

BLG 454E Learning From Data (Spring 2018)

Homework 3

Onat Sahin - 150150129

1 Question 1

The code for this question can be found in `pca.py` .

- a)** Main motivations for reducing a dataset's dimensionality are eliminating erroneous features and visualizing the data.
- b)** We can train two separate classifiers with the original data and the reduced data. Then we can compare the test results of both classifiers to see if the reduced data retained the important features.
- c)** It seems like the 2d data is projected on to the y-axis. However, it is mainly the x coordinate that determines the class. Therefore, performance of the given PCA is low. The data should be projected on to the x-axis.
- d)** There are several drawbacks of PCA. First of all, PCA assumes that the principle components of the given data are linear combinations of the original features and these principle components are orthogonal. These assumptions are not necessarily true all the time. If they are false, PCA fails to give an accurate representation of the original data. Also, PCA uses variance as a measurement of the importance of a feature, which may not be the case. Finally, if the number of dimensions the data is going to be reduced to is not carefully selected, some important features may be lost and PCA may not give an accurate representation of the original data.
- e)** Figure 1 shows my implementation of `pca`, which is included in "`pca.py`". Figure 2 shows the output of "`pca.py`". For some reason I could not understand, my output is the flipped version of the example output given in the question.

```

1 #Onat Şahin 150150129~
2 import numpy as np~
3 import matplotlib.pyplot as plt~
4 import random~
5 ~
6 def readData(dataFile): #Funtion to read the data to a list~
7     with open(dataFile) as file:~
8         list = [line.strip('\n') for line in file.readlines()]~
9         list = [line.split(',') for line in list]~
10        for instance in list:~
11            for i in range(len(instance)):~
12                instance[i] = int(instance[i])~
13        ~
14        features = []~
15        labels = []~
16        ~
17        for instance in list: #Remove the labels from features to a sepearte list~
18            features.append(instance[0:len(instance)-1])~
19            labels.append(instance[-1])~
20        ~
21        return features, labels~
22        ~
23 features, labels = readData('data.txt')~
24 covariance = np.cov(features, rowvar=False) #Finding covariance matrix of features~
25 eigenvals, eigenvecs = np.linalg.eig(covariance) #Finding eigenvalues and eigenvectors of the covariance matrix~
26 index = eigenvals.argsort()[::-1] #Sorting eigenvalues in descending order~
27 eigenvals = eigenvals[index]~
28 eigenvecs = eigenvecs[:,index] #Sorting eigenvectors according to corresponding eigenvalues~
29 ~
30 reduced = []~
31 for vec in eigenvecs:~
32     reduced.append(vec[0:2]) #Take the first two eigenvectors to reduce the data to two dimensions.~
33 ~
34 reduced = np.array(reduced)~
35 reduced = reduced.transpose()~
36 npfeatures = np.array(features)~
37 ~
38 x = []~
39 y = []~
40 ~
41 for i in range(len(features)):~
42     temp = np.dot(reduced, npfeatures[i].transpose()) #creating the new features~
43     x.append(temp[0])~
44     y.append(temp[1])~
45 ~
46 plt.scatter(x,y,s=1)~
47 ~
48 samples = random.sample(range(0, len(x)), 200)~
49 ~
50 for i in samples:~
51     plt.annotate(labels[i], (x[i],y[i]))~
52 ~
53 plt.show()~

```

Figure 1: My PCA implementation

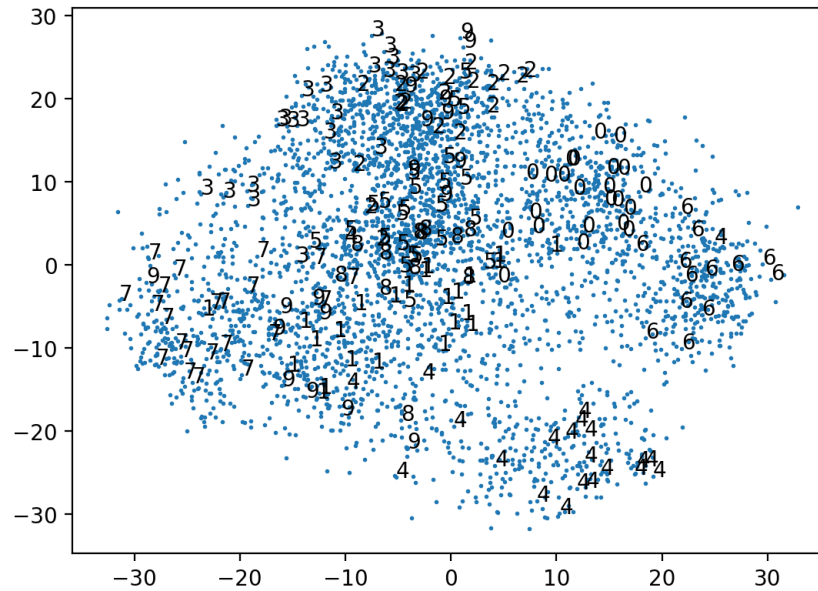


Figure 2: Output of my PCA implementation

2 Question 2

For this implementation, I assumed the term(rank) to be the percentage of the singular values used for the compressed image. For example, image obtained using 20 terms means the largest %20 of the singular values of the original image are used to obtain the output image. The outputs I got are not as clear as the example outputs given in the question because of the difference between the eigenvalue calculations of numpy "eig" function and the built-in svd function. Figure 3 shows the original while Figure 4, Figure 5, Figure 6 and Figure 7 shows the images obtained using 1, 5, 20 and 50 terms respectively.

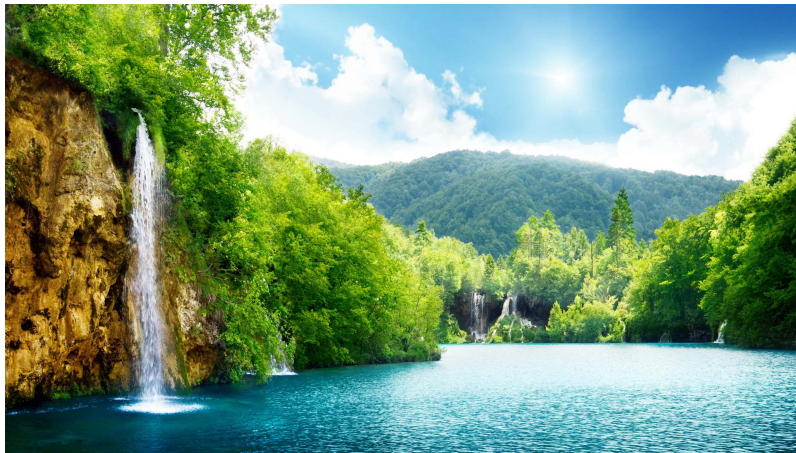


Figure 3: Original image

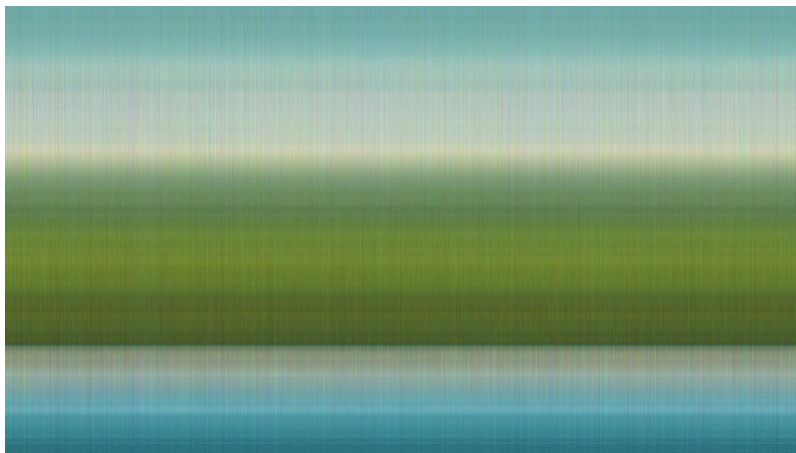


Figure 4: Image obtained using 1 term

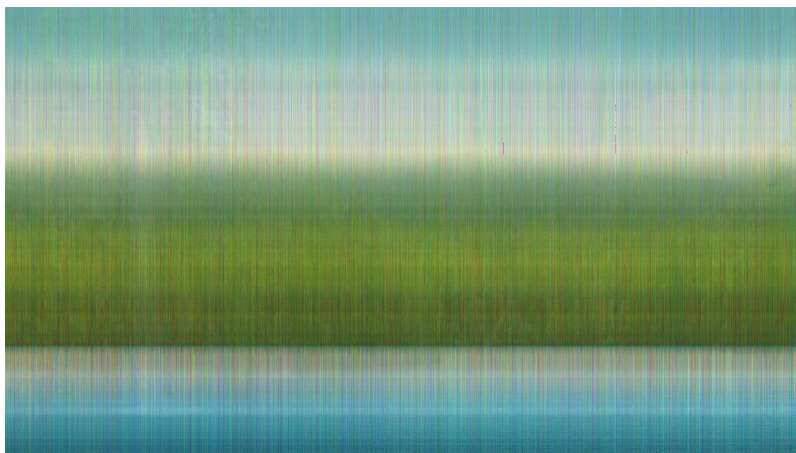


Figure 5: Image obtained using 5 terms

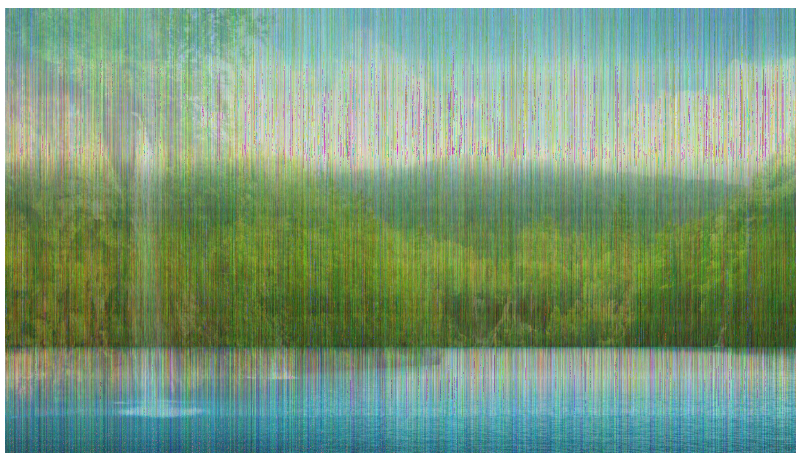


Figure 6: Image obtained using 20 terms

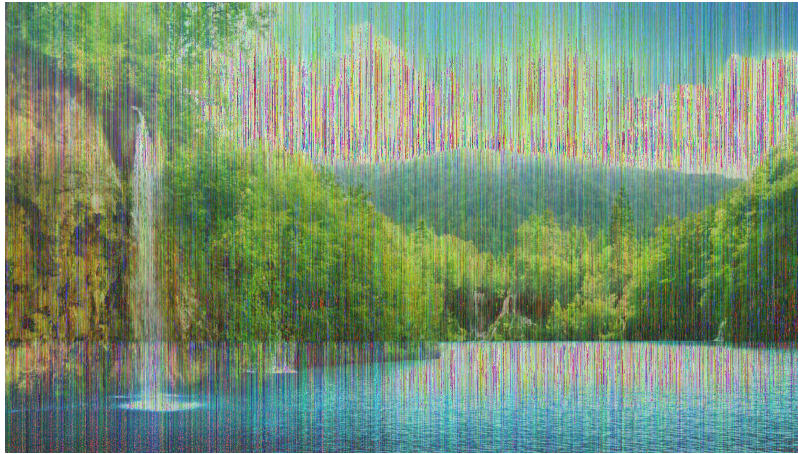


Figure 7: Image obtained using 50 terms