

Aman Chauhan - achauha3

Khantil Choksi - khchoksi

Q1 - PCA

Import the libraries

```
1 from sklearn.neighbors import KNeighborsClassifier  
2 import matplotlib.pyplot as plt  
3 import numpy as np  
4 import os  
5 %matplotlib inline
```

(A) 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 = [[float(x[i]) for i in range(class_field)] for x in data[1:]]  
9     data_x = np.asarray(data_x)  
10    data_y = [[str(x[i]) for i in range(class_field, class_field + 1)] for  
11        x in data[1:]]  
12    data_y = np.asarray(data_y)  
13    return headers, data_x, data_y
```

```
1 headers, train_x, train_y = data_and_headers('Data' + os.sep +  
2 'hw2q1_train.csv')  
2 headers, test_x, test_y = data_and_headers('Data' + os.sep +  
3 'hw2q1_test.csv')
```

```

1 print('Training Data')
2 print('Number of features - ' + str(train_x.shape[1]))
3 print('Number of target features - ' + str(train_y.shape[1]))
4 print('Number of observations - ' + str(train_x.shape[0]))
5 print('Number of observations in category R - ' +
6     str(train_y[train_y=='R'].shape[0]))
7 print('Number of observations in category M - ' +
8     str(train_y[train_y=='M'].shape[0]))
9 print()
10 print('Testing Data')
11 print('Number of features - ' + str(test_x.shape[1]))
12 print('Number of target features - ' + str(test_y.shape[1]))
13 print('Number of observations - ' + str(test_x.shape[0]))
14 print('Number of observations in category R - ' +
15     str(test_y[test_y=='R'].shape[0]))
16 print('Number of observations in category M - ' +
17     str(test_y[test_y=='M'].shape[0]))

```

```

1 Training Data
2 Number of features - 60
3 Number of target features - 1
4 Number of observations - 156
5 Number of observations in category R - 73
6 Number of observations in category M - 83
7
8 Testing Data
9 Number of features - 60
10 Number of target features - 1
11 Number of observations - 52
12 Number of observations in category R - 24
13 Number of observations in category M - 28

```

(B) Normalization and PCA

```

1 def normalize(data, minima, maxima):
2     normal = np.copy(data)
3     normal = (normal - minima) / (maxima - minima)
4     return normal

```

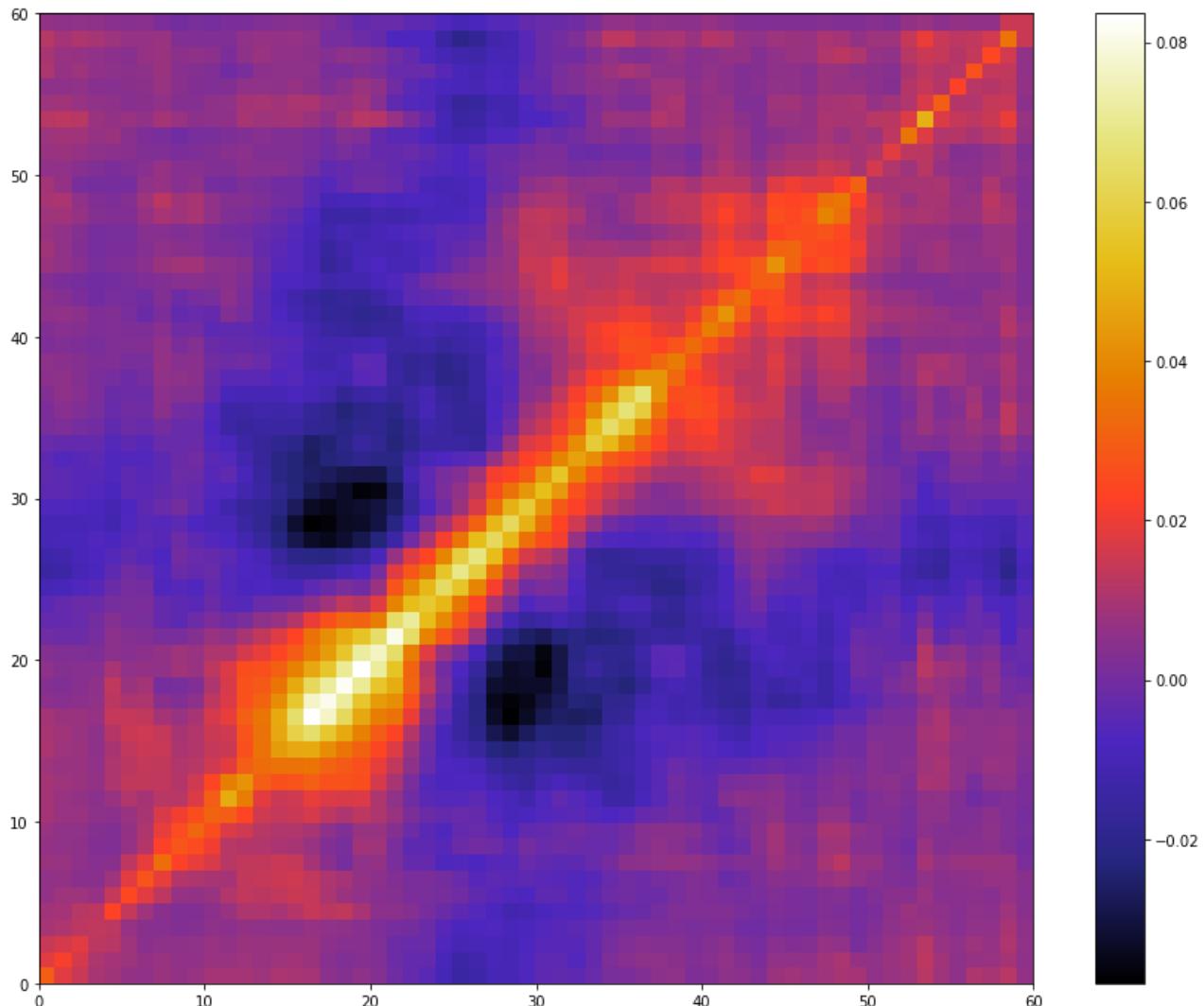
```

1 normal_train = normalize(train_x, np.amin(train_x, axis=0), np.amax(train_x,
2 axis=0))
3 normal_test = normalize(test_x, np.amin(train_x, axis=0), np.amax(train_x,
4 axis=0))

```

(i) Covariance of Training Dataset

```
1 covariance = np.cov(normal_train, rowvar=False)
2 fig, axes = plt.subplots(nrows=1, ncols=1)
3 fig.set_figheight(12)
4 fig.set_figwidth(15)
5 im = axes.pcolor(covariance, cmap='CMRmap')
6 fig.colorbar(im, ax=axes)
7 plt.show()
```



(ii) Eigenvalue and Eigenvectors

```
1 print('Size of covariance matrix - ' + str(covariance.shape))
2 w,v = np.linalg.eig(covariance)
3 print('Top 5 Eigenvalues - ' + ', '.join(['{:3f}'.format(x) for x in
w[:5]]))
```

```

1 | Size of covariance matrix - (60, 60)
2 | Top 5 Eigenvalues - 0.695, 0.457, 0.223, 0.191, 0.118

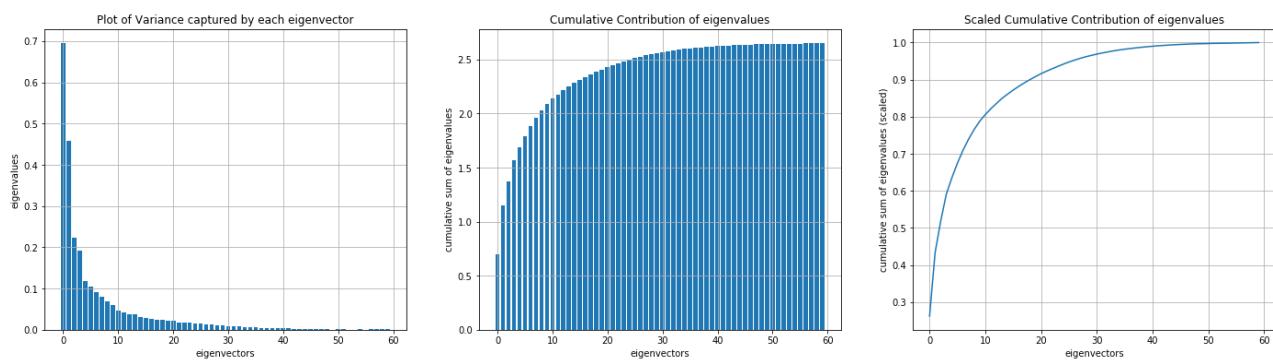
```

(iii) Plot of Eigenvalues

```

1 | fig, axes = plt.subplots(nrows=1, ncols=3)
2 | fig.set_figheight(6)
3 | fig.set_figwidth(24)
4 | axes[0].bar(np.arange(60), w)
5 | axes[0].xaxis.grid()
6 | axes[0].yaxis.grid()
7 | axes[0].set_xlabel('eigenvectors')
8 | axes[0].set_ylabel('eigenvalues')
9 | axes[0].set_title('Plot of Variance captured by each eigenvector')
10 | axes[1].bar(np.arange(60), np.cumsum(w))
11 | axes[1].xaxis.grid()
12 | axes[1].yaxis.grid()
13 | axes[1].set_xlabel('eigenvectors')
14 | axes[1].set_ylabel('cumulative sum of eigenvalues')
15 | axes[1].set_title('Cumulative Contribution of eigenvalues')
16 | axes[2].plot(np.arange(60), np.cumsum(w)/np.sum(w))
17 | axes[2].xaxis.grid()
18 | axes[2].yaxis.grid()
19 | axes[2].set_xlabel('eigenvectors')
20 | axes[2].set_ylabel('cumulative sum of eigenvalues (scaled)')
21 | axes[2].set_title('Scaled Cumulative Contribution of eigenvalues')
22 | plt.show()

```



Number of Eigenvectors - 10

This is because if we look at the scaled cumulative contribution of the eigen values, we can see that the first 10 eigenvalues cover almost 80% of all variance in the data set. So 10 seems to be a good choice for selecting the number of eigenvectors.

(iv) PCA with KNN

```

1 ncomps = [2,4,6,8,10,20,40,60]
2 cls = [KNeighborsClassifier(n_neighbors=3,
metric='euclidean').fit(np.matmul(normal_train, v[:, :ncomps[i]]),
np.ravel(train_y)) for i in range(len(ncomps))]

```

```

1 pred4 = cls[4].predict(np.matmul(normal_test, v[:, :ncomps[4]]))
2 true4 = np.ravel(test_y)
3 with open('q1iv.csv', 'w') as fp:
4     fp.write(', '.join(['Component' + str(i) for i in
range(1,ncomps[4]+1)])+', Class, Class\n')
5     transformed = np.matmul(normal_test, v[:, :ncomps[4]])
6     for i in range(len(transformed)):
7         fp.write(', '.join([str(x) for x in transformed[i]]))
8         fp.write(', ' + str(true4[i]))
9         fp.write(', ' + str(pred4[i]) + '\n')

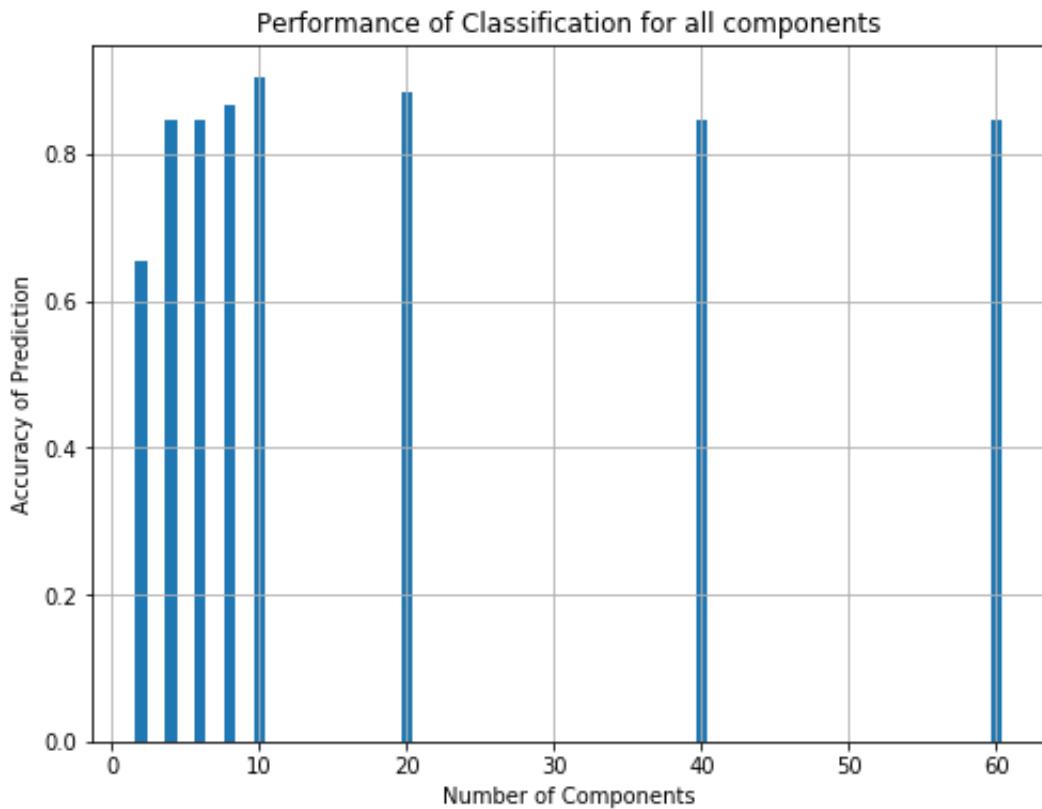
```

You can find the output of the above code in file q1iv.csv

```

1 fig, axes = plt.subplots(nrows=1, ncols=1)
2 fig.set_figheight(6)
3 fig.set_figwidth(8)
4 accuracies = [cls[i].score(np.matmul(normal_test, v[:, :ncomps[i]]),
np.ravel(test_y)) for i in range(len(ncomps))]
5 axes.bar(ncomps, accuracies)
6 axes.set_xlabel('Number of Components')
7 axes.set_ylabel('Accuracy of Prediction')
8 axes.set_title('Performance of Classification for all components')
9 plt.grid()
10 plt.show()
11 print('Accuracy - ')
12 print('Components\tAccuracy')
13 for i in range(len(ncomps)):
14     print('{}\t\t{:.4f}%'.format(ncomps[i], 100*accuracies[i]))

```



```

1 Accuracy -
2 Components Accuracy
3 2      65.3846%
4 4      84.6154%
5 6      84.6154%
6 8      86.5385%
7 10     90.3846%
8 20     88.4615%
9 40     84.6154%
10 60    84.6154%

```

Reasonable Number of PCA components - 10

The reasons are 2-fold - 1) We predicted it above when we plotted the eigenvalues and 2) It gives the best accuracy out of the given number of components.

(C) Standardization and PCA

```

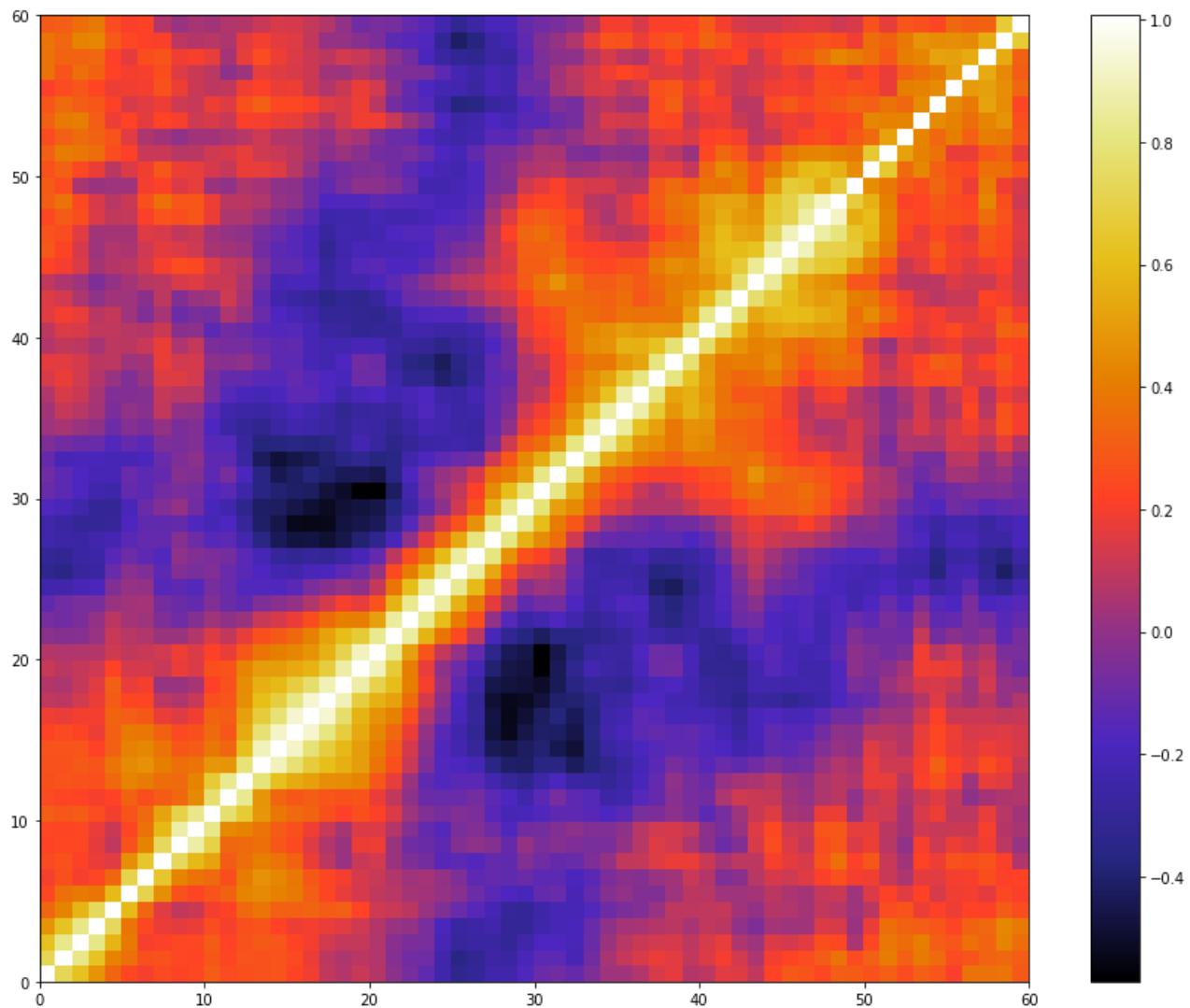
1 def standardize(data, mean, sd):
2     standard = np.copy(data)
3     standard = (standard - mean) / sd
4     return standard

```

```
1 standard_train = standardize(train_x, np.mean(train_x, axis=0),  
2 np.std(train_x, axis=0))  
3 standard_test = standardize(test_x, np.mean(train_x, axis=0),  
4 np.std(train_x, axis=0))
```

(i) Covariance of Training Dataset

```
1 covariance = np.cov(standard_train, rowvar=False)  
2 fig, axes = plt.subplots(nrows=1, ncols=1)  
3 fig.set_figheight(12)  
4 fig.set_figwidth(15)  
5 im = axes.pcolor(covariance, cmap='CMRmap')  
6 fig.colorbar(im, ax=axes)  
7 plt.show()
```



(ii) Eigenvalue and Eigenvectors

```

1 print('Size of covariance matrix - ' + str(covariance.shape))
2 w,v = np.linalg.eig(covariance)
3 print('Top 5 Eigenvalues - ' + ', '.join(['{:,.3f}'.format(x) for x in
w[:5]]))

```

```

1 Size of covariance matrix - (60, 60)
2 Top 5 Eigenvalues - 12.429, 11.558, 4.979, 3.345, 3.226

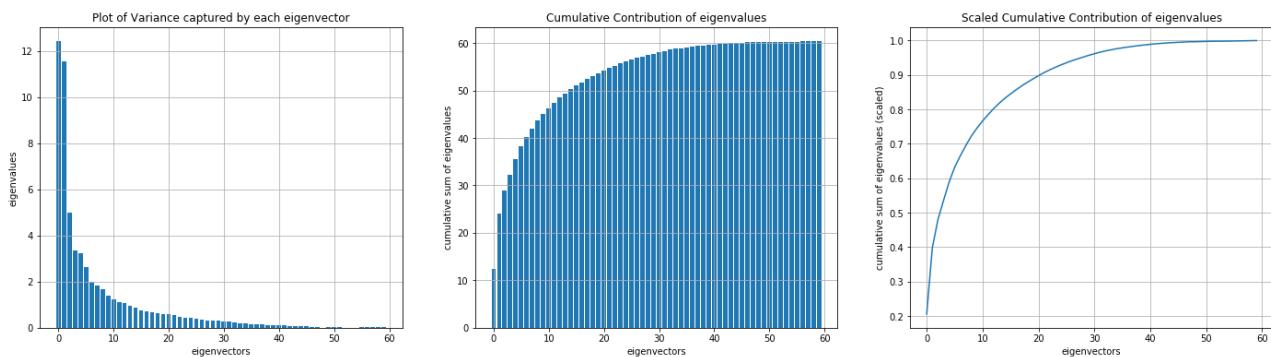
```

(iii) Plot of Eigenvalues

```

1 fig, axes = plt.subplots(nrows=1, ncols=3)
2 fig.set_figheight(6)
3 fig.set_figwidth(24)
4 axes[0].bar(np.arange(60), w)
5 axes[0].xaxis.grid()
6 axes[0].yaxis.grid()
7 axes[0].set_xlabel('eigenvectors')
8 axes[0].set_ylabel('eigenvalues')
9 axes[0].set_title('Plot of Variance captured by each eigenvector')
10 axes[1].bar(np.arange(60), np.cumsum(w))
11 axes[1].xaxis.grid()
12 axes[1].yaxis.grid()
13 axes[1].set_xlabel('eigenvectors')
14 axes[1].set_ylabel('cumulative sum of eigenvalues')
15 axes[1].set_title('Cumulative Contribution of eigenvalues')
16 axes[2].plot(np.arange(60), np.cumsum(w)/np.sum(w))
17 axes[2].xaxis.grid()
18 axes[2].yaxis.grid()
19 axes[2].set_xlabel('eigenvectors')
20 axes[2].set_ylabel('cumulative sum of eigenvalues (scaled)')
21 axes[2].set_title('Scaled Cumulative Contribution of eigenvalues')
22 plt.show()

```



Number of Eigenvectors - 20

This is because if we look at the scaled cumulative contribution of the eigen values, we can see that the first 20 eigenvalues cover almost 90% of all variance in the data set. Also, beyond that, all the eigenvectors combined can only contribute 10% more variance. So 20 seems to be a good choice for selecting the number of eigenvectors.

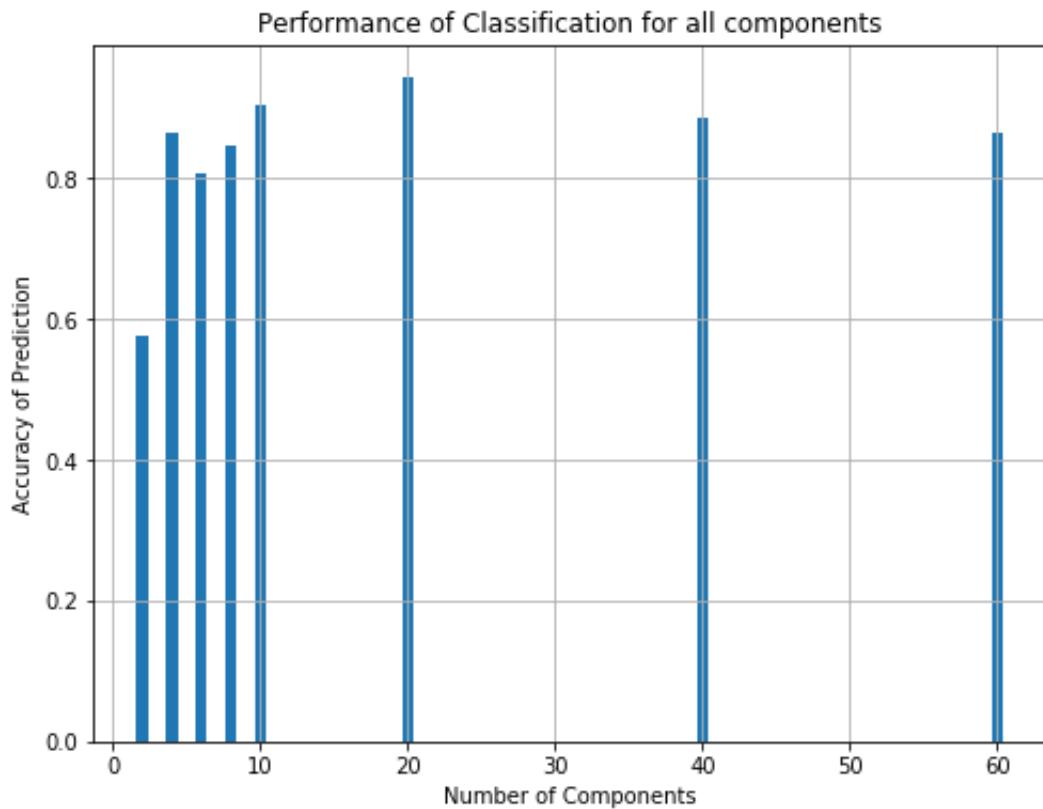
(iv) PCA with KNN

```
1 ncomps = [2, 4, 6, 8, 10, 20, 40, 60]
2 cls = [KNeighborsClassifier(n_neighbors=3,
metric='euclidean').fit(np.matmul(standard_train, v[:, :ncomps[i]]),
np.ravel(train_y)) for i in range(len(ncomps))]
```

```
1 pred4 = cls[4].predict(np.matmul(standard_test, v[:, :ncomps[4]]))
2 true4 = np.ravel(test_y)
3 with open('q1v.csv', 'w') as fp:
4     fp.write(', '.join(['Component' + str(i) for i in
range(1, ncomps[4]+1)]) + ', Class, Class\n')
5     transformed = np.matmul(standard_test, v[:, :ncomps[4]])
6     for i in range(len(transformed)):
7         fp.write(', '.join([str(x) for x in transformed[i]]))
8         fp.write(', ' + str(true4[i]))
9         fp.write(', ' + str(pred4[i]) + '\n')
```

You can find the output of the above code in file q1v.csv

```
1 fig, axes = plt.subplots(nrows=1, ncols=1)
2 fig.set_figheight(6)
3 fig.set_figwidth(8)
4 accuracies = [cls[i].score(np.matmul(standard_test, v[:, :ncomps[i]]),
np.ravel(test_y)) for i in range(len(ncomps))]
5 axes.bar(ncomps, accuracies)
6 axes.set_xlabel('Number of Components')
7 axes.set_ylabel('Accuracy of Prediction')
8 axes.set_title('Performance of Classification for all components')
9 plt.grid()
10 plt.show()
11 print('Accuracy - ')
12 print('Components\tAccuracy')
13 for i in range(len(ncomps)):
14     print('{}\t{:4f}%'.format(ncomps[i], 100*accuracies[i]))
```



```

1 | Accuracy -
2 | Components Accuracy
3 | 2      57.6923%
4 | 4      86.5385%
5 | 6      80.7692%
6 | 8      84.6154%
7 | 10     90.3846%
8 | 20     94.2308%
9 | 40     88.4615%
10| 60     86.5385%

```

Reasonable Number of PCA components - 20

The reason is simple - It gives the best accuracy out of the given number of components. If we look at the cumulative plot of eigen values above, we can see that 20 eigenvectors were able to cover almost 90% of the variance. As we keep increasing/decreasing the number of components from this point, the accuracy seems to fall gradually.

(D) Preference

Comparing both procedures, I believe that standardization is better than normalization for this dataset. Not only did PCA + KNN with standardized data give better accuracy, the eigenvalues were much smoother and more interpretable. We can clearly see how many attributes the top 10 eigen vectors were covering, giving a good intuition about the number of eigen vectors/number of PCA components to choose.

Q2 - Decision Trees

The data given is as follows -

Attr1	Attr2	Attr3	Attr4	Attr5	Class
3	BLUE	SMALL	LOW	COOL	F
7	BLUE	LARGE	HIGH	COOL	F
16	BLUE	LARGE	LOW	COOL	F
17	BLUE	LARGE	HIGH	COOL	F
27	BLUE	LARGE	HIGH	HOT	T
29	BLUE	SMALL	HIGH	HOT	T
33	BLUE	SMALL	LOW	HOT	F
34	BLUE	LARGE	HIGH	HOT	T
2	RED	LARGE	LOW	COOL	F
6	RED	SMALL	LOW	COOL	T
10	RED	SMALL	HIGH	COOL	T
11	RED	SMALL	LOW	COOL	F
25	RED	SMALL	HIGH	HOT	T
36	RED	LARGE	HIGH	HOT	T
45	RED	SMALL	LOW	HOT	F
50	RED	LARGE	LOW	HOT	F

(A) ID3

The entropy of Class label is (9F and 7T) -

$$H(\text{Class}) = -\left[\frac{9}{16}\log_2\left(\frac{9}{16}\right) + \frac{7}{16}\log_2\left(\frac{7}{16}\right)\right] = 0.9887$$

Root Split

Initial Attribute Specific Entropies (6 as split for Attr1) -

$$H(Class|Attr1) = \frac{1}{8}[-(\frac{0}{2}\log_2(\frac{0}{2}) + \frac{2}{2}\log_2(\frac{2}{2}))] + \frac{7}{8}[-(\frac{1}{2}\log_2(\frac{1}{2}) + \frac{1}{2}\log_2(\frac{1}{2}))] = 0.8750$$

$$IG(Class|Attr1) = H(Class) - H(Class|Attr1) = 0.9887 - 0.8750 = 0.1137$$

$$H(Class|Attr2) = \frac{1}{2}[-(\frac{3}{8}\log_2(\frac{3}{8}) + \frac{5}{8}\log_2(\frac{5}{8}))] + \frac{1}{2}[-(\frac{1}{2}\log_2(\frac{1}{2}) + \frac{1}{2}\log_2(\frac{1}{2}))] = 0.9772$$

$$IG(Class|Attr2) = H(Class) - H(Class|Attr1) = 0.9887 - 0.9772 = 0.0115$$

$$H(Class|Attr3) = \frac{1}{2}[-(\frac{1}{2}\log_2(\frac{1}{2}) + \frac{1}{2}\log_2(\frac{1}{2}))] + \frac{1}{2}[-(\frac{3}{8}\log_2(\frac{3}{8}) + \frac{5}{8}\log_2(\frac{5}{8}))] = 0.9772$$

$$IG(Class|Attr3) = H(Class) - H(Class|Attr1) = 0.9887 - 0.9772 = 0.0115$$

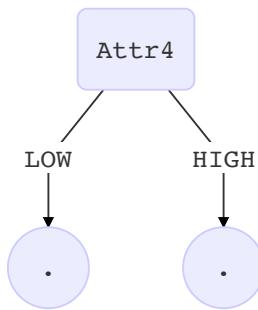
$$H(Class|Attr4) = \frac{1}{2}[-(\frac{7}{8}\log_2(\frac{7}{8}) + \frac{1}{8}\log_2(\frac{1}{8}))] + \frac{1}{2}[-(\frac{1}{4}\log_2(\frac{1}{4}) + \frac{3}{4}\log_2(\frac{3}{4}))] = 0.6774$$

$$IG(Class|Attr4) = H(Class) - H(Class|Attr1) = 0.9887 - 0.6774 = 0.3113$$

$$H(Class|Attr5) = \frac{1}{2}[-(\frac{5}{8}\log_2(\frac{5}{8}) + \frac{3}{8}\log_2(\frac{3}{8}))] + \frac{1}{2}[-(\frac{1}{4}\log_2(\frac{1}{4}) + \frac{3}{4}\log_2(\frac{3}{4}))] = 0.8828$$

$$IG(Class|Attr5) = H(Class) - H(Class|Attr1) = 0.9887 - 0.8828 = 0.1059$$

So we split at Attr4.



Level 1 Split

Attr4=high

Attr1	Attr2	Attr3	Attr4	Attr5	Class
7	BLUE	LARGE	HIGH	COOL	F
17	BLUE	LARGE	HIGH	COOL	F
27	BLUE	LARGE	HIGH	HOT	T
29	BLUE	SMALL	HIGH	HOT	T
34	BLUE	LARGE	HIGH	HOT	T
10	RED	SMALL	HIGH	COOL	T
25	RED	SMALL	HIGH	HOT	T
36	RED	LARGE	HIGH	HOT	T

$$H(Class|Attr4 = \text{high}) = -[\frac{1}{4}\log_2(\frac{1}{4}) + \frac{3}{4}\log_2(\frac{3}{4})] = 0.8113$$

Split at 25 for Attr1 -

$$H(Class|Attr4 = \text{high}, Attr1) = \frac{3}{8}[-(\frac{2}{3}\log_2(\frac{2}{3}) + \frac{1}{3}\log_2(\frac{1}{3}))] + \frac{5}{8}[-(\frac{0}{8}\log_2(\frac{0}{8}) + \frac{8}{8}\log_2(\frac{8}{8}))] = 0.3443$$

$$IG(Class|Attr4 = \text{high}, Attr1) = H(Class|Attr4 = \text{high}) - H(Class|Attr4 = \text{high}, Attr1) = 0.8113 - 0.3 = 0.4670$$

$$H(Class|Attr4 = \text{high}, Attr2) = \frac{3}{8}[-(\frac{3}{3}\log_2(\frac{3}{3}) + \frac{0}{3}\log_2(\frac{0}{3}))] + \frac{5}{8}[-(\frac{2}{5}\log_2(\frac{2}{5}) + \frac{3}{5}\log_2(\frac{3}{5}))] = 0.6068$$

$$IG(Class|Attr4 = \text{high}, Attr2) = H(Class|Attr4 = \text{high}) - H(Class|Attr4 = \text{high}, Attr2) = 0.8113 - 0.6068 = 0.2045$$

$$H(Class|Attr4 = \text{high}, Attr3) = \frac{3}{8}[-(\frac{0}{3}\log_2(\frac{0}{3}) + \frac{3}{3}\log_2(\frac{3}{3}))] + \frac{5}{8}[-(\frac{2}{5}\log_2(\frac{2}{5}) + \frac{3}{5}\log_2(\frac{3}{5}))] = 0.6068$$

$$IG(Class|Attr4 = \text{high}, Attr3) = H(Class|Attr4 = \text{high}) - H(Class|Attr4 = \text{high}, Attr3) = 0.8113 - 0.6068 = 0.2045$$

$$H(Class|Attr4 = \text{high}, Attr5) = \frac{3}{8}[-(\frac{2}{3}\log_2(\frac{2}{3}) + \frac{1}{3}\log_2(\frac{1}{3}))] + \frac{5}{8}[-(\frac{0}{8}\log_2(\frac{0}{8}) + \frac{8}{8}\log_2(\frac{8}{8}))] = 0.3443$$

$$IG(Class|Attr4 = \text{high}, Attr5) = H(Class|Attr4 = \text{high}) - H(Class|Attr4 = \text{high}, Attr5) = 0.8113 - 0.3443 = 0.4670$$

So we split on Attr1 at 25 for Attr4=high.

Attr4=low

Attr1	Attr2	Attr3	Attr4	Attr5	Class
3	BLUE	SMALL	LOW	COOL	F
16	BLUE	LARGE	LOW	COOL	F
33	BLUE	SMALL	LOW	HOT	F
2	RED	LARGE	LOW	COOL	F
6	RED	SMALL	LOW	COOL	T
11	RED	SMALL	LOW	COOL	F
45	RED	SMALL	LOW	HOT	F
50	RED	LARGE	LOW	HOT	F

$$H(Class|Attr4 = \text{low}) = -[\frac{1}{8}\log_2(\frac{1}{8}) + \frac{7}{8}\log_2(\frac{7}{8})] = 0.5436$$

Split at 11 for Attr1 -

$$H(Class|Attr4 = low, Attr1) = \frac{3}{8}[-(\frac{1}{3}\log_2(\frac{1}{3}) + \frac{2}{3}\log_2(\frac{2}{3}))] = 0.3443$$

$$IG(Class|Attr4 = low, Attr1) = H(Class|Attr4 = low) - H(Class|Attr4 = low, Attr1) = 0.5436 - 0.3443 = 0.1992$$

$$H(Class|Attr4 = low, Attr2) = \frac{3}{8}[-(\frac{3}{3}\log_2(\frac{3}{3}) + \frac{0}{3}\log_2(\frac{0}{3}))] + \frac{5}{8}[-(\frac{1}{5}\log_2(\frac{1}{5}) + \frac{4}{5}\log_2(\frac{4}{5}))] = 0.4512$$

$$IG(Class|Attr4 = low, Attr2) = H(Class|Attr4 = low) - H(Class|Attr4 = low, Attr2) = 0.5436 - 0.4512 = 0.0924$$

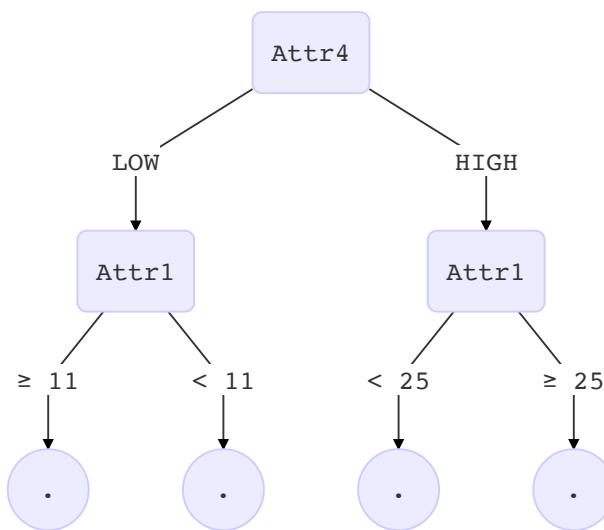
$$H(Class|Attr4 = low, Attr3) = \frac{3}{8}[-(\frac{0}{3}\log_2(\frac{0}{3}) + \frac{3}{3}\log_2(\frac{3}{3}))] + \frac{5}{8}[-(\frac{1}{5}\log_2(\frac{1}{5}) + \frac{4}{5}\log_2(\frac{4}{5}))] = 0.4512$$

$$IG(Class|Attr4 = low, Attr3) = H(Class|Attr4 = low) - H(Class|Attr4 = low, Attr3) = 0.5436 - 0.4512 = 0.0924$$

$$H(Class|Attr4 = low, Attr5) = \frac{3}{8}[-(\frac{0}{3}\log_2(\frac{0}{3}) + \frac{3}{3}\log_2(\frac{3}{3}))] + \frac{5}{8}[-(\frac{1}{5}\log_2(\frac{1}{5}) + \frac{4}{5}\log_2(\frac{4}{5}))] = 0.4512$$

$$IG(Class|Attr4 = low, Attr5) = H(Class|Attr4 = low) - H(Class|Attr4 = low, Attr5) = 0.5436 - 0.4512 = 0.0924$$

So we split at Attr1 at value 11



Level 2 Split

Attr4=high, Attr1 ≥ 25

Attr1	Attr2	Attr3	Attr4	Attr5	Class
27	BLUE	LARGE	HIGH	HOT	T
29	BLUE	SMALL	HIGH	HOT	T
34	BLUE	LARGE	HIGH	HOT	T
25	RED	SMALL	HIGH	HOT	T
36	RED	LARGE	HIGH	HOT	T

As all labels are same(T) this is a leaf node. So we don't proceed further on this branch.

Attr4=high, Attr1<25

Attr1	Attr2	Attr3	Attr4	Attr5	Class
7	BLUE	LARGE	HIGH	COOL	F
17	BLUE	LARGE	HIGH	COOL	F
10	RED	SMALL	HIGH	COOL	T

$$H(\text{Class} | \text{Attr4} = \text{high}, \text{Attr1} < 25) = -[\frac{1}{3}\log_2(\frac{1}{3}) + \frac{2}{3}\log_2(\frac{2}{3})] = 0.9183$$

$$H(\text{Class} | \text{Attr4} = \text{high}, \text{Attr1} < 25, \text{Attr2}) = \frac{1}{3}[-\frac{1}{1}\log_2(\frac{1}{1})] + \frac{2}{3}[-\frac{2}{2}\log_2(\frac{2}{2})] = 0$$

$$IG(\text{Class} | \text{Attr4} = \text{high}, \text{Attr1} < 25, \text{Attr2}) = 0.9183$$

$$H(\text{Class} | \text{Attr4} = \text{high}, \text{Attr1} < 25, \text{Attr3}) = \frac{1}{3}[-\frac{1}{1}\log_2(\frac{1}{1})] + \frac{2}{3}[-\frac{2}{2}\log_2(\frac{2}{2})] = 0$$

$$IG(\text{Class} | \text{Attr4} = \text{high}, \text{Attr1} < 25, \text{Attr3}) = 0.9183$$

$$H(\text{Class} | \text{Attr4} = \text{high}, \text{Attr1} < 25, \text{Attr5}) = \frac{1}{1}[-(\frac{2}{3}\log_2(\frac{2}{3}) + \frac{1}{3}\log_2(\frac{1}{3}))] = 0.9183$$

$$IG(\text{Class} | \text{Attr4} = \text{high}, \text{Attr1} < 25, \text{Attr5}) = 0$$

We split on Attr2.

Attr4=low, Attr1≥11

Attr1	Attr2	Attr3	Attr4	Attr5	Class
16	BLUE	LARGE	LOW	COOL	F
33	BLUE	SMALL	LOW	HOT	F
11	RED	SMALL	LOW	COOL	F
45	RED	SMALL	LOW	HOT	F
50	RED	LARGE	LOW	HOT	F

As all labels are same(F), this is a leaf node. So we don't proceed further along this branch.

Attr4=low, Attr1<11

Attr1	Attr2	Attr3	Attr4	Attr5	Class
3	BLUE	SMALL	LOW	COOL	F
2	RED	LARGE	LOW	COOL	F
6	RED	SMALL	LOW	COOL	T

$$H(Class|Attr4 = low, Attr1 < 11) = -[\frac{2}{3}\log_2(\frac{2}{3}) + \frac{1}{3}\log_2(\frac{1}{3})] = 0.9183$$

$$H(Class|Attr4 = high, Attr1 < 11, Attr2) = \frac{1}{3}[-\frac{1}{1}\log_2(\frac{1}{1})] + \frac{2}{3}[-(\frac{1}{2}\log_2(\frac{1}{2}) + \frac{1}{2}\log_2(\frac{1}{2}))] = 0.6667$$

$$IG(Class|Attr4 = high, Attr1 < 11, Attr2) = 0.2516$$

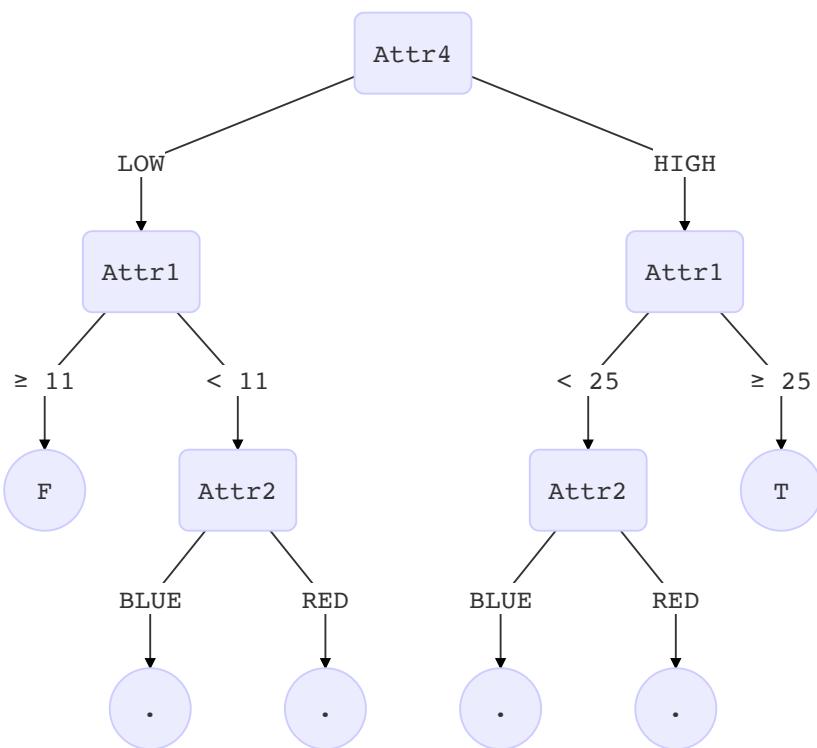
$$H(Class|Attr4 = high, Attr1 < 11, Attr3) = \frac{1}{3}[-\frac{1}{1}\log_2(\frac{1}{1})] + \frac{2}{3}[-(\frac{1}{2}\log_2(\frac{1}{2}) + \frac{1}{2}\log_2(\frac{1}{2}))] = 0.6667$$

$$IG(Class|Attr4 = high, Attr1 < 11, Attr3) = 0.2516$$

$$H(Class|Attr4 = high, Attr1 < 11, Attr5) = \frac{1}{1}[-(\frac{2}{3}\log_2(\frac{2}{3}) + \frac{1}{3}\log_2(\frac{1}{3}))] = 0.9183$$

$$IG(Class|Attr4 = high, Attr1 < 11, Attr5) = 0$$

We split on Attr2.



Level 3 Split

Attr4=high, Attr1<25, Attr2=red

Attr1	Attr2	Attr3	Attr4	Attr5	Class
10	RED	SMALL	HIGH	COOL	T

As all labels are same(T), this is a leaf node. So we don't proceed further along this branch.

Attr4=high, Attr1<25, Attr2=blue

Attr1	Attr2	Attr3	Attr4	Attr5	Class
7	BLUE	LARGE	HIGH	COOL	F
17	BLUE	LARGE	HIGH	COOL	F

As all labels are same(F), this is a leaf node. So we don't proceed further along this branch.

Attr4=low, Attr1<11, Attr2=blue

Attr1	Attr2	Attr3	Attr4	Attr5	Class
3	BLUE	SMALL	LOW	COOL	F

As all labels are same(F), this is a leaf node. So we don't proceed further along this branch.

Attr4=low, Attr1<11, Attr2=red

Attr1	Attr2	Attr3	Attr4	Attr5	Class
2	RED	LARGE	LOW	COOL	F
6	RED	SMALL	LOW	COOL	T

$$H(\text{Class} | \text{Attr4} = \text{low}, \text{Attr2} = \text{red}, \text{Attr1} < 11) = -[\frac{1}{2}\log_2(\frac{1}{2}) + \frac{1}{2}\log_2(\frac{1}{2})] = 1.0000$$

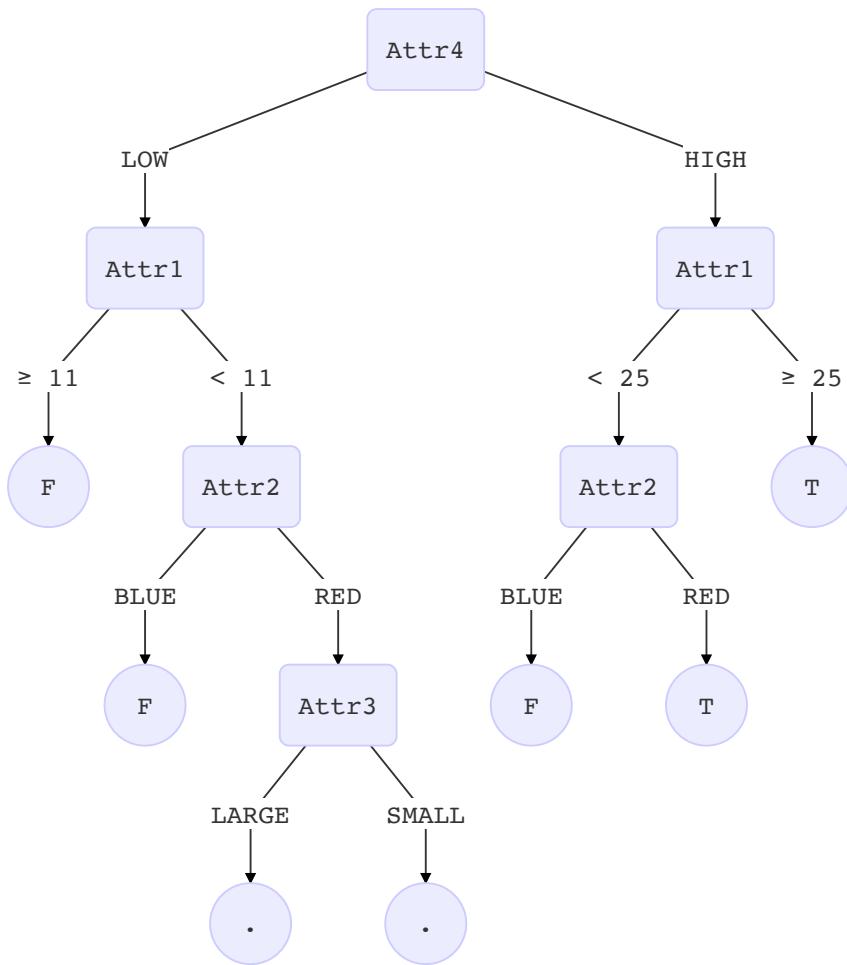
$$H(\text{Class} | \text{Attr4} = \text{low}, \text{Attr2} = \text{red}, \text{Attr1} < 11, \text{Attr3}) = \frac{1}{2}[-(\frac{1}{1}\log_2(\frac{1}{1})) + \frac{1}{2}[-(\frac{1}{1}\log_2(\frac{1}{1}))]] = 0.0000$$

$$IG(\text{Class} | \text{Attr4} = \text{low}, \text{Attr2} = \text{red}, \text{Attr1} < 11, \text{Attr3}) = 1.000$$

$$H(\text{Class} | \text{Attr4} = \text{low}, \text{Attr2} = \text{red}, \text{Attr1} < 11, \text{Attr5}) = \frac{2}{2}[-(\frac{1}{2}\log_2(\frac{1}{2}) + \frac{1}{2}\log_2(\frac{1}{2}))] = 1.0000$$

$$IG(\text{Class} | \text{Attr4} = \text{low}, \text{Attr2} = \text{red}, \text{Attr1} < 11, \text{Attr5}) = 0.0000$$

So we split on Attr3.



Level 4 Split

Attr4=low, Attr2=red, Attr1<11, Attr3=large

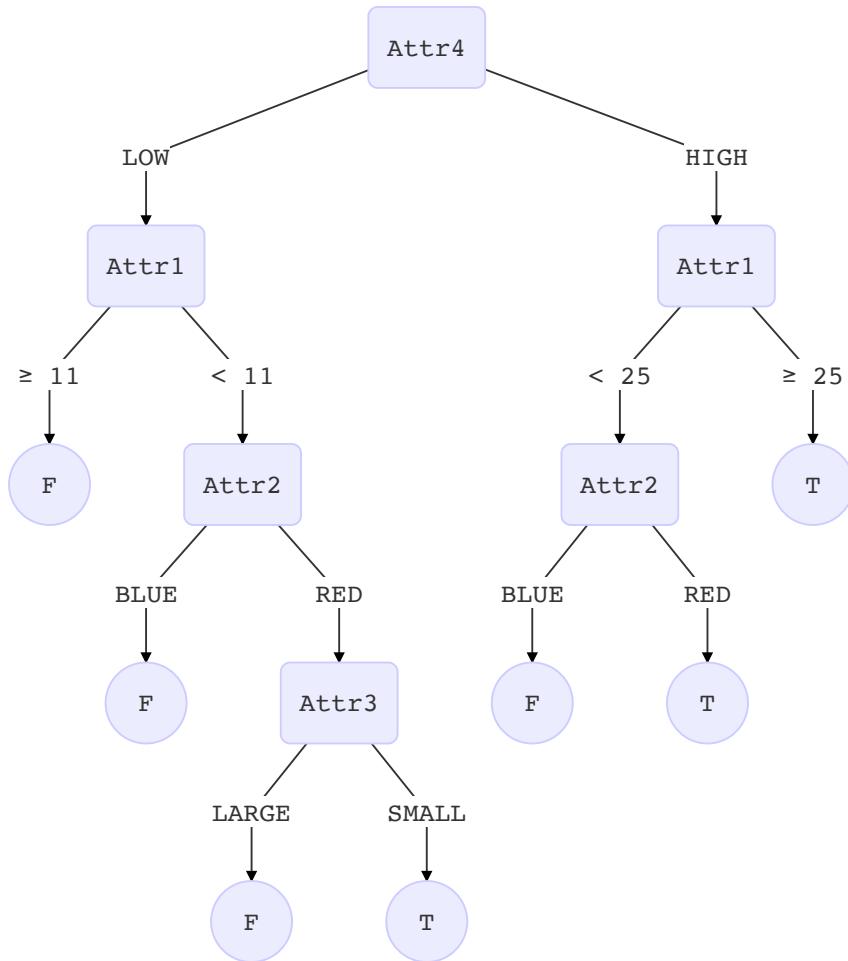
Attr1	Attr2	Attr3	Attr4	Attr5	Class
2	RED	LARGE	LOW	COOL	F

As all labels are same(F), this is a leaf node. So we don't proceed further along this branch.

Attr4=low, Attr2=red, Attr1<11, Attr3=small

Attr1	Attr2	Attr3	Attr4	Attr5	Class
6	RED	SMALL	LOW	COOL	T

As all labels are same(T), this is a leaf node. So we don't proceed further along this branch.



(B) GINI

Gini Score of Class attribute (9F and 7T) -

$$Gini = 1 - \left(\frac{9}{16}\right)^2 - \left(\frac{7}{16}\right)^2 = 0.4922$$

Root Split

Initial Attribute Specific Gini Splits (6 as split for Attr1) -

$$Gini_{split}(Attr1) = \frac{2}{16}[1 - (\frac{0}{2})^2 - (\frac{2}{2})^2] + \frac{14}{16}[1 - (\frac{7}{14})^2 + (\frac{7}{14})^2] = 0.4375$$

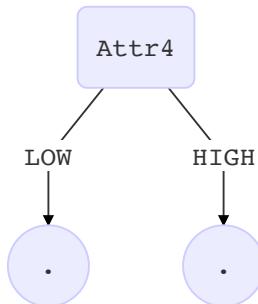
$$Gini_{split}(Attr2) = \frac{1}{2}[1 - (\frac{3}{8})^2 - (\frac{5}{8})^2] + \frac{1}{2}[1 - (\frac{4}{8})^2 + (\frac{4}{8})^2] = 0.4843$$

$$Gini_{split}(Attr3) = \frac{1}{2}[1 - (\frac{3}{8})^2 - (\frac{5}{8})^2] + \frac{1}{2}[1 - (\frac{4}{8})^2 + (\frac{4}{8})^2] = 0.4843$$

$$Gini_{split}(Attr4) = \frac{1}{2}[1 - (\frac{1}{8})^2 - (\frac{7}{8})^2] + \frac{1}{2}[1 - (\frac{1}{4})^2 + (\frac{3}{4})^2] = 0.2969$$

$$Gini_{split}(Attr5) = \frac{1}{2}[1 - (\frac{3}{8})^2 - (\frac{5}{8})^2] + \frac{1}{2}[1 - (\frac{1}{4})^2 + (\frac{3}{4})^2] = 0.4218$$

The lowest GINI split score is for Attr4. So we split on Attr4.



Level 1 Split

Attr4=low

Attr1	Attr2	Attr3	Attr4	Attr5	Class
3	BLUE	SMALL	LOW	COOL	F
16	BLUE	LARGE	LOW	COOL	F
33	BLUE	SMALL	LOW	HOT	F
2	RED	LARGE	LOW	COOL	F
6	RED	SMALL	LOW	COOL	T
11	RED	SMALL	LOW	COOL	F
45	RED	SMALL	LOW	HOT	F
50	RED	LARGE	LOW	HOT	F

$$Gini(Attr4 = low) = 1 - \left(\frac{1}{8}\right)^2 - \left(\frac{7}{8}\right)^2 = 0.2187$$

Splitting at 11 for Attr1

$$Gini_{split}(Attr1|Attr4 = low) = \frac{3}{8}[1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2] + \frac{5}{8}[1 - \left(\frac{5}{5}\right)^2] = 0.1667$$

$$Gini_{split}(Attr2|Attr4 = low) = \frac{3}{8}[1 - \left(\frac{3}{3}\right)^2] + \frac{5}{8}[1 - \left(\frac{4}{5}\right)^2 + \left(\frac{1}{5}\right)^2] = 0.2000$$

$$Gini_{split}(Attr3|Attr4 = low) = \frac{3}{8}[1 - \left(\frac{3}{3}\right)^2] + \frac{5}{8}[1 - \left(\frac{4}{5}\right)^2 + \left(\frac{1}{5}\right)^2] = 0.2000$$

$$Gini_{split}(Attr5|Attr4 = low) = \frac{3}{8}[1 - \left(\frac{3}{3}\right)^2] + \frac{5}{8}[1 - \left(\frac{4}{5}\right)^2 + \left(\frac{1}{5}\right)^2] = 0.2000$$

So we split at Attr1 at value 11

Attr4=high

Attr1	Attr2	Attr3	Attr4	Attr5	Class
7	BLUE	LARGE	HIGH	COOL	F
17	BLUE	LARGE	HIGH	COOL	F
27	BLUE	LARGE	HIGH	HOT	T
29	BLUE	SMALL	HIGH	HOT	T
34	BLUE	LARGE	HIGH	HOT	T
10	RED	SMALL	HIGH	COOL	T
25	RED	SMALL	HIGH	HOT	T
36	RED	LARGE	HIGH	HOT	T

$$Gini(Attr4 = high) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = 0.375$$

Splitting at 25 for Attr1

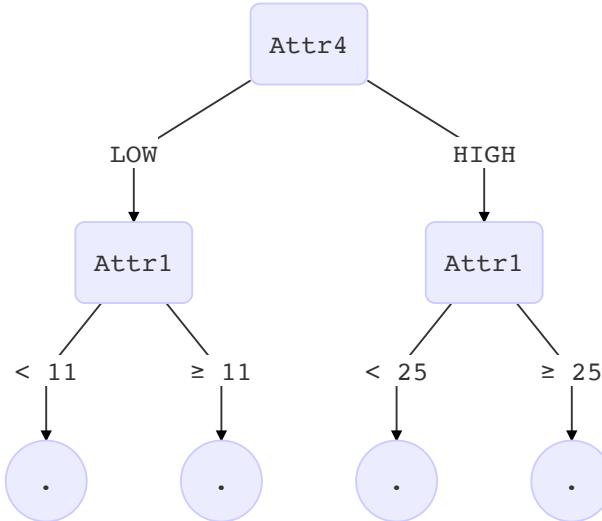
$$Gini_{split}(Attr1|Attr4 = high) = \frac{3}{8}[1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2] + \frac{5}{8}[1 - \left(\frac{5}{5}\right)^2] = 0.1667$$

$$Gini_{split}(Attr2|Attr4 = high) = \frac{3}{8}[1 - (\frac{3}{3})^2] + \frac{5}{8}[1 - (\frac{2}{5})^2 + (\frac{3}{5})^2] = 0.3000$$

$$Gini_{split}(Attr3|Attr4 = high) = \frac{3}{8}[1 - (\frac{3}{3})^2] + \frac{5}{8}[1 - (\frac{2}{5})^2 + (\frac{3}{5})^2] = 0.3000$$

$$Gini_{split}(Attr5|Attr4 = high) = \frac{3}{8}[1 - (\frac{1}{3})^2 - (\frac{2}{3})^2] + \frac{5}{8}[1 - (\frac{5}{5})^2] = 1.667$$

So we split at Attr1 at value 25



Level 2 Split

Attr4=low, Attr1<11

Attr1	Attr2	Attr3	Attr4	Attr5	Class
3	BLUE	SMALL	LOW	COOL	F
2	RED	LARGE	LOW	COOL	F
6	RED	SMALL	LOW	COOL	T

$$Gini(Attr4 = low, Attr1 < 11) = 1 - (\frac{1}{3})^2 - (\frac{2}{3})^2 = 0.4444$$

$$Gini_{split}(Attr2|Attr4 = low, Attr1 < 11) = \frac{1}{3}[1 - (\frac{1}{1})^2] + \frac{2}{3}[1 - (\frac{1}{2})^2 + (\frac{1}{2})^2] = 0.3333$$

$$Gini_{split}(Attr3|Attr4 = low, Attr1 < 11) = \frac{1}{3}[1 - (\frac{1}{1})^2] + \frac{2}{3}[1 - (\frac{1}{2})^2 + (\frac{1}{2})^2] = 0.3333$$

$$Gini_{split}(Attr5|Attr4 = low, Attr1 < 11) = \frac{1}{1}[1 - (\frac{1}{3})^2 + (\frac{2}{3})^2] = 0.4444$$

So we split on Attr2

Attr4=low, Attr1≥11

Attr1	Attr2	Attr3	Attr4	Attr5	Class
16	BLUE	LARGE	LOW	COOL	F
33	BLUE	SMALL	LOW	HOT	F
11	RED	SMALL	LOW	COOL	F
45	RED	SMALL	LOW	HOT	F
50	RED	LARGE	LOW	HOT	F

As all labels are same(F), this is a leaf node. So we don't proceed further along this branch.

Attr4=high, Attr1<25

Attr1	Attr2	Attr3	Attr4	Attr5	Class
7	BLUE	LARGE	HIGH	COOL	F
17	BLUE	LARGE	HIGH	COOL	F
10	RED	SMALL	HIGH	COOL	T

$$Gini(Attr4 = \text{high}, Attr1 < 25) = 1 - (\frac{1}{3})^2 - (\frac{2}{3})^2 = 0.4444$$

$$Gini_{\text{split}}(\text{Attr2} | Attr4 = \text{high}, Attr1 < 25) = \frac{1}{3}[1 - (\frac{1}{1})^2] + \frac{2}{3}[1 - (\frac{1}{1})^2] = 0.0000$$

$$Gini_{\text{split}}(\text{Attr3} | Attr4 = \text{high}, Attr1 < 25) = \frac{1}{3}[1 - (\frac{1}{1})^2] + \frac{2}{3}[1 - (\frac{1}{1})^2] = 0.0000$$

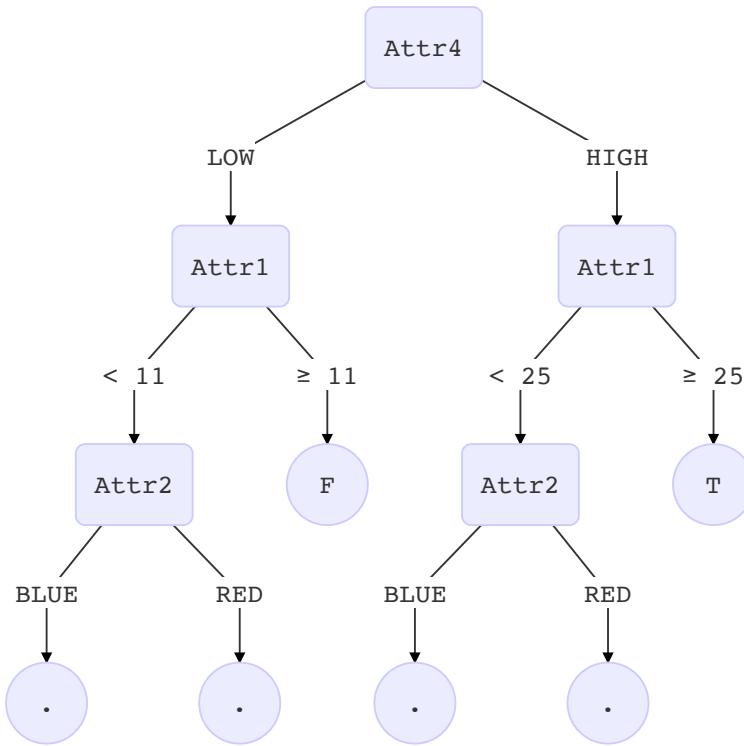
$$Gini_{\text{split}}(\text{Attr5} | Attr4 = \text{high}, Attr1 < 25) = \frac{1}{1}[1 - (\frac{1}{3})^2 + (\frac{2}{3})^2] = 0.4444$$

So we split on Attr2

Attr4=high, Attr1≥25

Attr1	Attr2	Attr3	Attr4	Attr5	Class
27	BLUE	LARGE	HIGH	HOT	T
29	BLUE	SMALL	HIGH	HOT	T
34	BLUE	LARGE	HIGH	HOT	T
25	RED	SMALL	HIGH	HOT	T
36	RED	LARGE	HIGH	HOT	T

As all labels are same(T), this is a leaf node. So we don't proceed further along this branch.



Level 3 Split

Attr4=low, Attr1<11, Attr2=blue

Attr1	Attr2	Attr3	Attr4	Attr5	Class
3	BLUE	SMALL	LOW	COOL	F

As all labels are same(F), this is a leaf node. So we don't proceed further along this branch.

Attr4=low, Attr1<11, Attr2=red

Attr1	Attr2	Attr3	Attr4	Attr5	Class
2	RED	LARGE	LOW	COOL	F
6	RED	SMALL	LOW	COOL	T

$$Gini(Attr4 = \text{low}, Attr1 < 11, Attr2 = \text{red}) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.4444$$

$$Gini_{split}(Attr3|Attr4 = low, Attr1 < 11, Attr2 = red) = \frac{1}{2}[1 - (\frac{1}{1})^2] + \frac{1}{2}[1 - (\frac{1}{1})^2] = 0.0000$$

$$Gini_{split}(Attr5|Attr4 = low, Attr1 < 11, Attr2 = red) = \frac{1}{1}[1 - (\frac{1}{2})^2 + (\frac{1}{2})^2] = 0.5000$$

So we split on Attr3

Attr4=high, Attr1<25, Attr2=blue

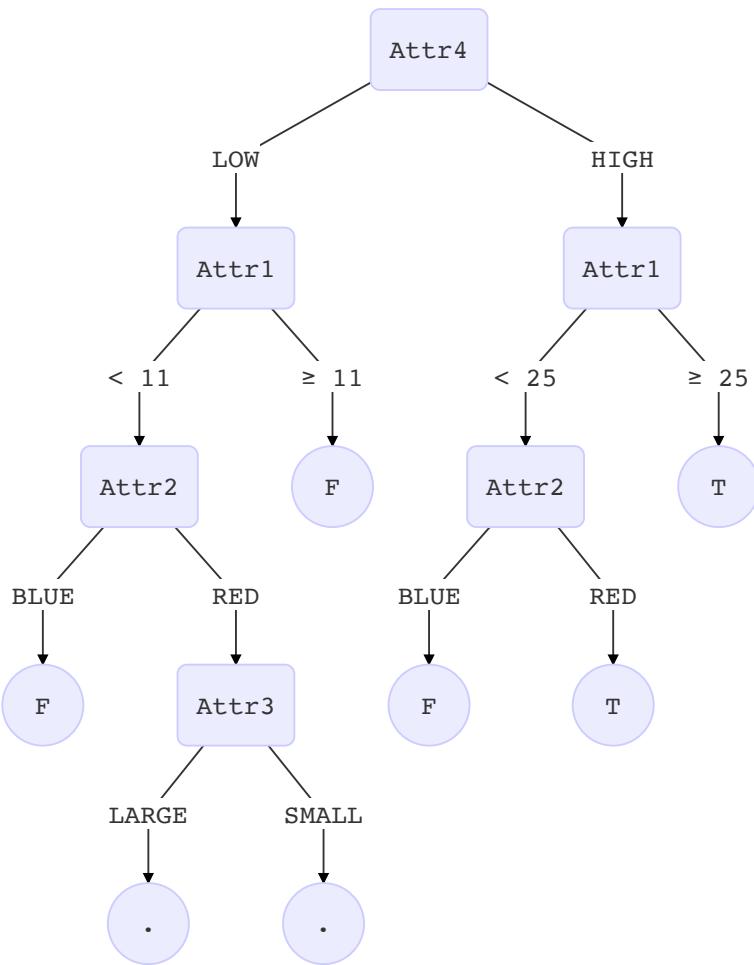
Attr1	Attr2	Attr3	Attr4	Attr5	Class
7	BLUE	LARGE	HIGH	COOL	F
17	BLUE	LARGE	HIGH	COOL	F

As all labels are same(F), this is a leaf node. So we don't proceed further along this branch.

Attr4=high, Attr1<25, Attr2=red

Attr1	Attr2	Attr3	Attr4	Attr5	Class
10	RED	SMALL	HIGH	COOL	T

As all labels are same(T), this is a leaf node. So we don't proceed further along this branch.



Split Level 4

Attr4=low, Attr1<11, Attr2=red, Attr3=large

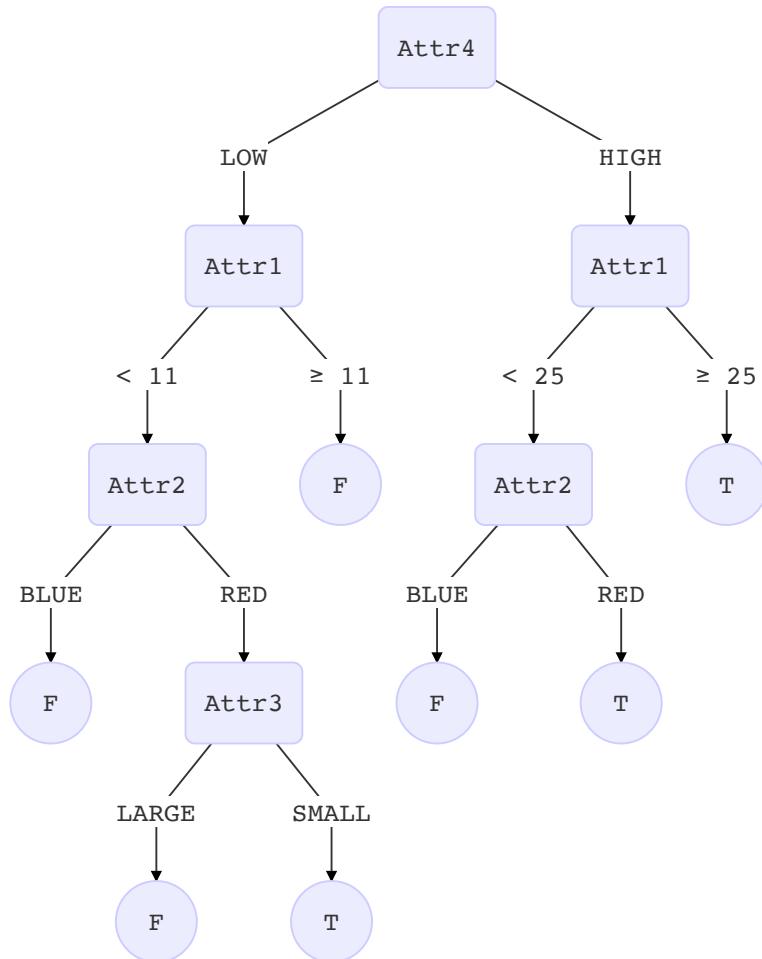
Attr1	Attr2	Attr3	Attr4	Attr5	Class
2	RED	LARGE	LOW	COOL	F

As all labels are same(F), this is a leaf node. So we don't proceed further along this branch.

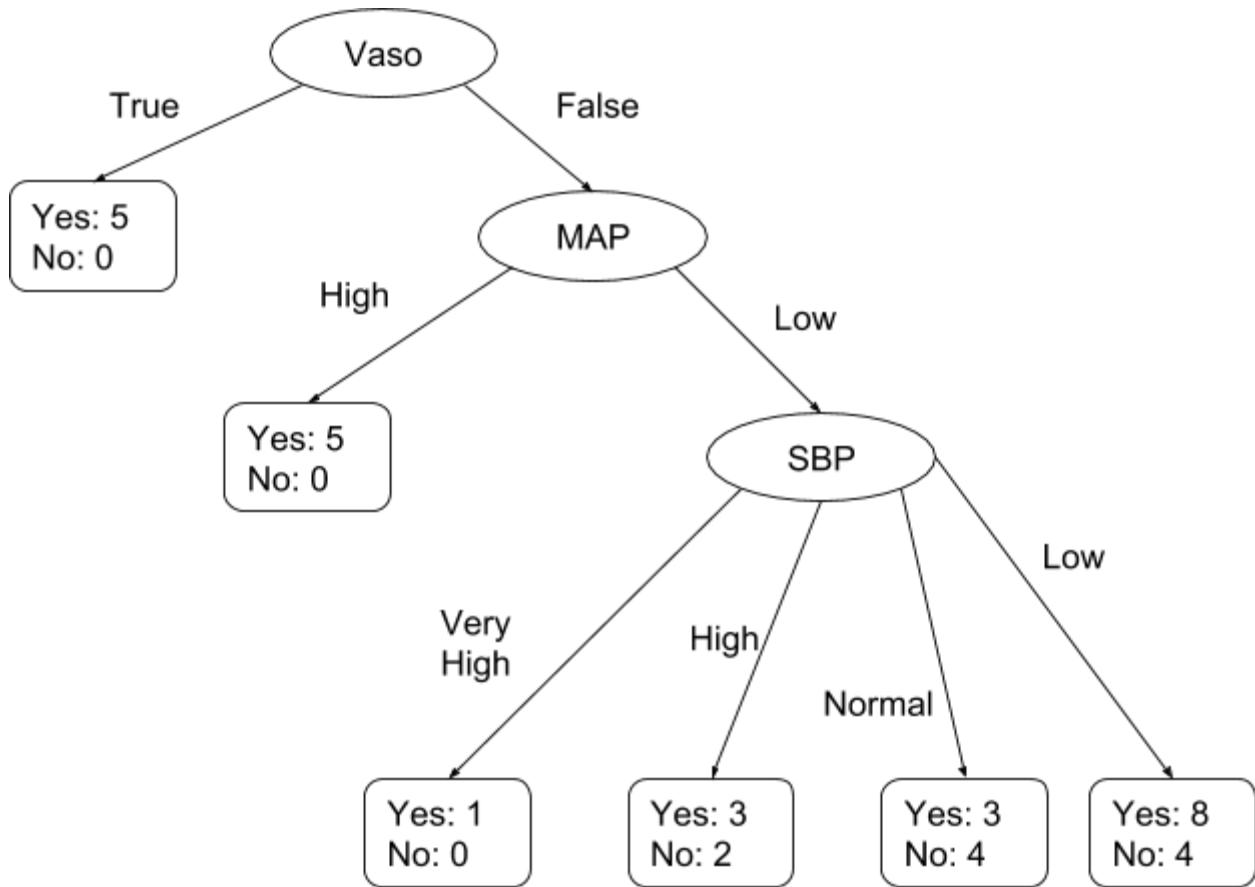
Attr4=low, Attr1<11, Attr2=red, Attr3=small

Attr1	Attr2	Attr3	Attr4	Attr5	Class
6	RED	SMALL	LOW	COOL	T

As all labels are same(T), this is a leaf node. So we don't proceed further along this branch.



Problem 3:



(a) Post-pruning based on optimistic errors

(i) Optimistic Errors:

Before Splitting	After Splitting
$\frac{10}{25}$	$\frac{9}{25}$

- (ii) Based upon the optimistic errors, the subtree should not be pruned. As after splitting based on SBP, the optimistic error reduces. Therefore, we keep the given decision tree retained.
- (iii) Classified provided testing dataset ("hw2q3_test.csv"):

Vaso	MAP	SBP	Sepsis Shock	Predicted Class (Based on given retained decision tree)
FALSE	Low	High	Yes	Yes
FALSE	Low	High	Yes	Yes
FALSE	Low	High	Yes	Yes
FALSE	Low	High	No	Yes
TRUE	High	High	Yes	Yes
FALSE	Low	Low	Yes	Yes
FALSE	Low	Low	Yes	Yes
FALSE	Low	Low	Yes	Yes
TRUE	Low	Low	Yes	Yes
FALSE	High	Low	Yes	Yes
FALSE	Low	Normal	Yes	No
FALSE	Low	Normal	Yes	No
FALSE	Low	Normal	No	No
FALSE	Low	Normal	No	No
FALSE	Low	Normal	No	No
TRUE	High	Normal	Yes	Yes
FALSE	Low	Very High	Yes	Yes
TRUE	Low	Very High	Yes	Yes

Performance Matrix

		Predicted Class	
		Class = Yes	Class = No
Actual Class	Class = Yes	13 (TP) (a)	2 (FN) (b)
	Class = No	1 (FP) (c)	4 (TN) (d)

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{13 + 4}{13 + 2 + 1 + 4} = \frac{17}{20} = 0.85$$

$$\text{Recall} = \frac{a}{a + b} = \frac{13}{13 + 2} = \frac{13}{15} = 0.867$$

$$\text{Precision} = \frac{a}{a + c} = \frac{13}{13 + 1} = \frac{13}{14} = 0.93$$

$$\text{Specificity (TNR)} = \frac{TN}{N} = \frac{TN}{TN + FP} = \frac{4}{4 + 1} = \frac{4}{5} = 0.8$$

$$\text{Sensitivity (TPR)} = \frac{TP}{P} = \frac{TP}{TP + FN} = \frac{13}{13 + 2} = \frac{13}{15} = 0.867$$

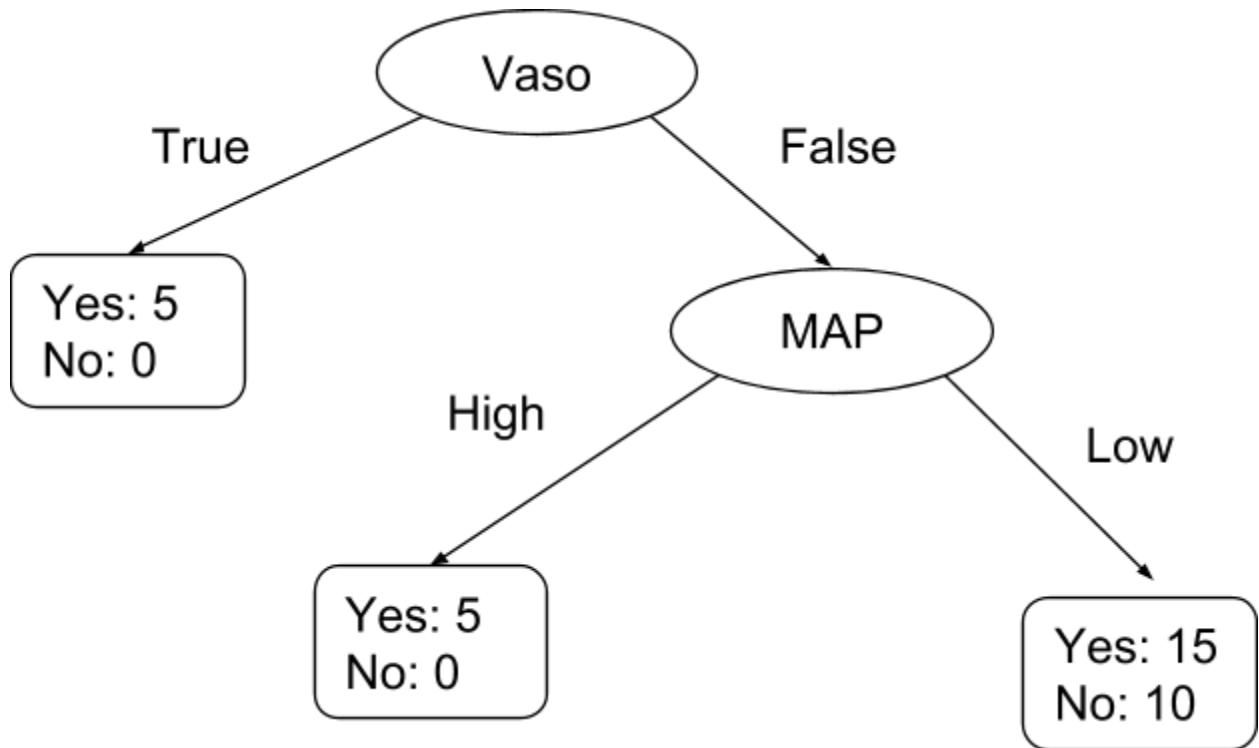
$$\text{F1 Measure} = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} = 2 * \frac{0.93 * 0.867}{0.93 + 0.867} = 0.8965$$

(b) Post-pruning based on pessimistic errors

(i) Pessimistic Errors:

Before Splitting	After Splitting
$\frac{(10 + 1 * 0.5)}{25} = \frac{10.5}{25}$	$\frac{(9 + 4 * 0.5)}{25} = \frac{11}{25}$

(ii) Based upon the pessimistic errors, the subtree should be pruned. As after splitting based on SBP, the pessimistic error increases. Therefore, we prune the given decision tree from SBP. Pruned Decision Tree is as shown below:



(iii) Classified provided testing dataset ("hw2q3_test.csv"):

Vaso	MAP	SBP	Sepsis Shock	Predicted Class (Has Septic Shock)
FALSE	Low	High	Yes	Yes
FALSE	Low	High	Yes	Yes
FALSE	Low	High	Yes	Yes
FALSE	Low	High	No	Yes
TRUE	High	High	Yes	Yes
FALSE	Low	Low	Yes	Yes
FALSE	Low	Low	Yes	Yes
FALSE	Low	Low	Yes	Yes
TRUE	Low	Low	Yes	Yes

FALSE	High	Low	Yes	Yes
FALSE	Low	Normal	Yes	Yes
FALSE	Low	Normal	Yes	Yes
FALSE	Low	Normal	No	Yes
FALSE	Low	Normal	No	Yes
FALSE	Low	Normal	No	Yes
FALSE	Low	Normal	No	Yes
TRUE	High	Normal	Yes	Yes
FALSE	Low	Very High	Yes	Yes
TRUE	Low	Very High	Yes	Yes

Performance Matrix

		Predicted Class	
		Class = Yes	Class = No
Actual Class	Class = Yes	15 (TP) (a)	0 (FN) (b)
	Class = No	5 (FP) (c)	0 (TN) (d)

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{15 + 0}{15 + 0 + 5 + 0} = \frac{15}{20} = 0.75$$

$$\text{Recall} = \frac{a}{a + b} = \frac{15}{15 + 0} = \frac{15}{15} = 1$$

$$\text{Precision} = \frac{a}{a + c} = \frac{15}{20} = \frac{15}{20} = 0.75$$

$$\text{Specificity (TNR)} = \frac{TN}{N} = \frac{TN}{TN + FP} = \frac{0}{0 + 5} = \frac{0}{5} = 0$$

$$\text{Sensitivity (TPR)} = \frac{TP}{P} = \frac{TP}{TP + FN} = \frac{15}{15 + 0} = \frac{15}{15} = 1$$

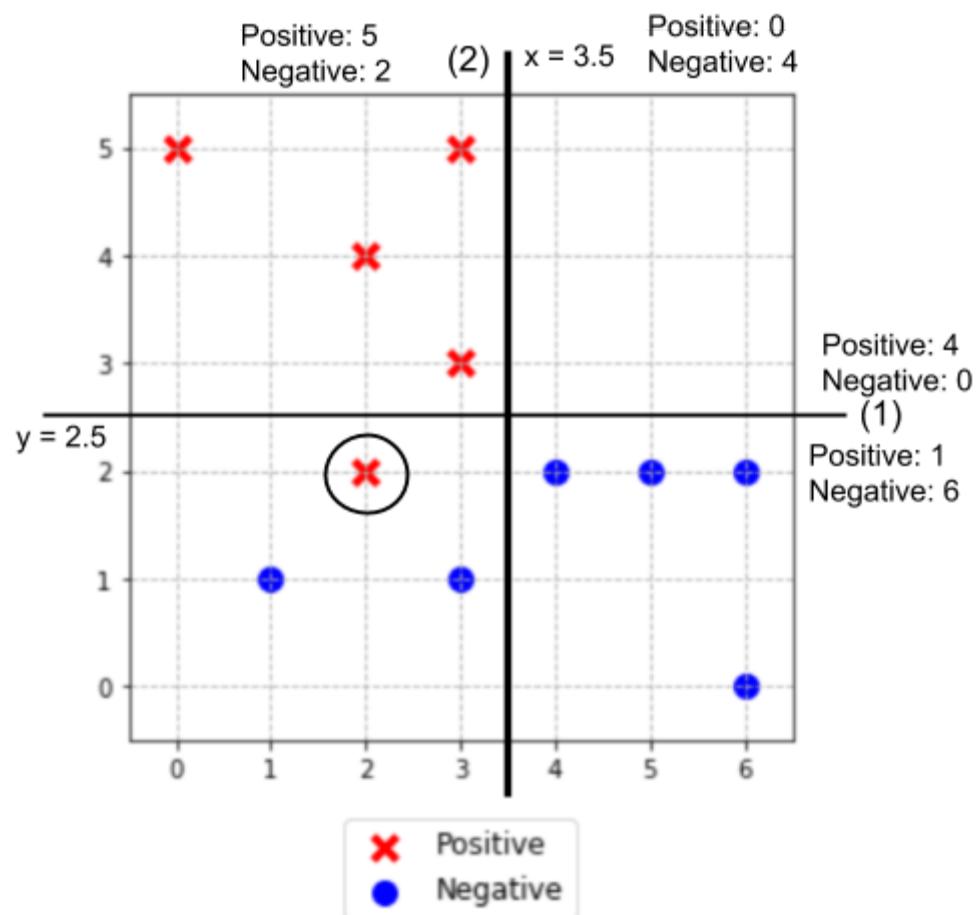
$$\text{F1 Measure} = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} = 2 * \frac{0.75 * 1}{0.75 + 1} = 0.857$$

(c) Performance Comparison

	Before Pruning	After Pruning
Accuracy	0.85	0.75
Recall	0.867	1
Precision	0.93	0.75

Here, as the problem mentioned that the sepsis deaths could be prevented with early diagnosis and intervention. So, it is very important to measure TP and TNs are very costly for this example. So, we more care about recall. So, after pruning we have higher recall. So problem 3 (b) (ii) decision tree would be better model for the task of septic shock prediction.

Problem 4:



Q4 (C)

Import Libraries

```
1 | import math
```

Iteration 1

```
1 | init_weight = 1 / 11
2 | print(f'Initial Weight: {init_weight}')
```

```
1 | Initial Weight: 0.09090909090909091
```

```
1 | epsilon = 1 * init_weight
2 | print(f'Epsilon: {epsilon}')
```

```
1 | Epsilon: 0.09090909090909091
```

```
1 | alpha = 0.5 * math.log((1-epsilon)/epsilon)
2 | print(f'Alpha: {alpha}')
```

```
1 | Alpha: 1.151292546497023
```

Correctly Classified Points

```
1 | exp_correct = init_weight * math.exp(-1*alpha)
2 | print(f'Correctly Classified : {exp_correct}')
```

```
1 | Correctly Classified : 0.028747978728803445
```

```
1 | exp_incorrect = init_weight * math.exp(alpha)
2 | print(f'Incorrectly Classified : {exp_incorrect}')
```

```
1 | Incorrectly Classified : 0.2874797872880345
```

Normalization

```
1 | z = 1 * exp_incorrect + 10 * exp_correct
2 | print(f'Normalization Factor: {z}')
```

```
1 | Normalization Factor: 0.5749595745760689
```

```
1 | w2_correct = exp_correct / z
2 | print(f'W2_Correct: {w2_correct}')
```

```
1 | W2_Correct: 0.04999999999999996
```

```
1 | w2_incorrect = exp_incorrect / z
2 | print(f'W2_InCorrect: {w2_incorrect}')
```

```
1 | W2_InCorrect: 0.5000000000000001
```

Weighted Error

```
1 | weighted_error = w2_incorrect * 1
2 | print(f'Weighted Error: {weighted_error}')
```

```
1 | Weighted Error: 0.5000000000000001
```

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,
      'hypermetrope':1, 'no':0, 'yes':1, 'normal':1, 'reduced':0}
4 d2={'Yes':1, 'No':0}
5 X = [[d1[X[i][j]] for j in range(len(X[0]))] for i in range(len(X))]
6 X = np.asarray(X)
7 Y = [d2[Y[i]] for i in range(len(Y))]
8 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],
19             dt[i], np.exp(nb.feature_log_prob_)[i][j]))
20     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
10    np.asarray(indexes)[folds[i]['test']])))
11    print('\tActual -\t' + ', '.join(np.ravel(Y[folds[i]['test']])))
12    print('\tPredict -\t' + ', '.join(ypred))
13 #     print('\tActual -\t' + ', '.join(['Yes' if x==1 else 'No' for x in
14 Y[folds[i]['test']])))
15 #     print('\tPredict -\t' + ', '.join(['Yes' if x==1 else 'No' for x in
16 ypred]))
17     for j in range(len(ypred)):
18       if ypred[j]!=Y[folds[i]['test']][j]:
19         cnt+=1
20
21     dot_data = export_graphviz(dt.tree_
22 'fold'+str(i)+'.dot',feature_names = subheaders)
23 print('\nDecision Tree 5-fold CV accuracy - ' + str((24-cnt)*100/24) + '%')

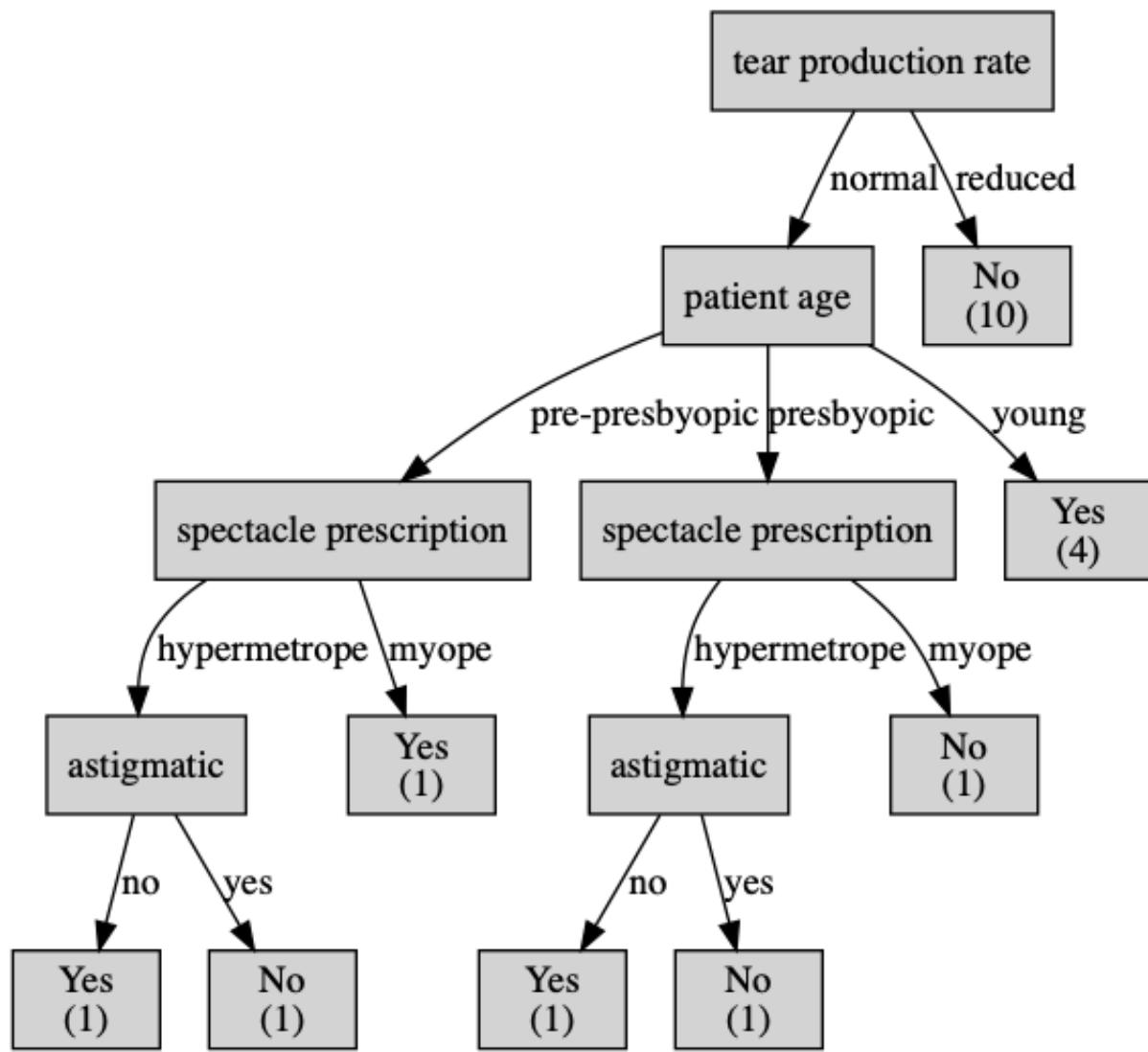
```

```

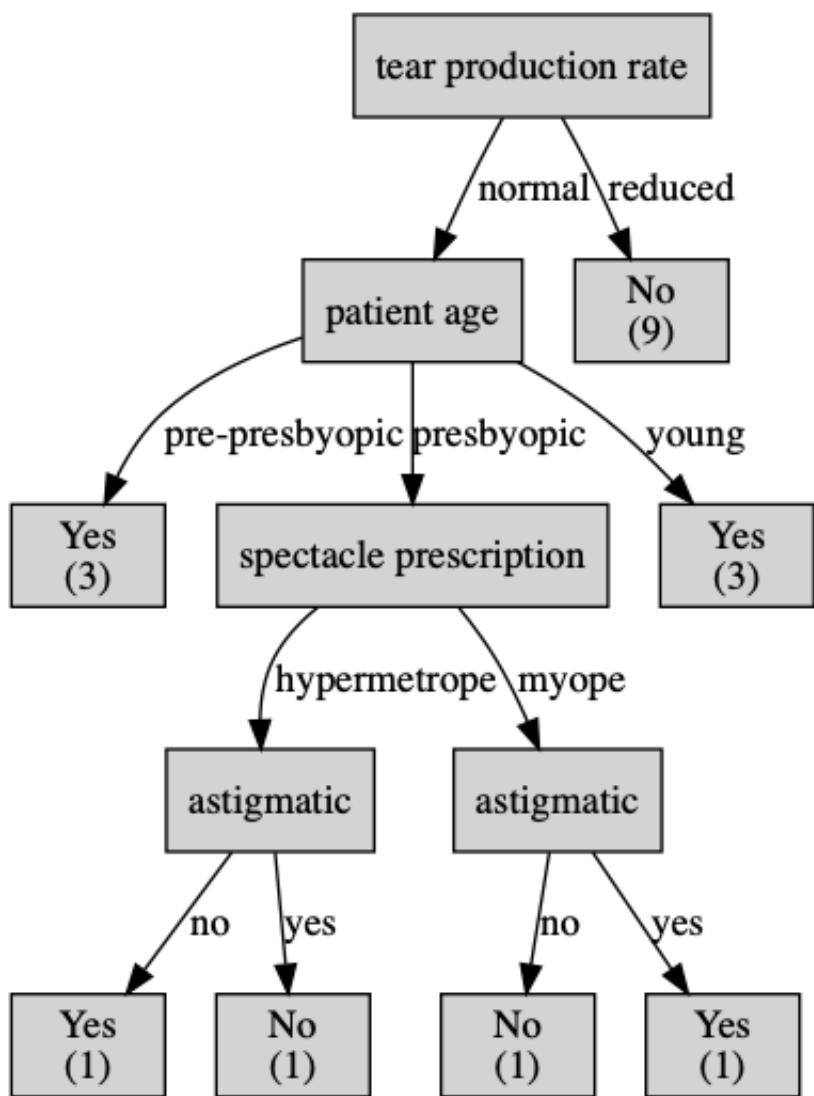
1 Fold 1
2   Test IID -  5, 10, 15, 20
3   Actual -    No, Yes, No, Yes
4   Predict -   No, Yes, No, No
5 Fold 2
6   Test IID -  1, 6, 11, 16, 21
7   Actual -    No, Yes, No, No, No
8   Predict -   No, Yes, No, Yes, No
9 Fold 3
10  Test IID -  2, 7, 12, 17, 22
11  Actual -    Yes, No, Yes, No, Yes
12  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%

```

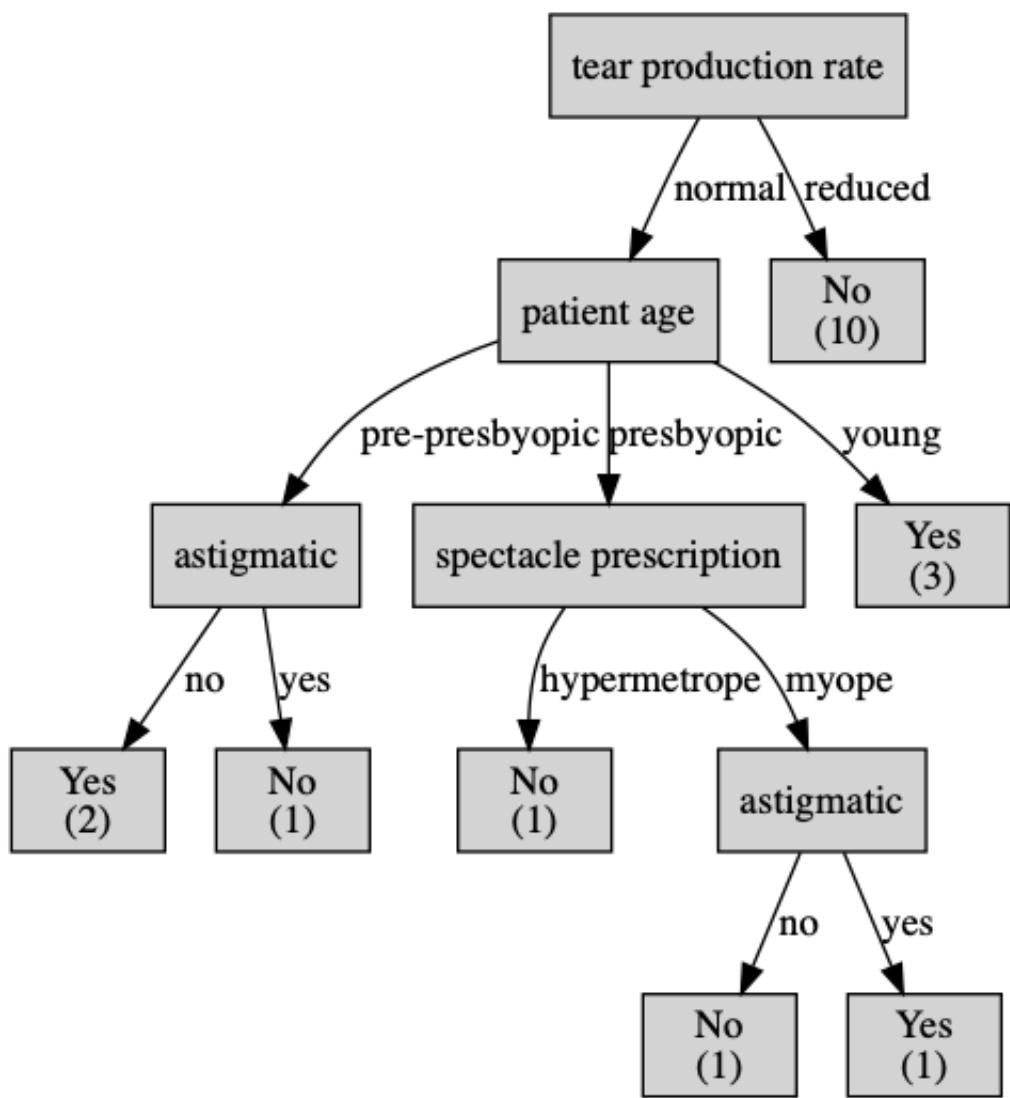
Fold 1



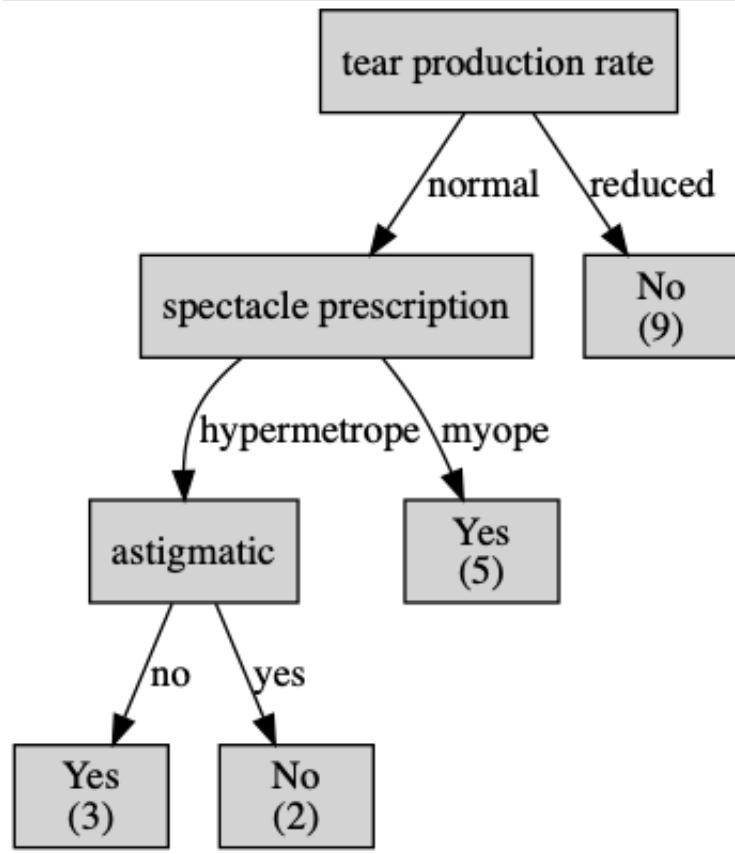
Fold 2



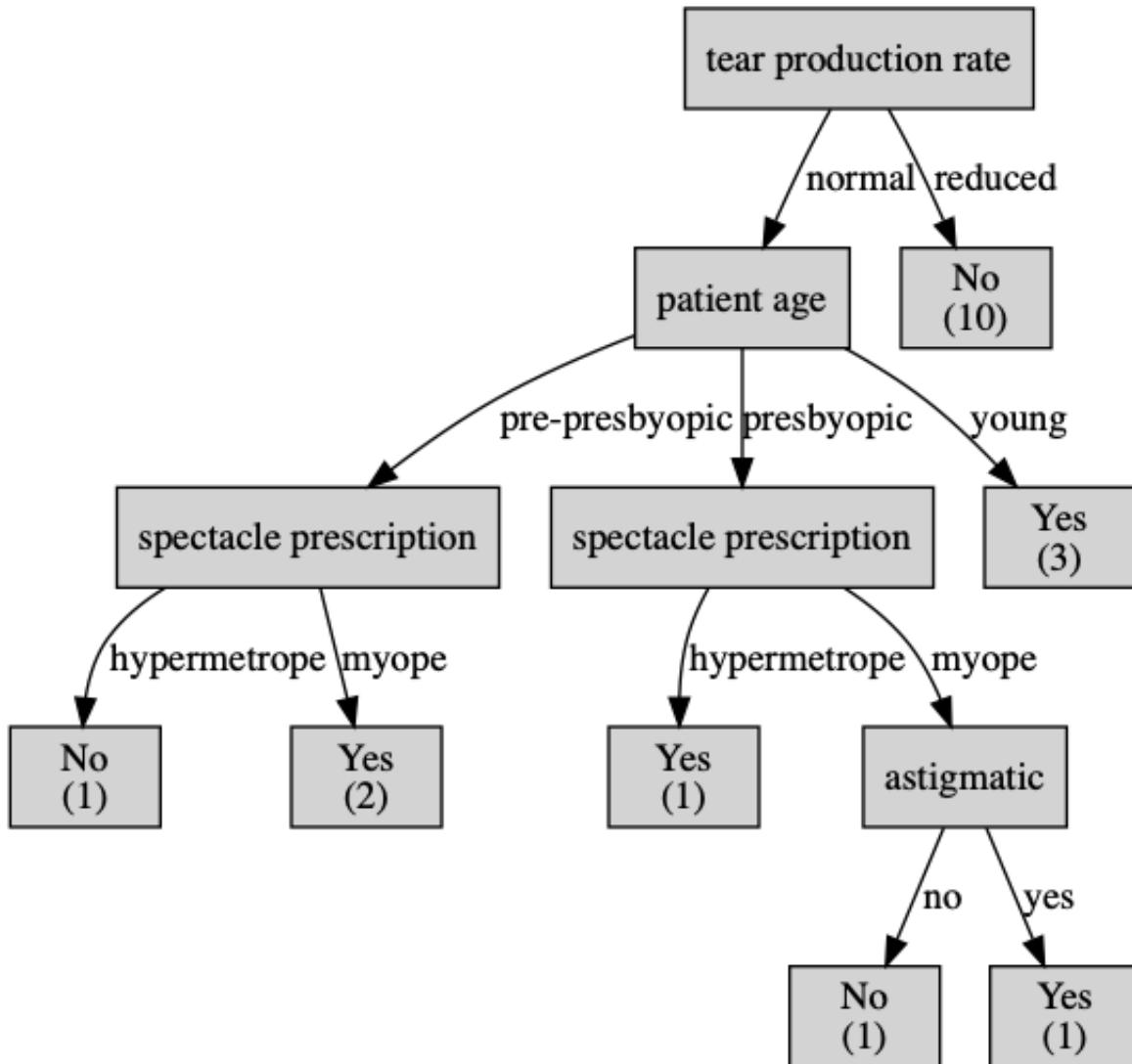
Fold 3



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,
      'hypermetrope':1, 'no':0, 'yes':1, 'normal':1, 'reduced':0}
4 d2={'Yes':1, 'No':0}
5 X = [[d1[X[i][j]] for j in range(len(X[0]))] for i in range(len(X))]
6 X = np.asarray(X)
7 Y = [d2[Y[i]] for i in range(len(Y))]
8 Y = np.asarray(Y)
  
```

```
1 nb = MultinomialNB(alpha=1)
2 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],
9 np.exp(nb.feature_log_prob_[i][j])))
```

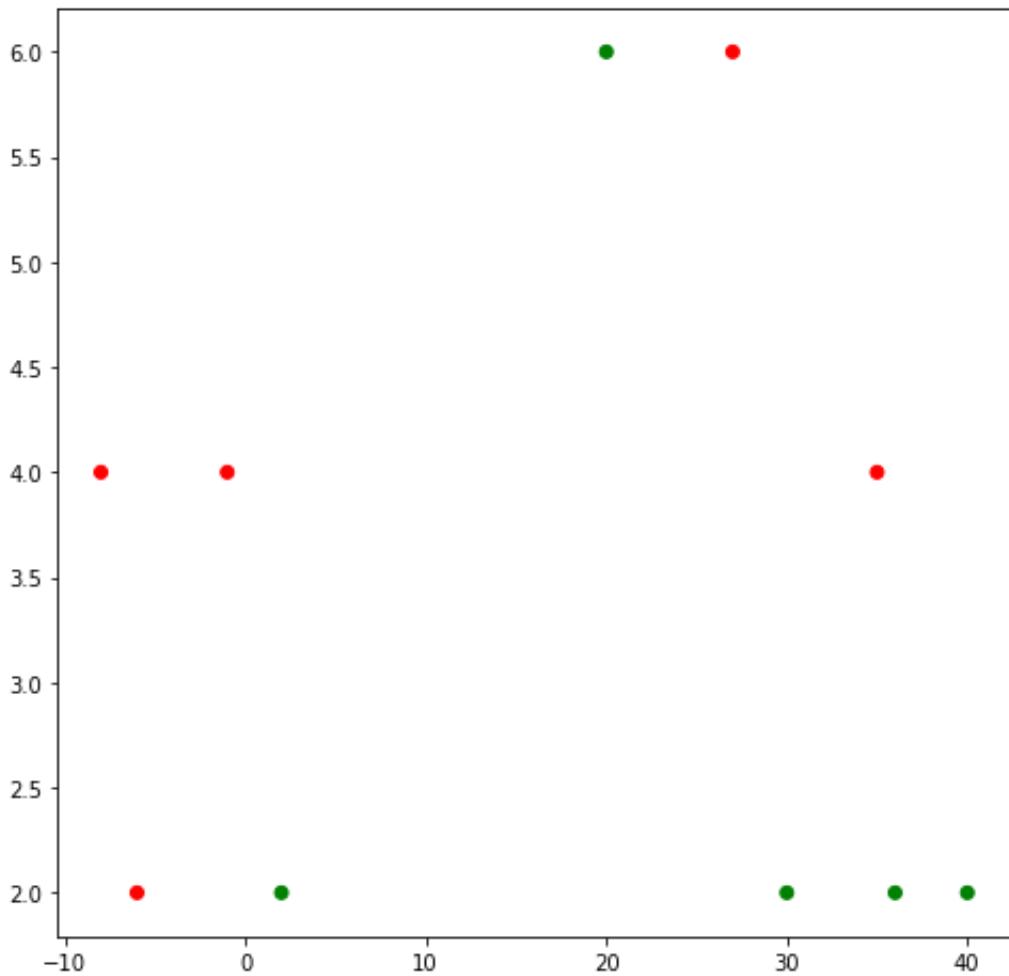
Model Details -

```
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
```

Q6

```
1 import matplotlib
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 ids = [1,2,3,4,5,6,7,8,9,10]
6 x = [27,-6,2,36,-8,40,35,30,20,-1]
7 y = [6,2,2,2,4,2,4,2,6,4]
8 label = [-1,-1,1,1,-1,1,-1,1,1,-1]
9 colors = ['red','green']
10
11 fig = plt.figure(figsize=(8,8))
12 plt.scatter(x, y, c=label, cmap=matplotlib.colors.ListedColormap(colors))
13
14 # cb = plt.colorbar()
15 # loc = np.arange(0,max(label),max(label)/float(len(colors)))
16 # cb.set_ticks(loc)
17 # cb.set_ticklabels(colors)
```

```
1 | <matplotlib.collections.PathCollection at 0x11b259128>
```



```
1 | import math
```

```
1 |
2 | ids = [1,2,3,4,5,6,7,8,9,10]
3 | x = [27,-6,2,36,-8,40,35,30,20,-1]
4 | y = [6,2,2,2,4,2,4,2,6,4]
5 | label = [-1,-1,1,1,-1,1,-1,1,1,-1]
6 | dataset = list(zip(ids,x,y,label))
```

```
1 | dataset
```

```
1 [(1, 27, 6, -1),  
2 (2, -6, 2, -1),  
3 (3, 2, 2, 1),  
4 (4, 36, 2, 1),  
5 (5, -8, 4, -1),  
6 (6, 40, 2, 1),  
7 (7, 35, 4, -1),  
8 (8, 30, 2, 1),  
9 (9, 20, 6, 1),  
10 (10, -1, 4, -1)]
```

```
1 #Helper Method to find euclidian distance between two datapoints  
2 def find_euclidian_distance(x,y):  
3     sum = 0  
4     sum += (x[1] - y[1]) ** 2  
5     sum += (x[2] - y[2]) ** 2  
6     return math.sqrt(sum)
```

```
1 find_euclidian_distance(dataset[0],dataset[1])
```

```
1 33.24154027718932
```

Q6(A)

3 Nearest datapoints for data point ID : 5

```
1 dist = []  
2 for i in range(0,10):  
3     dist.append((i+1,find_euclidian_distance(dataset[4],dataset[i])))  
4 # dist  
5 sorted(dist, key=lambda x: x[1])
```

```
1 [(5, 0.0),
2 (2, 2.8284271247461903),
3 (10, 7.0),
4 (3, 10.198039027185569),
5 (9, 28.071337695236398),
6 (1, 35.05709628591621),
7 (8, 38.05259518088089),
8 (7, 43.0),
9 (4, 44.04543109109048),
10 (6, 48.041648597857254)]
```

```
1 ### ANSWER: Therefore 3 nearest neighbors for datapoint ID: 5 are {2,10,3}
```

```
1 dist = []
2 for i in range(0,10):
3     dist.append((i+1,find_euclidian_distance(dataset[9],dataset[i])))
4
5 sorted(dist, key=lambda x: x[1])
```

```
1 [(10, 0.0),
2 (3, 3.605551275463989),
3 (2, 5.385164807134504),
4 (5, 7.0),
5 (9, 21.095023109728988),
6 (1, 28.071337695236398),
7 (8, 31.064449134018133),
8 (7, 36.0),
9 (4, 37.05401462729781),
10 (6, 41.048751503547585)]
```

Answer: Therefore 3 nearest neighbors for datapoint ID: 10 are {3,2,5}

Q6 (B)

```
1 # Experiment
```

```

2 from sklearn.neighbors import NearestNeighbors
3 import numpy as np
4 list_tuples = [(27, 6),
5 (-6, 2),
6 (2, 2),
7 (36, 2),
8 (-8, 4),
9 (40, 2),
10 (35, 4),
11 (30, 2),
12 (20, 6),
13 (-1, 4)]
14 labels = np.array([-1,-1,1,1,-1,1,-1,1,1,-1])
15 X = np.array(list_tuples)
16 # X = np.array([[-1, -1], [-2, -1], [-3, -2], [1, 1], [2, 1], [3, 2]])
17 # np.concatenate((X[:i],X[(i+1):]))

```

```

1 #Leave one out cross validation
2 total_errors = 0
3 for i in range(0,10):
4     nbrs = NearestNeighbors(n_neighbors=1, algorithm='auto',
5 metric='euclidean').fit(np.concatenate((X[:i],X[(i+1):])))
6     temp_labels = np.concatenate((labels[:i],labels[(i+1):]))
7     distances, indices = nbrs.kneighbors([X[i]])
8
9     for nearest_index in indices:
10         # print(f'Nearest Index: {nearest_index[0]}')
11         print(f'Testing Point: {dataset[i]} -> Nearest Point:
12 {dataset[nearest_index[0]]}')
13         if temp_labels[nearest_index[0]] != dataset[i][3]:
14             total_errors += 1
15
16 print(f'Leave-one-out-cross-validation error: {total_errors/10}')

```

```

1 Testing Point: (1, 27, 6, -1) -> Nearest Point: (7, 35, 4, -1)
2 Testing Point: (2, -6, 2, -1) -> Nearest Point: (4, 36, 2, 1)
3 Testing Point: (3, 2, 2, 1) -> Nearest Point: (9, 20, 6, 1)
4 Testing Point: (4, 36, 2, 1) -> Nearest Point: (6, 40, 2, 1)
5 Testing Point: (5, -8, 4, -1) -> Nearest Point: (2, -6, 2, -1)
6 Testing Point: (6, 40, 2, 1) -> Nearest Point: (4, 36, 2, 1)
7 Testing Point: (7, 35, 4, -1) -> Nearest Point: (4, 36, 2, 1)
8 Testing Point: (8, 30, 2, 1) -> Nearest Point: (1, 27, 6, -1)
9 Testing Point: (9, 20, 6, 1) -> Nearest Point: (1, 27, 6, -1)
10 Testing Point: (10, -1, 4, -1) -> Nearest Point: (3, 2, 2, 1)
11 Leave-one-out-cross-validation error: 0.7

```

Leave-one-out-cross-validation error: 0.7

Q6 (C)

```

1 # Create Folds:
2 def createfolds(num_folds):
3     folds = {i:{'train':[], 'test':[]} for i in range(1,num_folds+1)}
4     for i in range(1,num_folds+1):
5         for j in range(1,11):
6             if j % 5 == i-1:
7                 folds[i]['test'].append(j)
8             else:
9                 folds[i]['train'].append(j)
10    return folds

```

```

1 folds = createfolds(5)
2 folds
3

```

```

1 {1: {'test': [5, 10], 'train': [1, 2, 3, 4, 6, 7, 8, 9]}, 
2 2: {'test': [1, 6], 'train': [2, 3, 4, 5, 7, 8, 9, 10]}, 
3 3: {'test': [2, 7], 'train': [1, 3, 4, 5, 6, 8, 9, 10]}, 
4 4: {'test': [3, 8], 'train': [1, 2, 4, 5, 6, 7, 9, 10]}, 
5 5: {'test': [4, 9], 'train': [1, 2, 3, 5, 6, 7, 8, 10]}}

```

```
1 total_errors = 0
```

```
2  for i in range(1,6):
3      test_data = []
4      test_labels = []
5      train_data = []
6      train_labels = []
7
8      for index in folds[i]['train']:
9          # print(f'Index: {index}')
10         train_data.append(list_tuples[index-1])
11         train_labels.append(labels[index-1])
12
13     train_data_nparray = np.array(train_data)
14     print(f'Train Data: {train_data_nparray}')
15     # print(f'{train_labels}')
16
17     for index in folds[i]['test']:
18         # print(f'Index: {index}')
19         test_data.append(list_tuples[index-1])
20         test_labels.append(labels[index-1])
21
22     test_data_nparray = np.array(test_data)
23     # print(f'Test Data: {test_data_nparray}')
24
25     nbrs = NearestNeighbors(n_neighbors=3, algorithm='auto',
26                             metric='euclidean').fit(train_data_nparray)
27
28     distances, indices = nbrs.kneighbors(test_data_nparray)
29
30     # print(f'indices: {indices}\n')
31
32     k = 0
33     for nearest_indices in indices:
34         label_total = 0
35         for near_indice in nearest_indices:
36             # print(f'near_indice: {near_indice}')
37             label_total += train_labels[near_indice]
38
39             if label_total > 0:
40                 predicted_label = 1
41             if label_total < 0:
42                 predicted_label = -1
43
44             # print(f'predicted_label: {predicted_label}')
45             if predicted_label != test_labels[k]:
46                 total_errors += 1
```

```
47     k += 1
48
49 print(f'5-fold-cross-validation error: {total_errors/10}')
```

```
1 Train Data: [[27  6]
2   [-6  2]
3   [ 2  2]
4   [36  2]
5   [40  2]
6   [35  4]
7   [30  2]
8   [20  6]]
9 indices: [[1 2 7]
10  [2 1 7]]
11
12 near_indice: 1
13 near_indice: 2
14 near_indice: 7
15 predicted_label: 1
16 near_indice: 2
17 near_indice: 1
18 near_indice: 7
19 predicted_label: 1
20 Train Data: [[-6  2]
21   [ 2  2]
22   [36  2]
23   [-8  4]
24   [35  4]
25   [30  2]
26   [20  6]
27   [-1  4]]
28 indices: [[5 6 4]
29  [2 4 5]]
30
31 near_indice: 5
32 near_indice: 6
33 near_indice: 4
34 predicted_label: 1
35 near_indice: 2
36 near_indice: 4
37 near_indice: 5
38 predicted_label: 1
39 Train Data: [[27  6]
40   [ 2  2]
41   [36  2]]
```

```
42 [-8 4]
43 [40 2]
44 [30 2]
45 [20 6]
46 [-1 4]]
47 indices: [[3 7 1]
48 [2 5 4]]
49
50 near_indice: 3
51 near_indice: 7
52 near_indice: 1
53 predicted_label: -1
54 near_indice: 2
55 near_indice: 5
56 near_indice: 4
57 predicted_label: 1
58 Train Data: [[27 6]
59 [-6 2]
60 [36 2]
61 [-8 4]
62 [40 2]
63 [35 4]
64 [20 6]
65 [-1 4]]
66 indices: [[7 1 3]
67 [0 5 2]]
68
69 near_indice: 7
70 near_indice: 1
71 near_indice: 3
72 predicted_label: -1
73 near_indice: 0
74 near_indice: 5
75 near_indice: 2
76 predicted_label: -1
77 Train Data: [[27 6]
78 [-6 2]
79 [ 2 2]
80 [-8 4]
81 [40 2]
82 [35 4]
83 [30 2]
84 [-1 4]]
85 indices: [[5 4 6]
86 [0 6 5]]
87
```

```
88 near_indice: 5
89 near_indice: 4
90 near_indice: 6
91 predicted_label: 1
92 near_indice: 0
93 near_indice: 6
94 near_indice: 5
95 predicted_label: -1
96 5-fold-cross-validation error: 0.7
```

5-fold cross-validation error: 0.7

Q6 (D)

Based on the results of (B) and (C) we can not determine

1. Data points are very less.
2. For 5-fold cross-validation: testing set is size of 2 and training set size is of 8. Which is almost same as for leave-one-out (where ratio for testing is 1 and training 9).
3. So, We can not determine base