

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')
3 headers, test_x, test_y = data_and_headers('Data' + os.sep +
4         '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 - ' +
      str(train_y[train_y=='R'].shape[0]))
6 print('Number of observations in category M - ' +
      str(train_y[train_y=='M'].shape[0]))
7 print()
8 print('Testing Data')
9 print('Number of features - ' + str(test_x.shape[1]))
10 print('Number of target features - ' + str(test_y.shape[1]))
11 print('Number of observations - ' + str(test_x.shape[0]))
12 print('Number of observations in category R - ' +
      str(test_y[test_y=='R'].shape[0]))
13 print('Number of observations in category M - ' +
      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

```

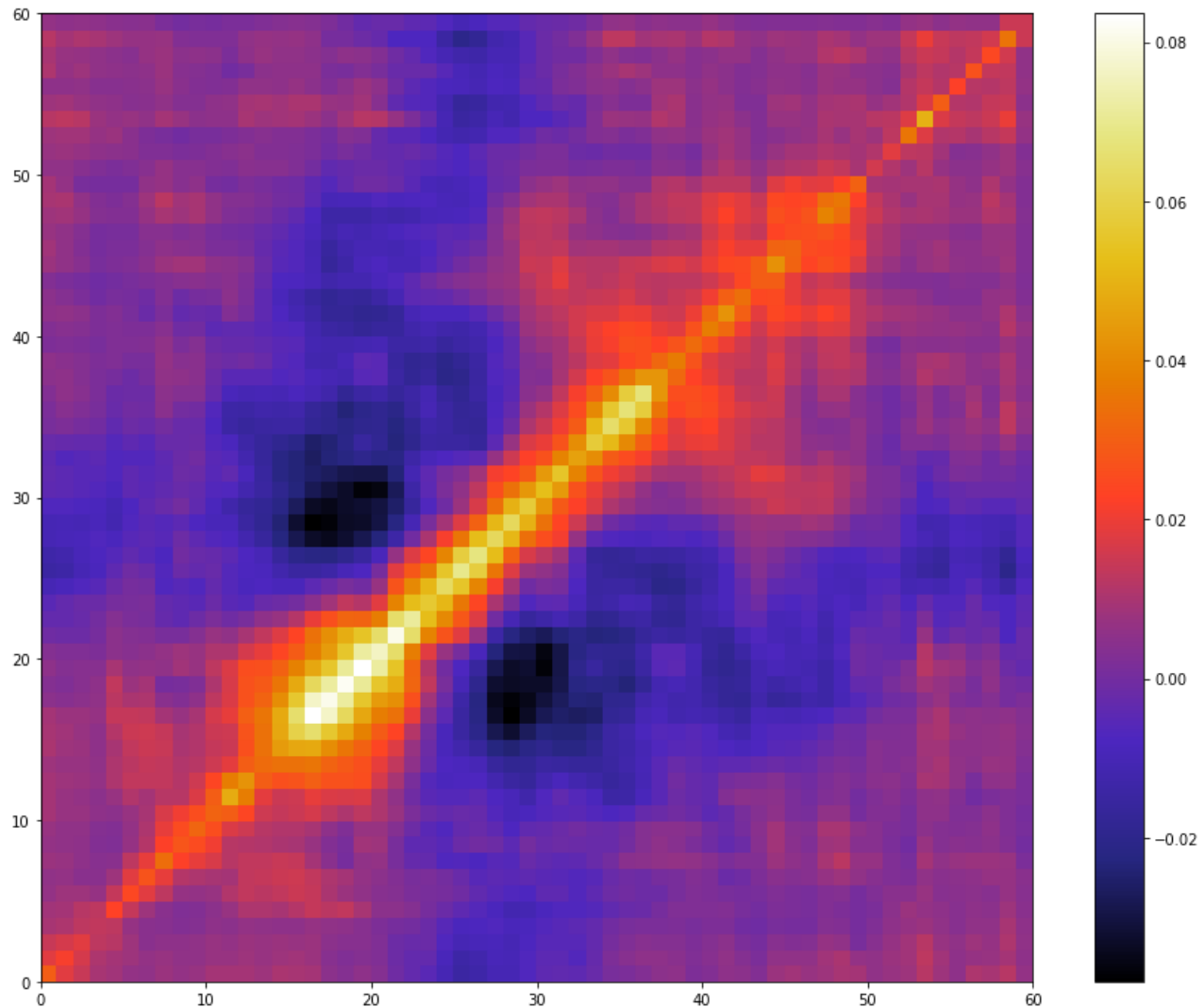
```

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



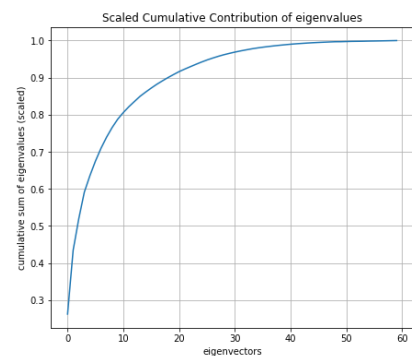
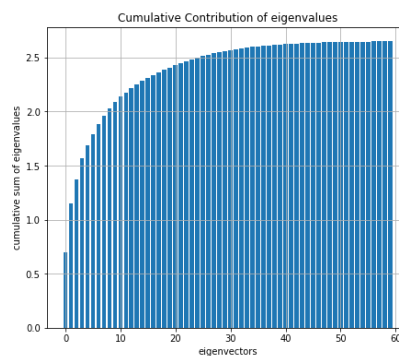
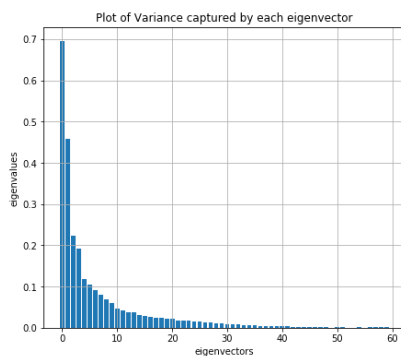
(ii) Eigenvalue and Eigenvectors

```
1 print('Size of covariance matrix - ' + str(covariance.shape))
2 w,v = np.linalg.eig(covariance)
3 print('Top 5 Eigenvalues - ' + ', '.join(['{:0.3f}'.format(x) for x in
4     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')

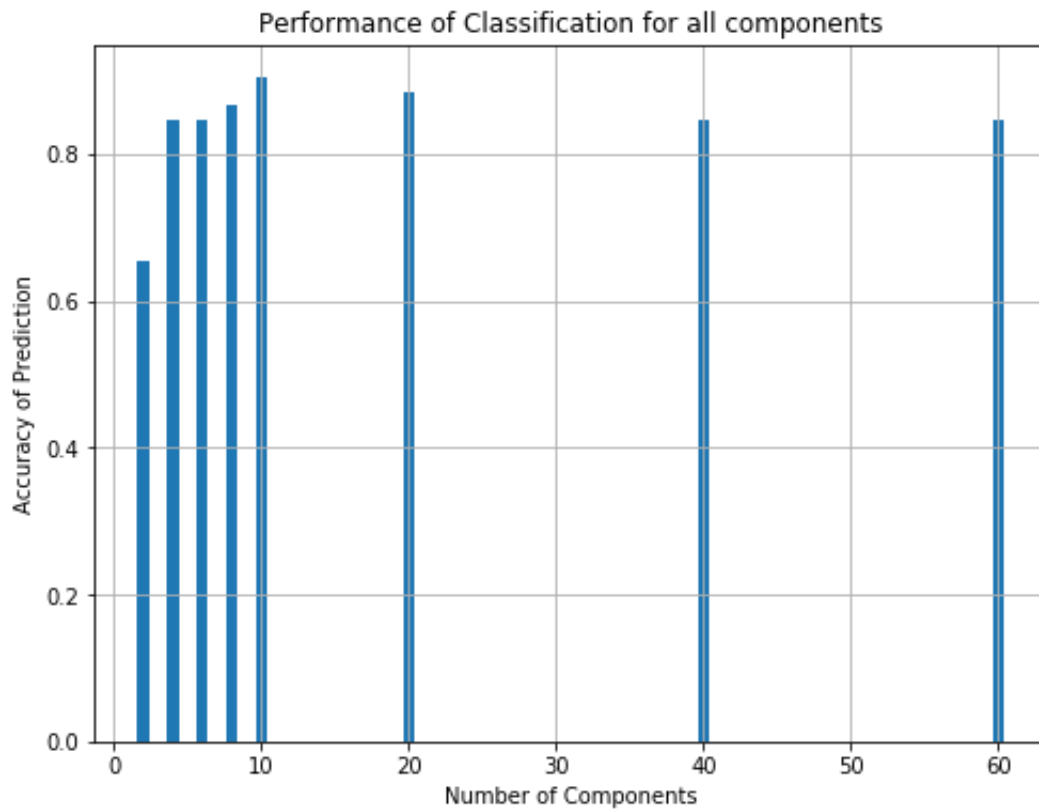
```

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),
  np.std(train_x, axis=0))
2 standard_test = standardize(test_x, np.mean(train_x, axis=0),
  np.std(train_x, axis=0))

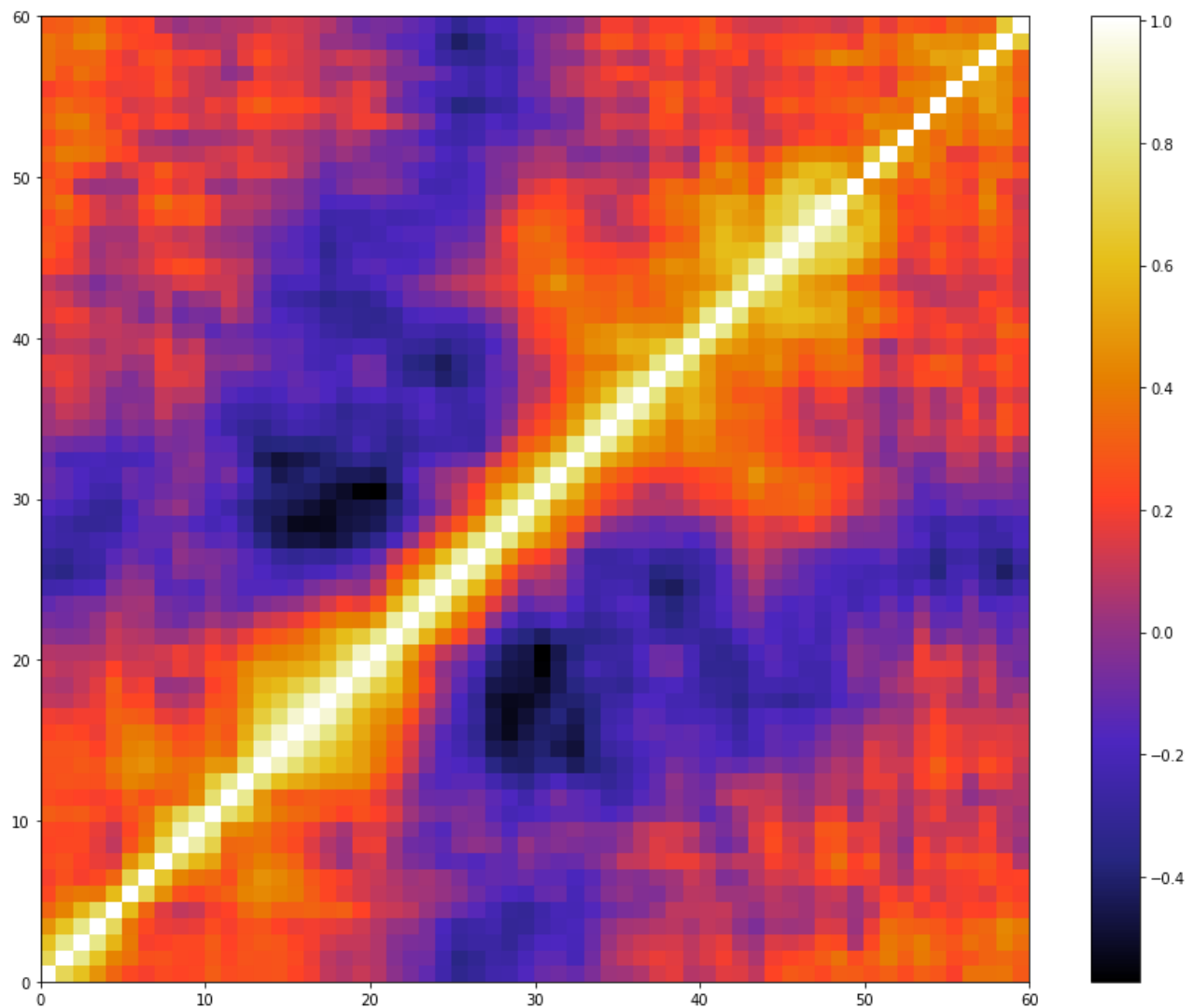
```

(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(['{:0.3f}'.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

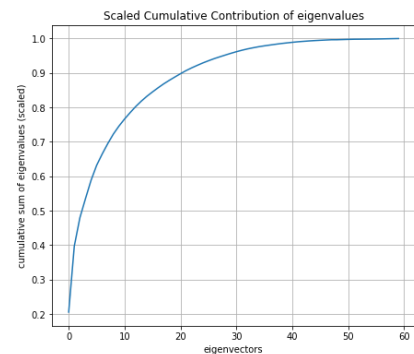
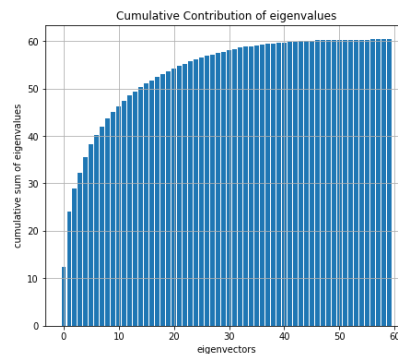
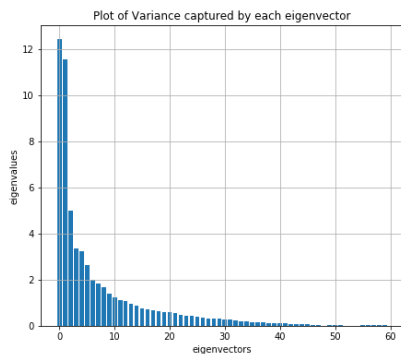
```

(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.

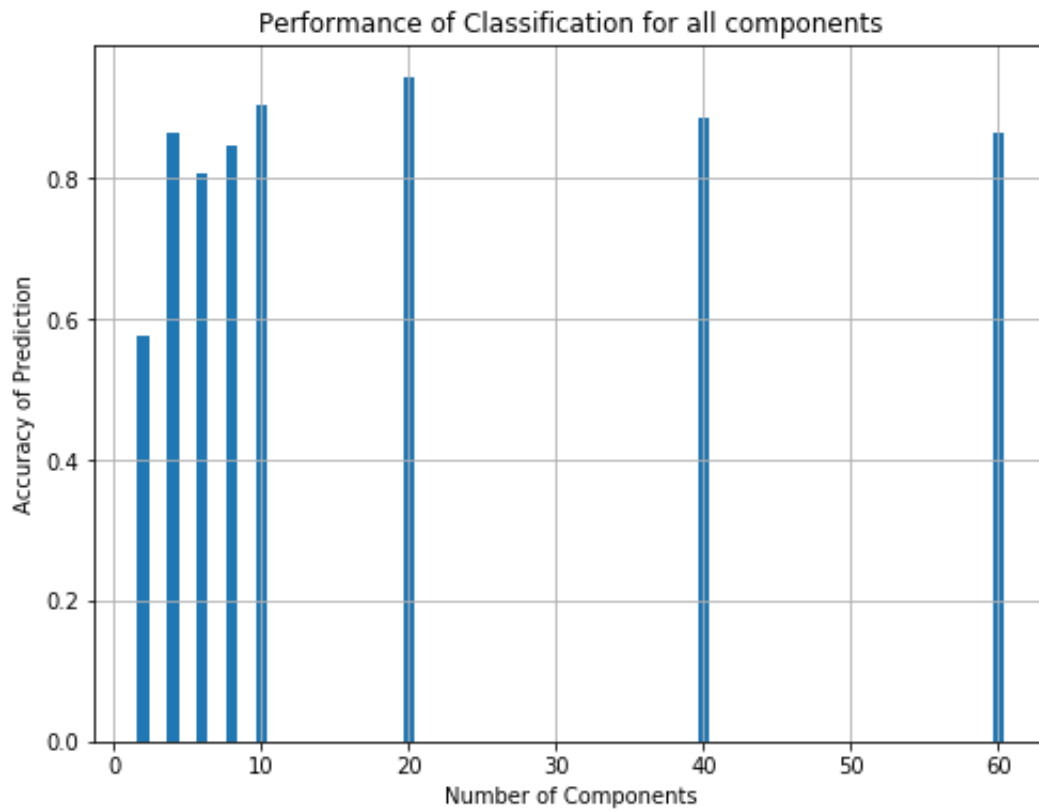
(iv) PCA with KNN

```
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\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.