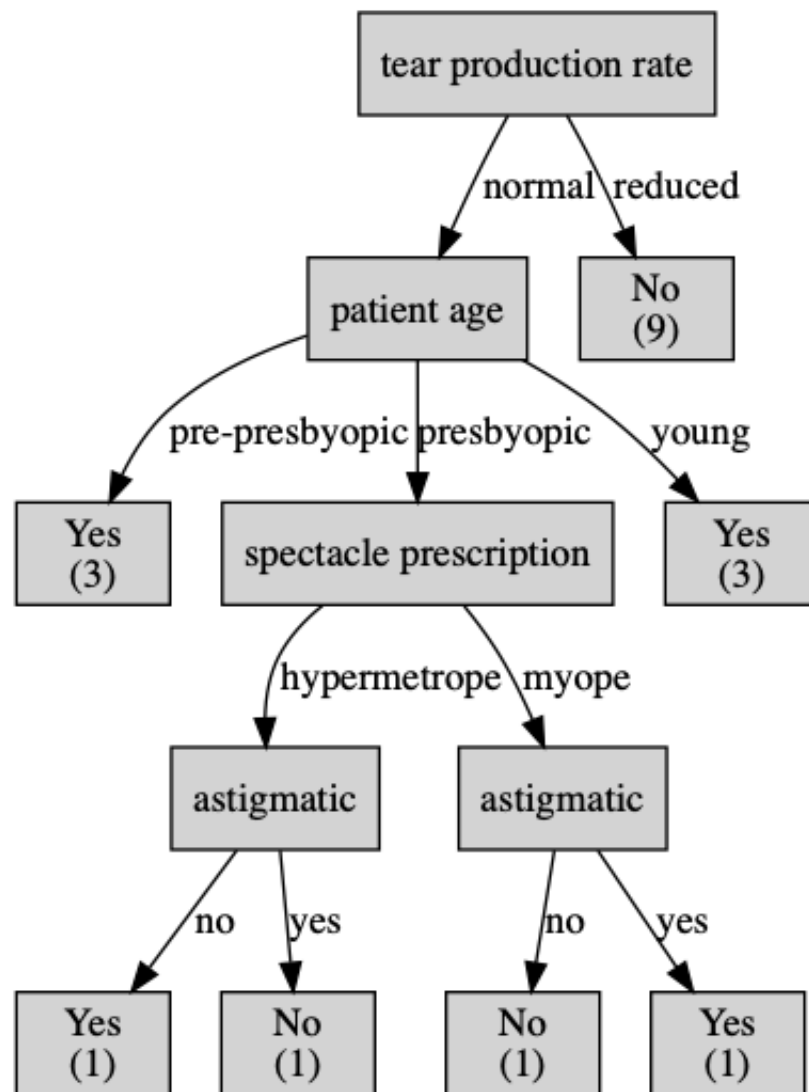
1  Fold 1
2      Test IID -   5, 10, 15, 20
3      Actual -    No, Yes, No, Yes
4      Predict -   No, Yes, No, No
5  Fold 2
6      Test IID -   1, 6, 11, 16, 21
7      Actual -    No, Yes, No, No, No
8      Predict -   No, Yes, No, Yes, No
9  Fold 3
10     Test IID -   2, 7, 12, 17, 22
11     Actual -    Yes, No, Yes, No, Yes
12     Predict -   Yes, No, No, No, No
13  Fold 4
14     Test IID -   3, 8, 13, 18, 23
15     Actual -    No, Yes, No, No, No
16     Predict -   No, No, No, Yes, No
17  Fold 5
18     Test IID -   4, 9, 14, 19, 24
19     Actual -    Yes, No, Yes, No, No
20     Predict -   Yes, No, No, No, Yes
21
22  Decision Tree 5-fold CV accuracy - 66.6666666666667%

```

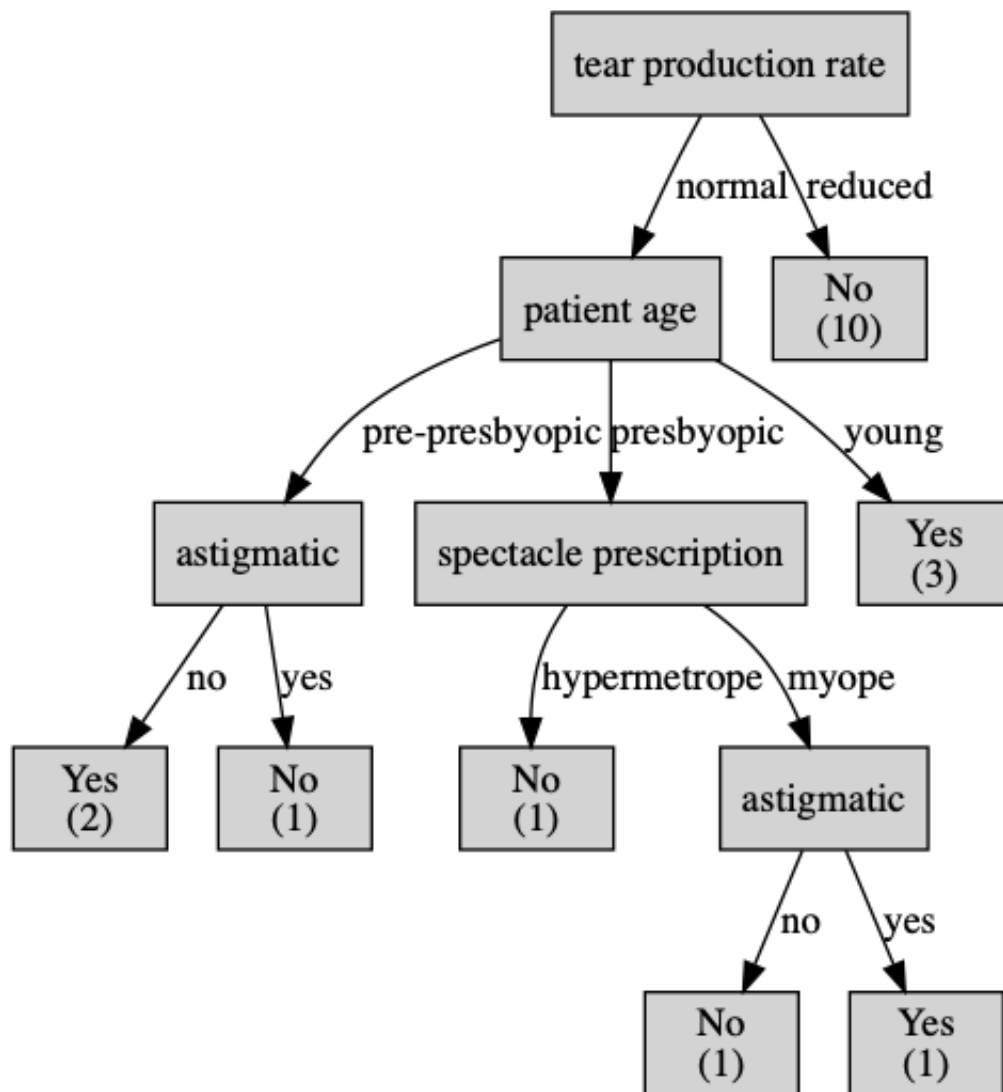
Fold 1



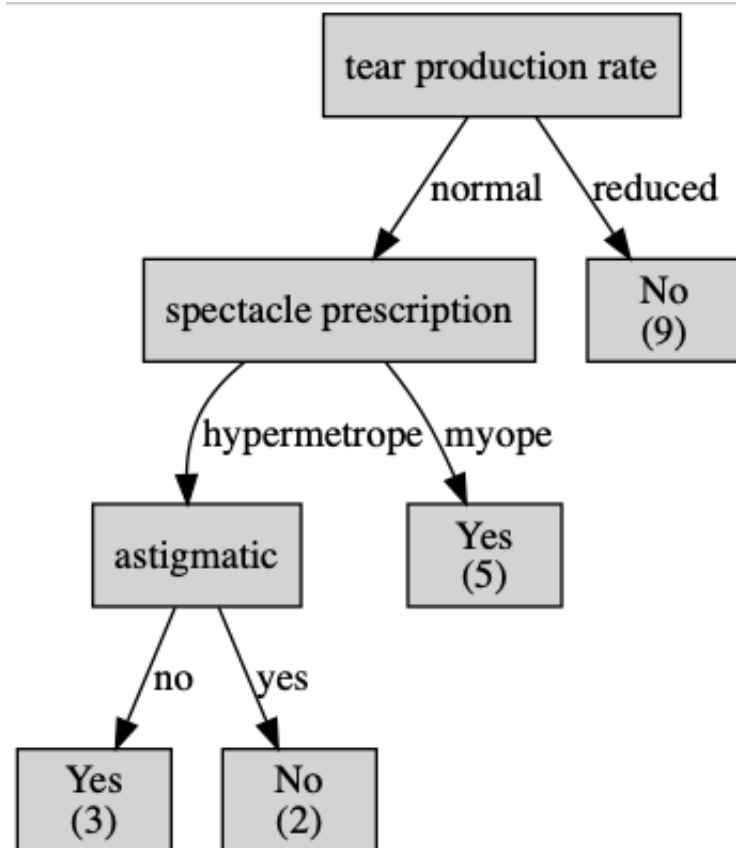
Fold 2



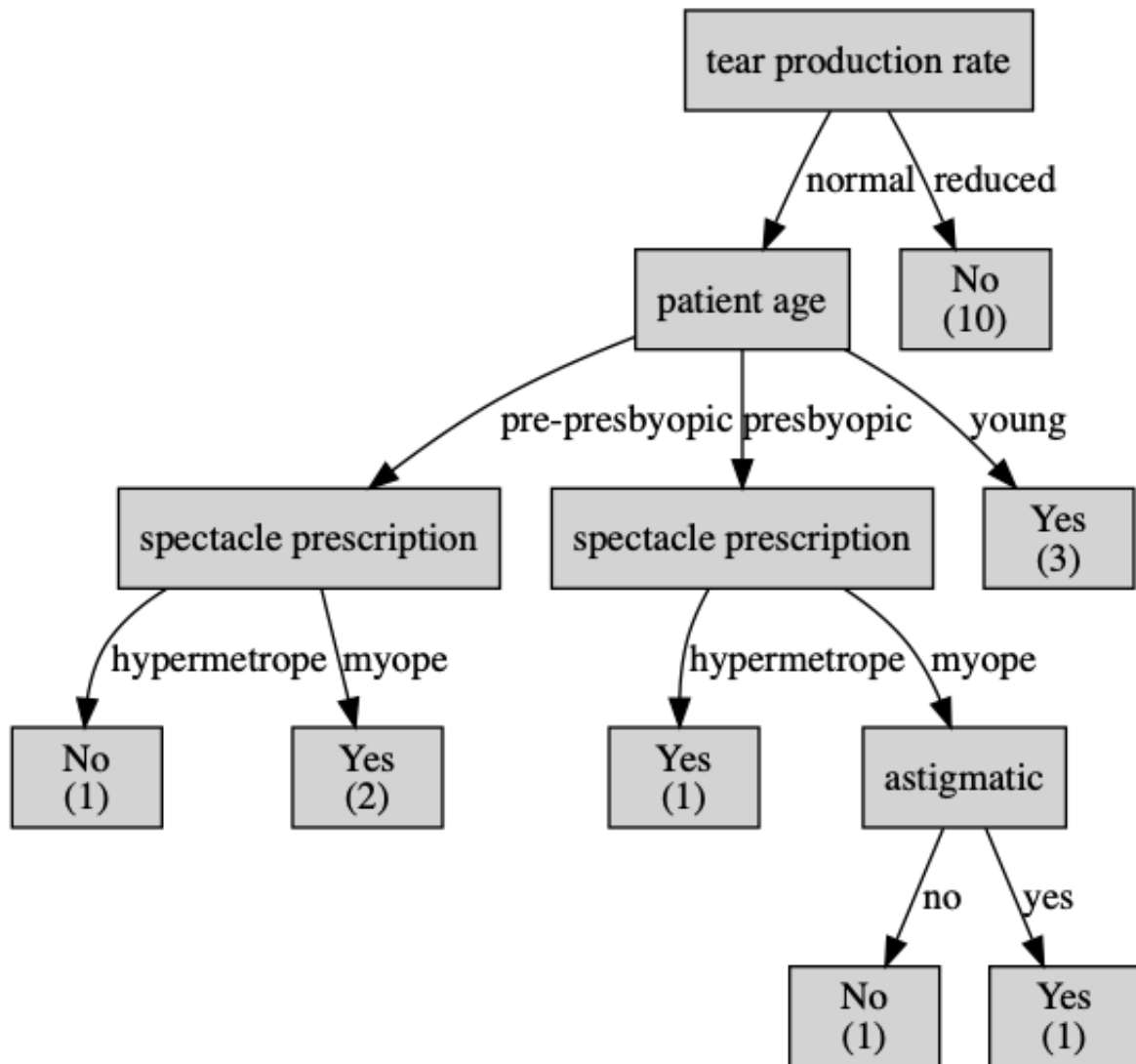
Fold 3



Fold 4



Fold 5



(B) Choosing Model

Based on the above 5-fold CV accuracy, it seems that Naive Bayes is the better model for this dataset.

Naive Bayes for Full Data

```

1 X = X.tolist()
2 Y = np.ravel(Y).tolist()
3 d1={'presbyopic':2, 'pre-presbyopic':1, 'young':0, 'myope':0,
4     'hypermetrope':1, 'no':0, 'yes':1, 'normal':1, 'reduced':0}
5 d2={'Yes':1, 'No':0}
6 X = [[d1[X[i][j]] for j in range(len(X[0]))] for i in range(len(X))]
7 X = np.asarray(X)
8 Y = [d2[Y[i]] for i in range(len(Y))]
9 Y = np.asarray(Y)

```

```

1 nb = MultinomialNB(alpha=1)
2 nb=nb.fit(X,Y)
3 print('Final Model for Naive Bayes')
4 print('Probabilities - ')
5 dt={0:'Yes', 1:'No'}
6 for i in range(len(nb.feature_log_prob_)):
7     for j in range(len(subheaders)):
8         print('\tP({}|Class={}) = {:.3f}'.format(subheaders[j], dt[i],
np.exp(nb.feature_log_prob_)[i][j]))

```

Model Details -

```

1 Final Model for Naive Bayes
2 Probabilities -
3     P(patient age|Class=Yes) = 0.450
4     P(spectacle prescription|Class=Yes) = 0.225
5     P(astigmatic|Class=Yes) = 0.225
6     P(tear production rate|Class=Yes) = 0.100
7     P(patient age|Class=No) = 0.286
8     P(spectacle prescription|Class=No) = 0.179
9     P(astigmatic|Class=No) = 0.179
10    P(tear production rate|Class=No) = 0.357

```