

IE406 Machine Learning

Lab Assignment - 6

Group 28

201801015: Shantanu Tyagi
201801076: Shivani Nandani
201801407: Pratvi Shah
201801408: Arkaprabha Baerjee

Question 1

Given the following data use Principal Component Analysis to reduce the feature dimension from 3 to 1 also show eigen values. (Do manual calculation)

X ₁	X ₂	X ₃
8	14	3
13	9	6
4	3	15
7	2	1
19	8	4
5	18	11

Answer

In order to reduce dimensions from 3 to 1 we need to follow the below mentioned steps:

1. Calculate the mean and standard deviation for each feature/column(X_i) of A where,

$$A = \begin{bmatrix} 8 & 14 & 3 \\ 13 & 9 & 6 \\ 4 & 3 & 15 \\ 7 & 2 & 1 \\ 19 & 8 & 4 \\ 5 & 18 & 11 \end{bmatrix}$$

$$M_i(\text{Mean}) = \frac{\sum_{j=1}^6 A_{ji}}{6}$$

$$S_i(\text{Standard deviation}) = \sqrt{\frac{\sum_{j=1}^6 (A_{ji} - M_i)^2}{6}}$$

$$M = [9.33 \quad 9.00 \quad 6.67]$$

$$S = [5.185 \quad 5.6568 \quad 4.8534]$$

2. After we calculate the mean and variance of each column we standardize the elements of each column by subtracting the mean and dividing by the standard deviation

$$B_i = \frac{A_i - M}{S} \text{ for all } i \in [1, 6]$$

$$B = \begin{bmatrix} -0.257 & 0.88 & -0.755 \\ 0.707 & 0.00 & -0.137 \\ -1.028 & -1.06 & 1.717 \\ -0.449 & -1.23 & -1.1675 \\ 1.864 & -0.177 & -0.549 \\ -0.835 & 1.591 & 0.892 \end{bmatrix}$$

3. After reducing the columns to zero mean we can move on to calculate the covariance matrix which for our case will be of 3x3 dimension as there are 3 features

$$C = \text{Covariance}(B)$$

$$C_{ij} = \frac{B' \times B}{n - 1} \text{ where } n=6, B' = \text{transpose}(B)$$

$$C = \begin{bmatrix} 1.2 & -0.0477 & -0.5827 \\ -0.0477 & 1.2 & 0.0947 \\ -0.5827 & 0.0947 & 1.2 \end{bmatrix}$$

4. Now we find the eigen vectors and eigen values corresponding to this covariance matrix C. This was calculated using `numpy.linalg.eig(C)` and were found to be:

$$\text{Eigen values} = [1.79971072 \quad 0.61531314 \quad 1.18497614]$$

$$\text{Eigen vectors} = \begin{bmatrix} -0.69401025 & 0.70241371 & -0.15800237 \\ 0.16586365 & -0.05756869 & -0.98446691 \\ 0.70059904 & 0.70943698 & 0.07655164 \end{bmatrix}$$

5. As we need to reduce the 3 dimension to a single one we just retain the eigen vector corresponding to the maximum eigen value, because higher the eigen value higher is the variance in the direction of that eigen vector. Therefore, we only retain the second eigen vector in order to form the projection matrix W:

$$W = \begin{bmatrix} -0.69401025 \\ 0.16586365 \\ 0.70059904 \end{bmatrix}$$

6. Now in order to transform 3D to single dimension we project the standardized matrix(B) onto the projection matrix W:

$$B \times W = \begin{bmatrix} -0.257 & 0.88 & -0.755 \\ 0.707 & 0.00 & -0.137 \\ -1.028 & -1.06 & 1.717 \\ -0.449 & -1.23 & -1.1675 \\ 1.864 & -0.177 & -0.549 \\ -0.835 & 1.591 & 0.892 \end{bmatrix}_{6 \times 3} \times \begin{bmatrix} -0.69401025 \\ 0.16586365 \\ 0.70059904 \end{bmatrix}_{3 \times 1} = \begin{bmatrix} -0.20423593 \\ -0.58697404 \\ 1.74081123 \\ -0.71095189 \\ -1.70802694 \\ 1.46937756 \end{bmatrix}_{6 \times 1}$$

$$7. \text{ Final answer of PCA} = \begin{bmatrix} -0.20423593 \\ -0.58697404 \\ 1.74081123 \\ -0.71095189 \\ -1.70802694 \\ 1.46937756 \end{bmatrix}_{6 \times 1}$$

Observation/ Justification

For the above case we had to reduce the dimensions from 3 to 1 hence only the eigen-vector with highest eigen values was chosen while constructing projection matrix. But had it been 3D to 2D transformation we would choose the eigen vectors corresponding to the 2 highest eigen-values i.e., eigen-vectors 1 and 3 in our case. In order to regenerate the 3D data points post PCA we can just do $(\text{ans} \times W^T)$.

Question 2

A classic application of PCA is to project the 3-D point cloud onto a plane that could still retain the information or essence of the point distribution. Given a 3-D point cloud, estimate an optimal plane, onto which if the 3D data points when projected would still retain the essential information. We provide you with 5 different point clouds P_1.txt to P_5.txt containing the 3D coordinates of points in space. Perform PCA on these point clouds to obtain their projection on a 2D plane and visualize the results using python libraries (like matplotlib).

Answer

For this question we first load the 3D data points from each file. Then we standardize the data-points in order to perform principal component analysis. By performing PCA we reduce the 3D data to 2D by considering only those two features that cover maximum variance. Post this we regenerate the 3D data points by performing inverse transform to calculate the reconstruction error incurred due to feature selection(PCA).

Code

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 %matplotlib inline
5 from sklearn.decomposition import PCA
6 from sklearn.preprocessing import StandardScaler
7 import seaborn as sns
8 from sklearn.metrics import mean_squared_error as mse
9 from mpl_toolkits import mplot3d
10
11 #Function to standardize data
12 def standardize(df1):
13     scaler=StandardScaler()
14     df1=scaler.fit_transform(df1)
15     return pd.DataFrame(df1)
16
17 #Function to perform PCA
18 def pca_func(data):
19     pca = PCA(n_components=2,svd_solver='full')
20     pca=pca.fit(data)
21
22     #Transform data to 2D
23     data_encoded=pca.transform(data)
24
25     #Regenerate 3D data from 2 features
26     data_decoded=pd.DataFrame(pca.inverse_transform(data_encoded))
27     comp=pca.components_
28
29     #Principal components
30     print('Components:')
31     print(comp)
32
33     #Plot original data
34     fig = plt.figure(figsize=[8,5])
35     ax = plt.axes(projection='3d')
36     ax.plot3D(data[0],data[1], data[2], 'gray')
37     for i in range(len(pca.components_)):
38         exp_var = pca.explained_variance_[i]
39         component = pca.components_[i]
40         p_component_with_length = component*3*np.sqrt(exp_var)
41         PC = list(zip(pca.mean_,pca.mean_ + p_component_with_length))
42         ax.plot(PC[0],PC[1],PC[2], 'k')
43     plt.title('Original normalized data')
44     ax.set_xlabel('x')
45     ax.set_ylabel('y')
46     ax.set_zlabel('z')
47     plt.show()
48
49     #Plot regained data after performing PCA
50     fig = plt.figure(figsize=[8,5])
51     ax = plt.axes(projection='3d')
52     ax.plot3D(data_decoded[0],data_decoded[1], data_decoded[2], 'gray')
53     for i in range(len(pca.components_)):
54         exp_var = pca.explained_variance_[i]
```

```

55     component = pca.components_[i]
56     p_component_with_length = component*3*np.sqrt(exp_var)
57     PC = list(zip(pca.mean_,pca.mean_ + p_component_with_length))
58     ax.plot(PC[0],PC[1],PC[2], 'k')
59 plt.title('Regained data')
60 ax.set_xlabel('x')
61 ax.set_ylabel('y')
62 ax.set_zlabel('z')
63 plt.show()
64
65 #Plot 2D figure
66 plt.scatter(data_encoded[:,0],data_encoded[:,1])
67 plt.axhline(y=0,color='k')
68 plt.axvline(x=0,color='k')
69 plt.xlabel('PC1')
70 plt.ylabel('PC2')
71 plt.title('2D transformed data')
72 plt.grid(True)
73 plt.show()
74
75 #Calculate error
76 rmse = mse(data, data_decoded, squared=False)
77 print('RMSE between original and regained data after PCA: ',rmse)
78
79
80 #Read data
81 data1 = pd.read_csv('P_1.txt', sep=" ", header=None)
82 data2 = pd.read_csv('P_2.txt', sep=" ", header=None)
83 data3 = pd.read_csv('P_3.txt', sep=" ", header=None)
84 data4 = pd.read_csv('P_4.txt', sep=" ", header=None)
85 data5 = pd.read_csv('P_5.txt', sep=" ", header=None)
86
87 #Standardize the data for performing PCA
88 data1=standardize(data1)
89 data2=standardize(data2)
90 data3=standardize(data3)
91 data4=standardize(data4)
92 data5=standardize(data5)
93 cols=data1.columns
94
95 ### P_1
96 pca_func(data1)
97 ###P_2
98 pca_func(data2)
99 ###P_3
100 pca_func(data3)
101 ###P_4
102 pca_func(data4)
103 ###P_5
104 pca_func(data5)

```

Listing 1: Question 2

Result

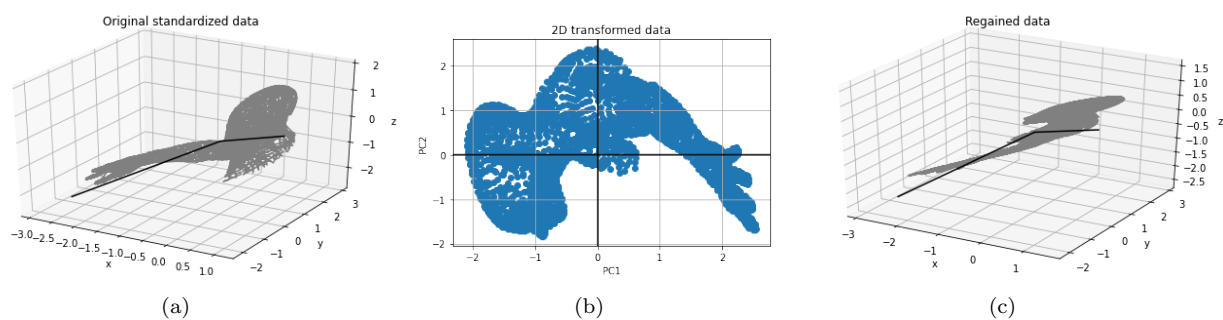


Figure 1: P_1.txt

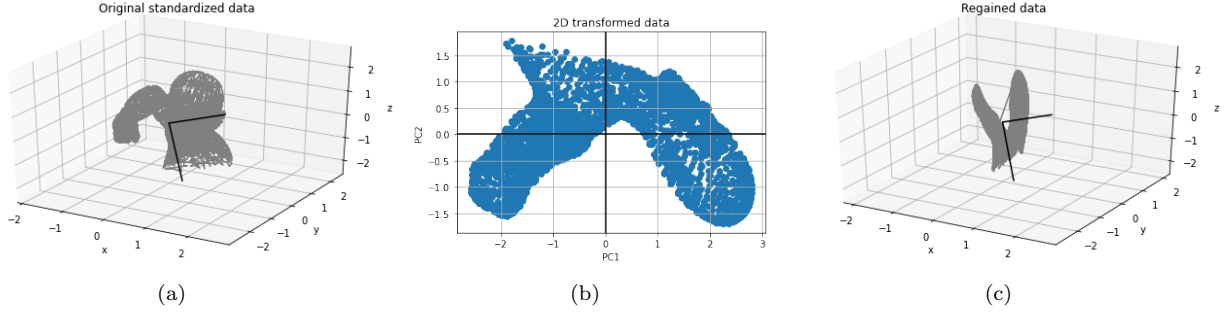


Figure 2: P_2.txt

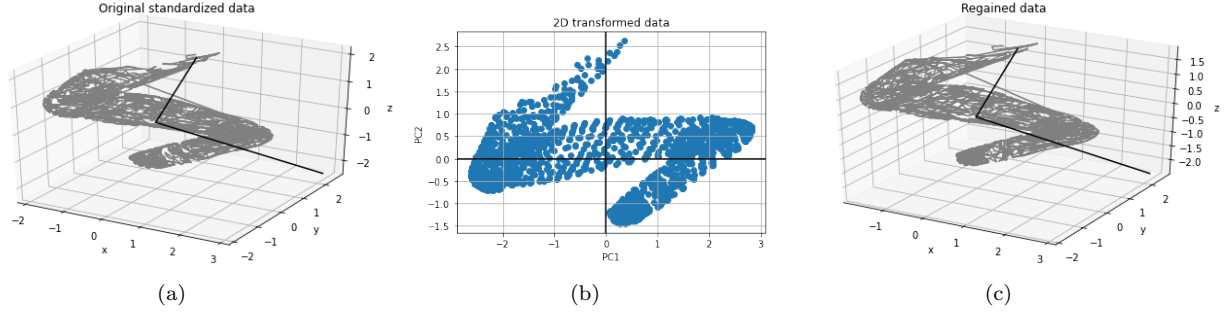


Figure 3: P_3.txt

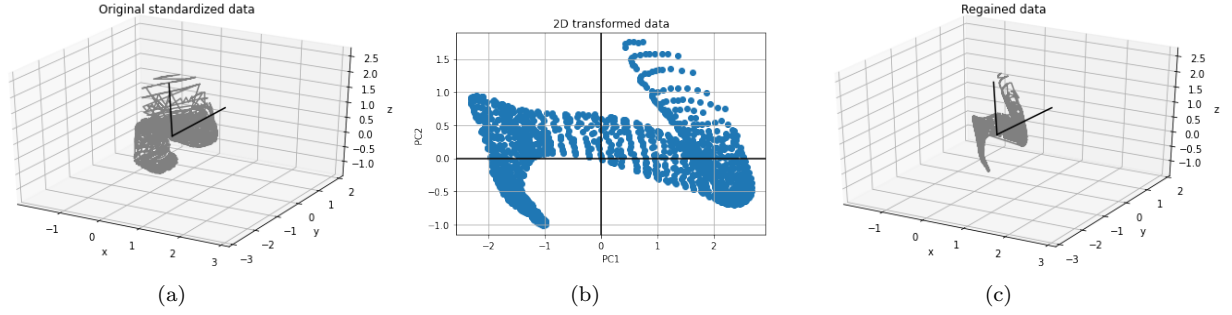


Figure 4: P_4.txt

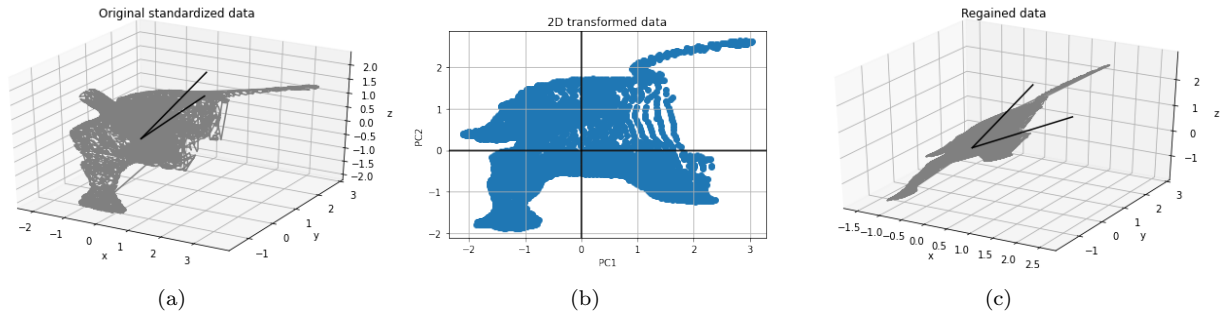


Figure 5: P_5.txt

In the above plots the black vectors represent the principal components. After reconstructing the final 3D plot post PCA we calculate the reconstruction error by evaluating RMSE and R^2 score for each point clouds.

File Name	RMSE	R^2 score
P_1	0.22275	0.95038
P_2	0.33047	0.89078
P_3	0.217206	0.952821
P_4	0.20048	0.959806
P_5	0.438327	0.807869

Table 1: RMSE and R^2 score

Observation/ Justification

We could draw the following observations:

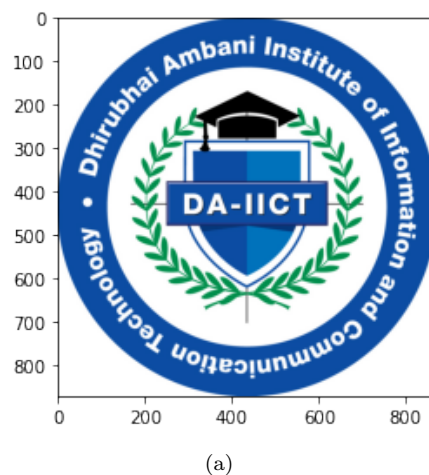
- The first principal component we obtain are majorly in descending order of the explained variance along that direction.
- The reconstructed 3D image post PCA shows that the points majorly lie on the plane formed by PC1 and PC2 and hence paves the way for the reconstruction error.
- RMSE values are low and R^2 values are considerable good (above 0.8) suggesting that performing PCA did not have a significant loss in information hence, indicating the goodness of fit.

Question 3

Another classic example of PCA comes in image compression. The human eye cannot perceive the minute and frequent changes in the image. A typical smartphone camera takes a 5MP image on average. We will utilize PCA and analyze how the reduction in the no. of principal components or dimensions affects the visual quality of the image. Given an image, apply PCA to investigate the effect of reducing the dimensions on the visual quality of the image. You will use the given image and show the analysis.

Answer

For this question we first load the image file. Then we split it into its RGB channels and then standardize the data in order to perform principal component analysis on each channel. By performing PCA we reduce the dimensions significantly by considering only those two features that cover maximum variance. Post this we regenerate the image by performing inverse transform to calculate the reconstruction error incurred due to feature selection (PCA).



(a)
Figure 6: Original image

Code

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.decomposition import PCA
5 from sklearn.preprocessing import StandardScaler
6 import cv2
7 from scipy.stats import stats
8 import matplotlib.image as mpimg
9 from sklearn.metrics import mean_squared_error as mse
10
11 img = cv2.cvtColor(cv2.imread('da_image.jpg'), cv2.COLOR_BGR2RGB)
12 plt.imshow(img)
13 plt.show()
14 print(img.shape)
15
16 #Splitting into channels
17 blue, green, red = cv2.split(img)
```

```

18
19 # Plotting the images
20 fig = plt.figure(figsize = (15, 7.2))
21 fig.add_subplot(131)
22 plt.title("Blue Channel")
23 plt.imshow(blue)
24 fig.add_subplot(132)
25 plt.title("Green Channel")
26 plt.imshow(green)
27 fig.add_subplot(133)
28 plt.title("Red Channel")
29 plt.imshow(red)
30 plt.show()
31
32 # ensuring data is between 0 to 1
33 df_blue = blue/255
34 df_green = green/255
35 df_red = red/255
36
37 blue_variance_ratios = []
38 green_variance_ratios = []
39 red_variance_ratios = []
40
41 # Doing PCA for all 3 channels
42 for i in range(10,870,10):
43     pca_b = PCA(n_components=i)
44     pca_b.fit(df_blue)
45     pca_g = PCA(n_components=i)
46     pca_g.fit(df_green)
47     pca_r = PCA(n_components=i)
48     pca_r.fit(df_red)
49     trans_pca_b = pca_b.transform(df_blue)
50     trans_pca_g = pca_g.transform(df_green)
51     trans_pca_r = pca_r.transform(df_red)
52     blue_variance_ratios.append(sum(pca_b.explained_variance_ratio_))
53     green_variance_ratios.append(sum(pca_g.explained_variance_ratio_))
54     red_variance_ratios.append(sum(pca_r.explained_variance_ratio_))
55
56 x_line = []
57 for i in range(10,870,10):
58     x_line.append(i)
59 plt.title("Blue Channel")
60 plt.xlabel("Dimensions Considered")
61 plt.ylabel("Variance Explained")
62 plt.plot(x_line, blue_variance_ratios)
63 plt.grid()
64 plt.show()
65 plt.title("Green Channel")
66 plt.xlabel("Dimensions Considered")
67 plt.ylabel("Variance Explained")
68 plt.plot(x_line, green_variance_ratios, color='green')
69 plt.grid()
70 plt.show()
71 plt.title("Red Channel")
72 plt.xlabel("Dimensions Considered")
73 plt.ylabel("Variance Explained")
74 plt.plot(x_line, red_variance_ratios, color='red')
75 plt.grid()
76 plt.show()
77
78 def CompressAndDisplayImage(dimensions):
79     pca_b = PCA(n_components=dimensions)
80     pca_b.fit(df_blue)
81     trans_pca_b = pca_b.transform(df_blue)
82     pca_g = PCA(n_components=dimensions)
83     pca_g.fit(df_green)
84     trans_pca_g = pca_g.transform(df_green)
85     pca_r = PCA(n_components=dimensions)
86     pca_r.fit(df_red)
87     trans_pca_r = pca_r.transform(df_red)
88
89     print(f"Blue Channel Variance: {sum(pca_b.explained_variance_ratio_)}")
90     print(f"Green Channel Variance: {sum(pca_g.explained_variance_ratio_)}")
91     print(f"Red Channel Variance: {sum(pca_r.explained_variance_ratio_)}")
92

```

```

93 b_arr = pca_b.inverse_transform(trans_pca_b)
94 g_arr = pca_g.inverse_transform(trans_pca_g)
95 r_arr = pca_r.inverse_transform(trans_pca_r)
96
97 img_reduced= (cv2.merge((b_arr, g_arr, r_arr)))
98
99 fig = plt.figure(figsize = (10, 7.2))
100 plt.title(f'Post PCA ({dimensions} dimensions)')
101 plt.imshow(img_reduced)
102 plt.show()
103
104 print(mse(df_blue , b_arr , squared=False))
105 print(mse(df_green , g_arr , squared=False))
106 print(mse(df_red , r_arr , squared=False))
107
108 CompressAndDisplayImage(200)
109 CompressAndDisplayImage(100)
110 CompressAndDisplayImage(50)
111 CompressAndDisplayImage(10)
112 CompressAndDisplayImage(5)
113 CompressAndDisplayImage(2)

```

Listing 2: Question 3

Result

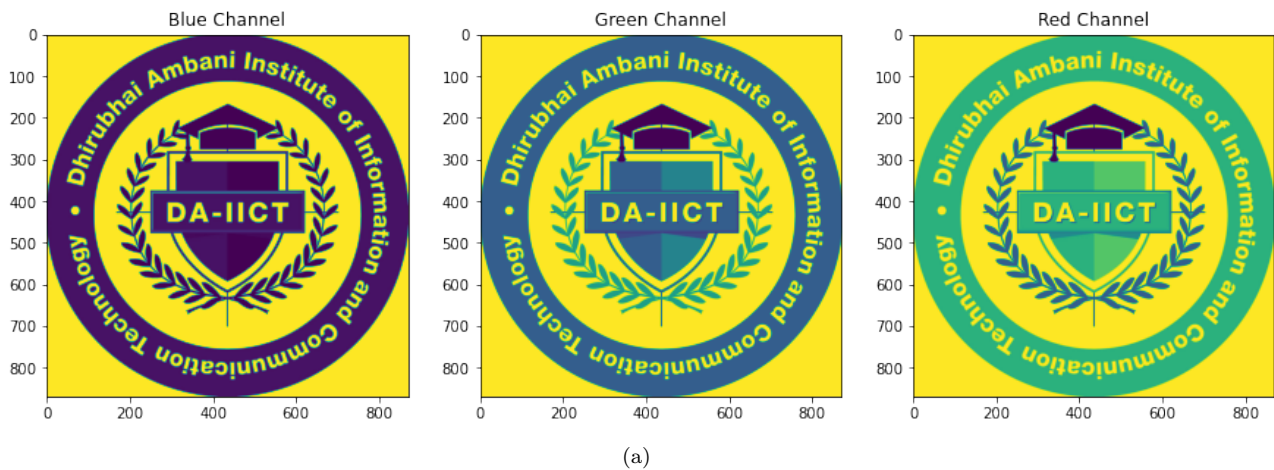


Figure 7: RGB channels of the image

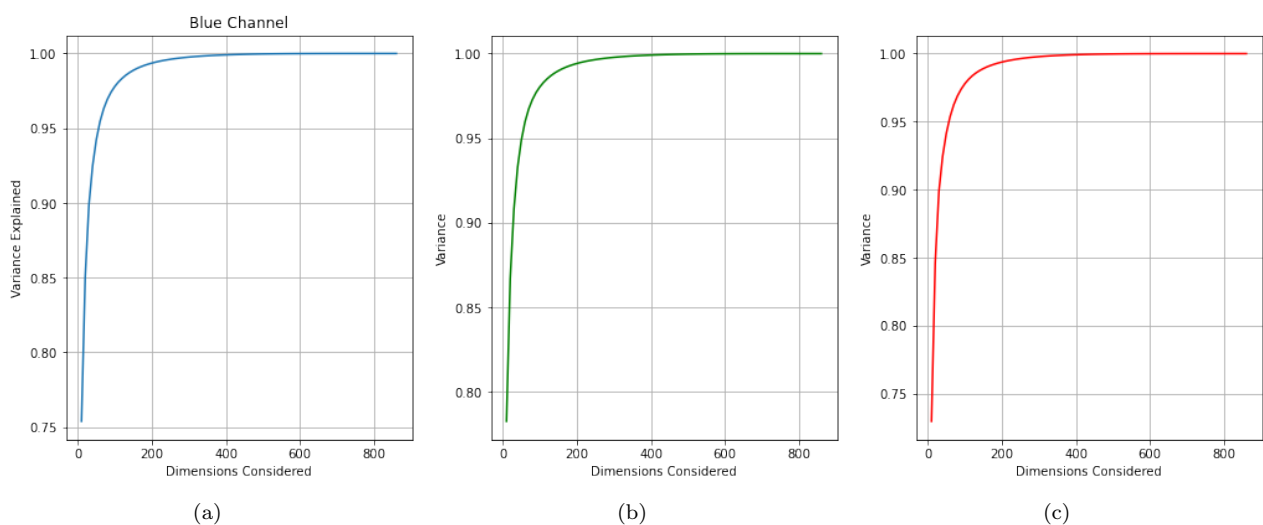


Figure 8: Variance VS Dimensions plots of each color channel

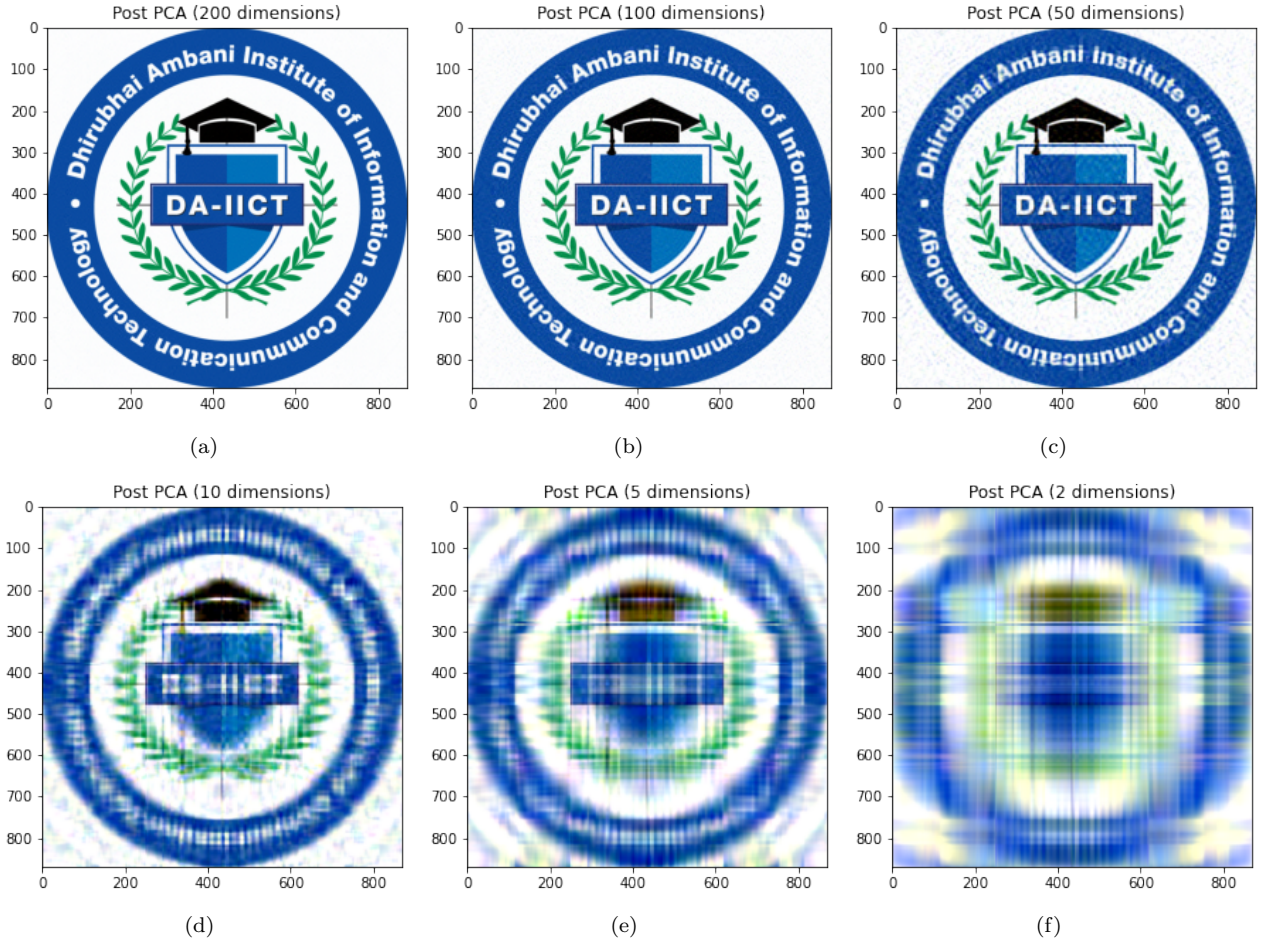


Figure 9: Reconstructed images considering lesser dimensions

Dimensions	Blue Variance Explained	Green Variance Explained	Red Variance Explained
200	0.9935648226837049	0.9941616565664149	0.9937959902748918
100	0.9779432517901917	0.9802967864880041	0.977847744796052
50	0.9423032980246341	0.9487310465998462	0.9413689828250666
10	0.7538841839537315	0.7825686583827102	0.7297300272813159
5	0.628864418099541	0.6555120231603824	0.5761178869817015
2	0.3832129939545568	0.42068274651318016	0.33822686728695284

Table 2: Variance Explained

Dimensions	Blue RMSE	Green RMSE	Red RMSE
200	0.03628799070104658	0.025171460512205462	0.0175091265028994
100	0.0670588181377301	0.046157278873273956	0.03294640724137046
50	0.10835934101538819	0.07451383636470789	0.05323306621979239
10	0.22337024018093732	0.15282446891207815	0.11343895031462171
5	0.2756272089262313	0.19334849728941825	0.14168458679812998
2	0.35666386847049236	0.2529646395655778	0.1778214097607088

Table 3: RMSE

Observation/ Justification

We could draw the following observations:

- The first principal component we obtain are majorly in descending order of the explained variance along that direction.

- We observe that 100 dimensions give us nearly 98% variance on all color channels. As the dimensions are reduced below 50, we see a significant difference in the reconstructed image.
- RMSE values are low for dimensions more than 50 suggesting that performing PCA did not have a significant loss in information hence, indicating the goodness of fit.

Question 4

Load the Data matrix `faceimages.mat` given to you. There are 400 face images of size 112×92 .

- Last column of the data is the class label. Column 1 to 10304 are 112×92 pixel.
- Each row represents one image.
- Construct the data matrix and Mean-center the data.
- Construct covariance matrix for the data.
- Solve eigenvalue problem to find projection matrix.
 - Find low dimensional representation of Face images
 - Draw energy curve of PCA for the given data
 - Reconstruct any face image with 5, 10, 50 and 100 principal components.
 - Compute the reconstruction error for the same.

Answer

For this question, we load the `faceimages.mat` and remove the labels column. The remaining data (400×10304) is mean-centered. Using inbuilt functions, we calculate covariance, eigenvalues and eigenvectors. We then perform PCA for varying number of components using `sklearn.decomposition.PCA` and calculate the error (in terms of RMSE) for the reconstructed images.

Code

```

1 import h5py
2 import numpy as np
3 from numpy.core.einsumfunc import _parse_possible_contraction
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 from sklearn.base import _pprint
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.decomposition import PCA
9
10 f = h5py.File('faceimages.mat', 'r')
11 data = f.get('data')
12 data = np.array(data).T
13 labels = data[:, -1]
14 data = data[:, :-1]
15
16 print(data.shape)
17
18 print(np.unique(labels))
19
20 data = data - np.mean(data, axis=0)
21 print(data.shape)
22
23 pca = PCA()
24 pca.fit_transform(data.T)
25 pca_variance_ratio_ = pca.explained_variance_ratio_
26 print(pca_variance_ratio_.shape)
27
28 plt.figure(figsize=(8, 6))
29 plt.plot(pca_variance_ratio_, 'b--')
30 plt.ylabel('Explained Variance ratio')
31 plt.xlabel('Principal components')
32 plt.show()
33

```

```

34 cov_matrix = (np.dot(data.T, data))/(np.size(data, 0)-1)
35 print(cov_matrix.shape)
36
37 print(cov_matrix)
38
39 eigenvalues, eigenvectors = np.linalg.eig(cov_matrix)
40 print(eigenvalues.shape, eigenvectors.shape)
41
42 print(eigenvalues)
43
44 print(eigenvectors)
45
46
47 def n_compPCA(feats, n_comp):
48     pca = PCA(n_components=n_comp)
49     print('Explained variance ratio: ', sum(pca.explained_variance_ratio_))
50     return pca
51
52
53 """n_components = 5"""
54
55 pca = n_compPCA(data, 5)
56
57 trans_data = pca.inverse_transform(pca.fit_transform(data))
58
59 print(trans_data.shape)
60
61 plt.imshow(trans_data[0].reshape((92, 112)).T, cmap=plt.cm.gray)
62 plt.show()
63
64 reconstruction_error = mean_squared_error(data, trans_data, squared=False)
65 print(reconstruction_error)
66
67 """n_components = 10"""
68
69 pca = n_compPCA(data, 10)
70
71 trans_data = pca.inverse_transform(pca.fit_transform(data))
72
73 print(trans_data.shape)
74
75 plt.imshow(trans_data[0].reshape((92, 112)).T, cmap=plt.cm.gray)
76 plt.show()
77
78 reconstruction_error = mean_squared_error(data, trans_data, squared=False)
79 print(reconstruction_error)
80
81 """n_components = 50"""
82
83 pca = n_compPCA(data, 50)
84
85 trans_data = pca.inverse_transform(pca.fit_transform(data))
86
87 print(trans_data.shape)
88
89 plt.imshow(trans_data[0].reshape((92, 112)).T, cmap=plt.cm.gray)
90 plt.show()
91
92 reconstruction_error = mean_squared_error(data, trans_data, squared=False)
93 print(reconstruction_error)
94
95 """n_components = 100"""
96
97 pca = n_compPCA(data, 100)
98
99 trans_data = pca.inverse_transform(pca.fit_transform(data))
100
101 print(trans_data.shape)
102
103 plt.imshow(trans_data[0].reshape((92, 112)).T, cmap=plt.cm.gray)
104 plt.show()
105
106 reconstruction_error = mean_squared_error(data, trans_data, squared=False)
107 print(reconstruction_error)
108

```

```

109 """n_components = 400"""
110
111 pca = n_compPCA(data, 400)
112
113 trans_data = pca.inverse_transform(pca.fit_transform(data))
114
115 print(trans_data.shape)
116
117 plt.imshow(trans_data[0].reshape((92, 112)).T, cmap=plt.cm.gray)
118 plt.show()
119
120 reconstruction_error = mean_squared_error(data, trans_data, squared=False)
121 print(reconstruction_error)

```

Listing 3: Question 4

Result

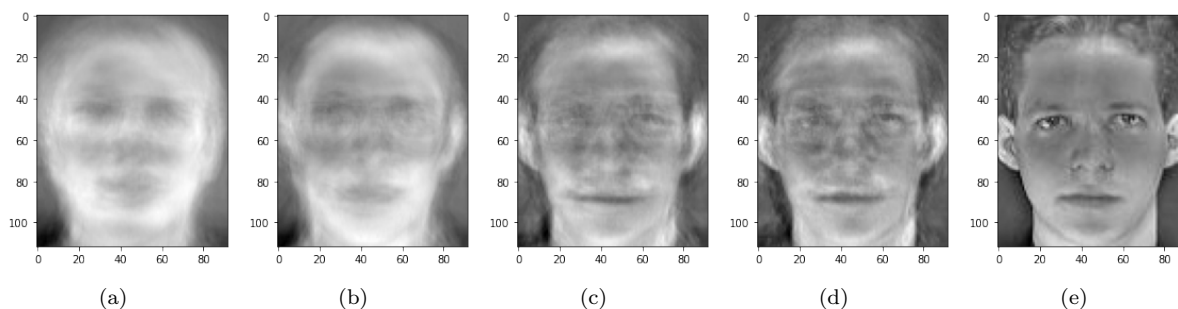


Figure 10: Reconstructed images for varying `n_components`

(a) `n_components=5`, (b) `n_components=10`, (c) `n_components=50`, (d) `n_components=100`, (e) `n_components=400`

n_components	Explained Variance Ratio	RMSE
5	0.48052	28.40129
10	0.59959	24.93487
50	0.81583	16.91369
100	0.88995	13.06288
400	0.99998	1.265848×10^{-13}

Table 4: Explained Variance Ratio and RMSE

Observation/ Justification

We could draw the following observations:

- From Figure 11 and Table 4, it is clear that as the number of components increase, the error reduces and the reconstructed image is closer to the actual image before PCA.
- With 400 components, which is much less than half the total number of features, the reconstruction error is almost negligible.

Question 5

Apply PCA on the MNIST dataset. Plot cumulative sum of variance with no of components and find the minimum no. of component for 85% variance.

Answer

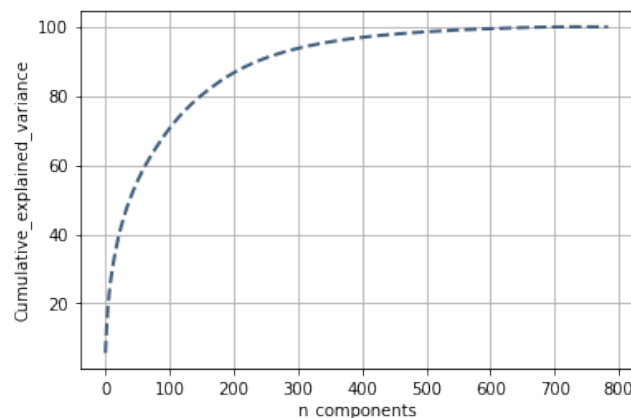
In this question the sklearn PCA function has been used to plot the cumulative variance plot and to find out the number of components required to explain 85% variance.

Code

```
1 import numpy as np
2 import pandas as pd
3 from matplotlib import pyplot as plt
4 from sklearn.datasets import fetch_openml
5 import seaborn as sns
6 %matplotlib inline
7 from sklearn.preprocessing import StandardScaler
8 from sklearn import decomposition
9
10
11 mnist = fetch_openml('mnist_784')
12 x = mnist.data
13 y = mnist.target
14 normalizedData = StandardScaler().fit_transform(x)
15 print(normalizedData.shape)
16
17 pca = decomposition.PCA()
18 pca.n_components = 784
19 pca_data = pca.fit_transform(normalizedData)
20 percentage_var_explained = pca.explained_variance_ / np.sum(pca.explained_variance_) * 100
21 cum_var_explained = np.cumsum(percentage_var_explained)
22
23 plt.figure(1, figsize=(6, 4))
24 plt.clf()
25 plt.plot(cum_var_explained, linewidth=2, linestyle='--', color = '#3F5D7D')
26 plt.axis('tight')
27 plt.grid()
28 plt.xlabel('n_components')
29 plt.ylabel('Cumulative_explained_variance')
30 plt.show()
31
32 pca = decomposition.PCA(n_components=0.85)
33 principalComponents = pca.fit_transform(normalizedData)
34 print("Net % Variance explained by Principal Components is :" + str(sum(pca.
35     explained_variance_ratio_)*100))
36 print("No of components to explain " + str(85)+"% Variance is" ,pca.explained_variance_ratio_.
37     shape)
```

Listing 4: Question 5

Result



(a)

Figure 11: Cumulative variance plot

Observation/ Justification

- One can observe the initial few principal components explain the most variance. Latter principal components explain very little variance.
- We require the first 186 principal components to explain 85% variance.