



AMIRKABIR UNIVERSITY OF TECHNOLOGY  
(TEHRAN POLYTECHNIC)  
DEPARTMENT OF COMPUTER ENGINEERING

STATISTICAL PATTERN RECOGNITION

## Assignment IV

*Soroush Mahdi  
99131050*

supervised by  
Dr. Mohammad Rahmati

February 17, 2021



## Contents

<b>1 PCA vs. LDA: Comparison of Two Different Linear Projections</b>	<b>1</b>
1.1 a . . . . .	1
1.2 b . . . . .	1
1.3 c . . . . .	1
1.4 d . . . . .	2
1.5 e . . . . .	2
1.6 f . . . . .	2
<b>2 Understanding the Behavior of Clustering Techniques</b>	<b>3</b>
2.1 a . . . . .	3
2.2 b . . . . .	3
2.3 c . . . . .	4
2.4 d . . . . .	4
2.5 e . . . . .	4
2.6 f . . . . .	4
2.7 g . . . . .	4
<b>3 Why Naples-Style Pizzas Are So Delicious?</b>	<b>5</b>
3.1 a . . . . .	5
3.2 b . . . . .	6
3.3 c . . . . .	6
<b>4 Interpretation of 1000 Genomes Project Results</b>	<b>8</b>
4.1 a . . . . .	8
4.2 b . . . . .	8
4.3 c . . . . .	9
4.4 d . . . . .	10
4.5 e . . . . .	10
4.6 f . . . . .	10
<b>5 How Twitter Reacts to a Crisis?</b>	<b>11</b>
5.1 a . . . . .	11
5.2 b . . . . .	13
<b>6 A Tribute to Stephen Hawking</b>	<b>14</b>
6.1 a . . . . .	14
6.2 b . . . . .	15
6.3 c . . . . .	16
6.4 d . . . . .	17
6.5 e . . . . .	18
6.6 f . . . . .	18
6.7 g . . . . .	19
<b>7 Are You Into Fashion?</b>	<b>19</b>
7.1 a . . . . .	19
7.2 b . . . . .	21
7.3 c . . . . .	21
7.4 d . . . . .	22
7.5 e . . . . .	23
7.6 f . . . . .	24
7.7 g . . . . .	25
7.8 h . . . . .	26



<b>8 Some Explanatory Questions</b>	<b>27</b>
8.1 a . . . . .	27
8.2 b . . . . .	27
8.3 c . . . . .	27
8.4 d . . . . .	28
8.5 e . . . . .	28

# 1 PCA vs. LDA: Comparison of Two Different Linear Projections

## 1.1 a

pca is unsupervised and we know that the variation of values along the line that pca projects points on it should be maximal. so he first principal component will be aligned along the direction of major continuity.

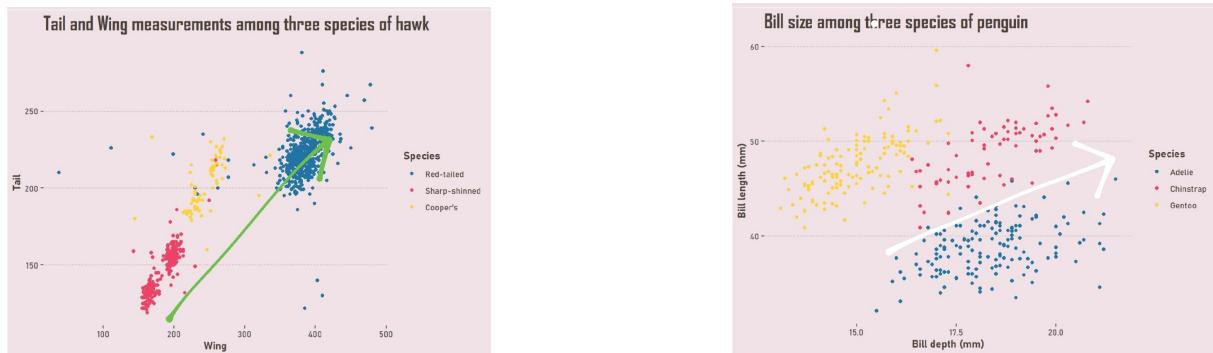


Figure 1: first pca direction for blue and red classes

## 1.2 b

FLD selects a projection that maximizes the class separation and it is supervised.

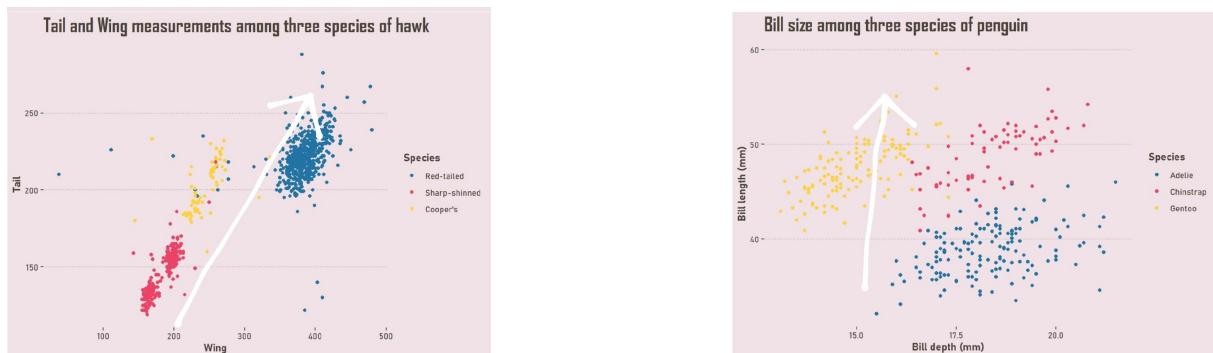


Figure 2: first fld direction for blue and red classes

## 1.3 c



Figure 3: first pca direction for yellow and red classes

**1.4 d**


Figure 4: first fld direction for yellow and red classes

**1.5 e**


Figure 5: first pca direction for all classes

**1.6 f**


Figure 6: first fld direction for all classes

## 2 Understanding the Behavior of Clustering Techniques

### 2.1 a

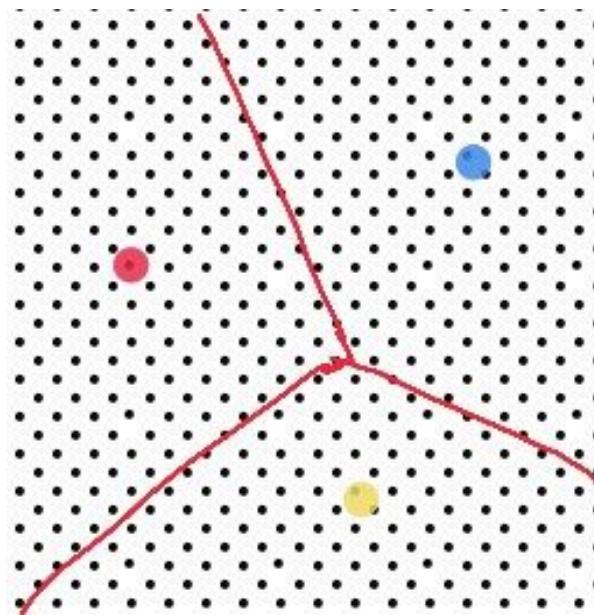


Figure 7: boundaries can rotate based on initial centroids

### 2.2 b

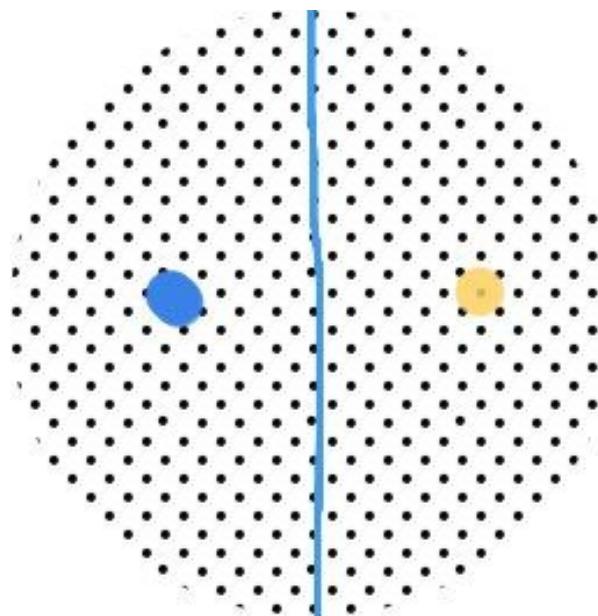


Figure 8: boundaries can rotate based on initial centroids

**2.3 c**

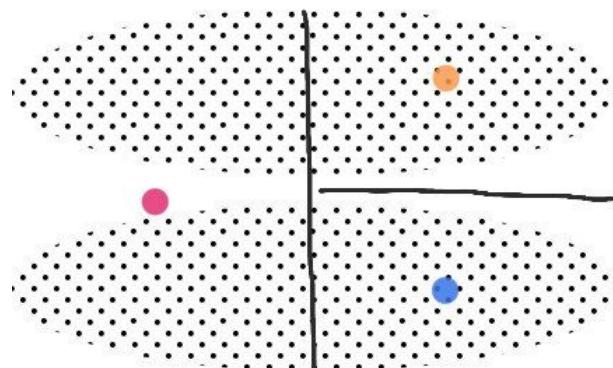


Figure 9: c

**2.4 d**

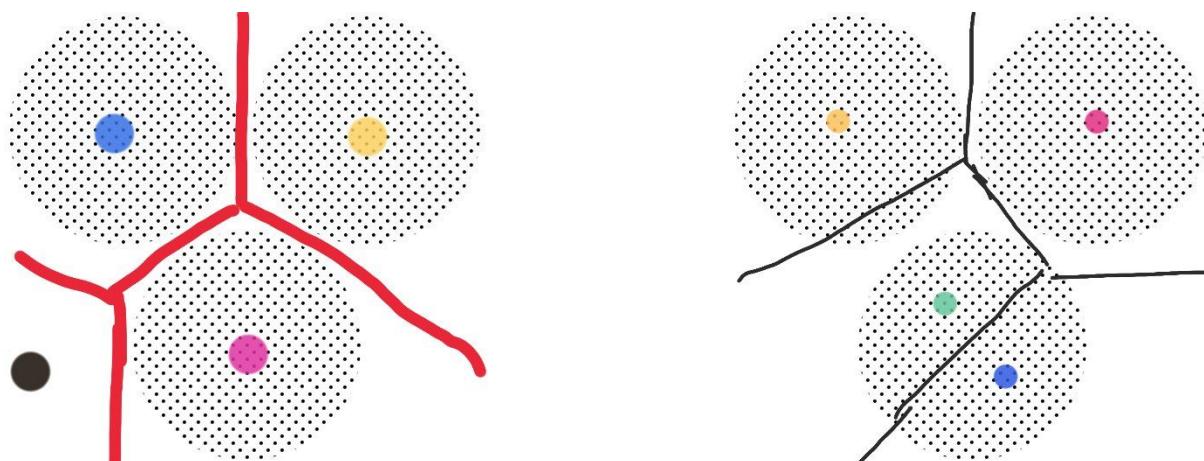


Figure 10: the left figure is global minimum

**2.5 e**



Figure 11: the left figure is global minimum

**2.6 f**

the diagram in the right.

**2.7 g**

it can't find the patterns that are created by empty spaces. in other words spaces that are surrounded by points.



### 3 Why Naples-Style Pizzas Are So Delicious?

#### 3.1 a

```
1 def pca(data , n_component):
2     """performs PCA on a given dataset and return reduced data
3
4     Args:
5         data (m*n numpy array): [m samples that each has n features]
6         n_component ([int]): [num of pca components to keep]
7
8     Returns:
9         [m * n_component numpy array]: [reduced data]
10    """
11    # computing mean of data set
12    mu = data.mean(axis = 0)
13    mu = mu.reshape(data.shape[1] , 1).T
14    #subtracting mean from data samples
15    data = data - mu
16    #computing scatter matrix
17    S = np.cov(data.T)*(data.shape[0]-1)
18    #computing eigen values and eigen vectors of scatter matrix
19    eigen_values , eigen_vecs = np.linalg.eig(S)
20    #selecting eigen vectors corresponding to first biggest n_component
21    ind = eigen_values.argsort()[-n_component:][::-1]
22    E = eigen_vecs[ : , ind ]
23    #computing reduced data
24    reduced_data = E.T @ data.T
25    return reduced_data.T
```

Listing 1: PCA implementation

```
1 #reading data
2 df = pd.read_csv('doughs.dat' , sep = ' ')
3 data = df.drop(['Restaurant'] , axis = 1).to_numpy()
4
5 reduced_data = pca(data , 3)
6
7 print(reduced_data)
```

Listing 2: performing pca on the given dataset

we droped restaurant column from data for performing pca since it's the label of the data

```
1 [[ 13.40370431  11.37083773   0.30190493]
2  [ 10.2126015   4.24606046   0.2822107 ]
3  [-3.64788876  12.49211826   0.33206521]
4  [ 13.89226015  5.2314273   0.45944424]
5  [ 20.14596971  3.38448304   0.76215955]
6  [ 16.21181324 -10.78931755  -0.2870618 ]
7  [  3.84156067   8.448928   0.56334086]
8  [ 15.87973653 -0.1018432  -1.24818735]
9  [  0.69173569  4.33845887  -0.11756393]
10 [  5.99252814 -9.26801967  -1.15644864]
11 [-13.00799405  4.91431839   1.16260885]
12 [ -5.55145301   1.7479005  -2.80424998]
13 [  5.44247276 -5.48477484   0.48696545]
14 [ -7.95719496   0.80218215  -0.04906952]
15 [ -0.67076025  -6.08928916   1.25486264]
16 [ -0.4312163   1.81151802   0.68236867]
17 [ -2.50757979   7.79374866   0.70397545]
18 [ -4.65523528   7.9756312  -0.46644912]
19 [ 23.51914958  -3.26915492  -0.43334213]
20 [ 17.39169082  -5.79685593  -0.16355244]
21 [ -21.99578986   7.83123773  -1.45754971]
22 [ -16.08178889  -9.23748228   1.81834199]
23 [ -7.08086097   1.22832774  -0.5428211 ]
24 [ -5.93294033  -10.87314559 -1.42558025]
25 [ -18.1445075   -8.01392006   0.3301058 ]
26 [ -2.6109593   1.64739583   0.9739266 ]
```



```
27 [-15.64833777 -5.25127303 -0.05930765]
28 [-4.2969673 -8.53637902 0.65459365]
29 [-4.90669916 -11.37183107 0.15695123]
30 [-11.49704962 8.81871243 -0.71464218]
```

Listing 3: first three principal components of the given data

based on the pca values it seems that they are sufficient to describe most of the variation in the dataset.

### 3.2 b

in this section i used scipy library for testing normality of data

```
1 reduced_data = pca(data , 6)
2
3 #testing normality with different PCAs
4 for i in range(6):
5     print('PCA num: ', i+1 , '----', stats.kstest(reduced_data[:, i], 'norm'))
```

Listing 4: testing normality of data with different PCAs

```
1 PCA num: 1 ---- KstestResult(statistic=0.5272552766410683, pvalue=2.6588368553952183e-08)
2 PCA num: 2 ---- KstestResult(statistic=0.45026164389662604, pvalue=4.52986908971393e-06)
3 PCA num: 3 ---- KstestResult(statistic=0.11110902105959519, pvalue=0.8527367641464626)
4 PCA num: 4 ---- KstestResult(statistic=0.20589309188960686, pvalue=0.13612676348380753)
5 PCA num: 5 ---- KstestResult(statistic=0.3456856468354019, pvalue=0.0010470771123993342)
6 PCA num: 6 ---- KstestResult(statistic=0.4252427561605906, pvalue=1.9402352972195527e-05)
```

Listing 5: output for testing normality of data

The D statistic is the absolute max distance (supremum) between the CDFs of the two samples. The closer this number is to 0 the more likely it is that the two samples were drawn from the same distribution. You reject the null hypothesis that the two samples were drawn from the same distribution if the p-value is less than your significance level.we consider this level as 0.05 so if p value is more than 0.05 then we can say the distribution is normal. based on these values we can say that the third and fourth PCAs from data are likely to have normal distribution because they have p values greater than 0.05 and smallest D statistics.

### 3.3 c

```
1 #reading data
2 df = pd.read_csv('doughs.dat' , sep = ' ')
3 data = df.drop(['Restaurant'] , axis = 1).to_numpy()
4
5 #adding a column to dataset to check if sample is from naples or not
6 df['Naples'] = df.apply(lambda x: 0 if x['Restaurant'] in [1,2,3,4] else 1 , axis = 1)
7
8 reduced_data = pca(data , 3)
9 print(reduced_data)
10 #plotting 3d scatter plot for first 3 PCAs
11 fig = plt.figure()
12 ax = fig.add_subplot(111, projection='3d')
13 scatter = ax.scatter(reduced_data[:, 0], reduced_data[:, 1], reduced_data[:, 2],
14                      marker='o' , c=df['Naples'] , cmap=plt.cm.Set1, edgecolor='k')
15 plt.xlabel('first pca')
16 plt.ylabel('sec pca')
17 ax.set_zlabel('third pca')
18 handles = scatter.legend_elements() [0]
19 labels = list(df['Naples'].unique())
```



```
19 legend1 = ax.legend(handles, labels, title="Naples")
20 plt.savefig('3c-3D.png')
21 plt.clf()
22
23 #plotting scatter plots for different pairs of first 3 PCAs
24 for i , j in [[0,1],[0,2],[1,2]]:
25     fig, ax = plt.subplots()
26     scatter = ax.scatter(reduced_data[:, i] ,reduced_data[:, j] ,c=df['Naples'],
27                           cmap=plt.cm.Set1,edgecolor='k')
28     handles = scatter.legend_elements() [0]
29     labels = list(df['Naples'].unique())
30     legend1 = ax.legend(handles, labels,loc="upper left", title="Naples")
31     plt.xlabel('pca num: '+str(i+1))
32     plt.ylabel('pca num: '+str(j+1))
33     ax.add_artist(legend1)
34     plt.savefig('3c-2D' + str(i+1) + str(j+1)+'.png')
35     plt.clf()
```

Listing 6: plotting different scatter plots

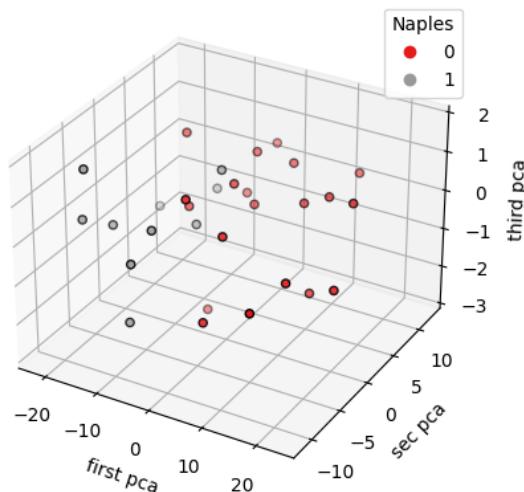


Figure 12: first 3 PCAs

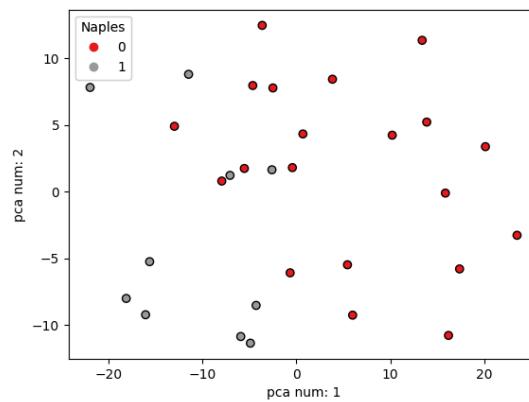


Figure 13: first 2 PCAs

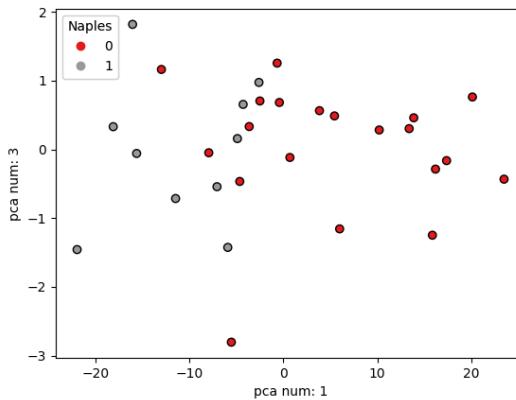


Figure 14: first and third PCAs

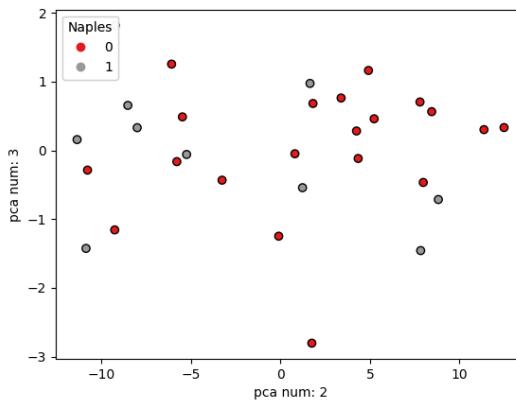


Figure 15: second and third PCAs

in 3D plot we can see that these three PCAs are enough to discriminate between doughs from Naples and doughs from other areas.

## 4 Interpretation of 1000 Genomes Project Results

### 4.1 a

we have 10101 features so we will have 10101 eigen values for scatter matrix and we will choose K of greatest eigen values and eigen vectors corresponding to them which K is less than or equal to 10101 so each vector would have K dimensions.

### 4.2 b

```

1 #reading dataset
2 df = pd.read_csv('1KGP.txt', delim_whitespace=True , header = None)
3 df[2] = df[2].astype("category")
4
5 #computing Y matrix
6 Y = df.drop([0,1,2] , axis = 1)
7 Y = Y.apply(lambda x: x!=x.mode()[0]).to_numpy()
8 Y = Y.astype('int')
9 #changing mean to zero
10 Y = Y - Y.mean(axis = 0).reshape(10101,1).T

```



```
11
12 #performing pca on Y
13 pca = PCA(n_components=2)
14 Y_reduced = pca.fit_transform(Y)
15
16 #plotting data based on 2 first PCAs and population
17 fig, ax = plt.subplots()
18 scatter = ax.scatter(Y_reduced[:, 0], Y_reduced[:, 1], c=df[2].cat.codes, cmap=plt.cm
    .Set1, edgecolor='k')
19 handles = scatter.legend_elements()[0]
20 labels = list(df[2].cat.categories)
21 legend1 = ax.legend(handles, labels, loc="lower right", title="population")
22 plt.xlabel('first pca')
23 plt.ylabel('second pca')
24 ax.add_artist(legend1)
25 plt.savefig('4b.png')
```

Listing 7: plotting data based on 2 first PCAs and population

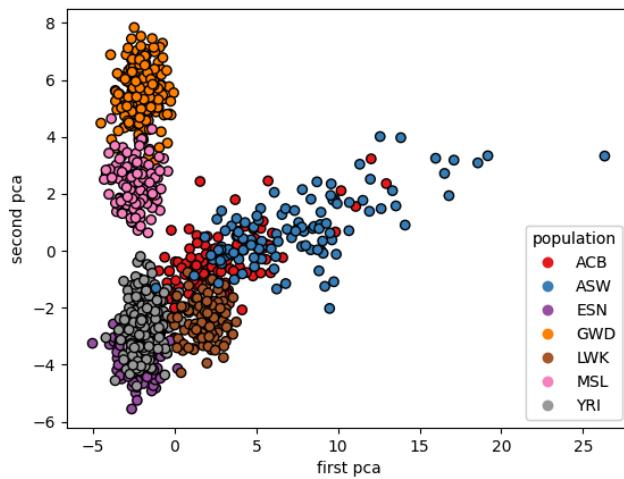


Figure 16: first 2 PCAs for data

### 4.3 c

we can see that these components are able to separate samples based on their locations. the ESN and YRI populations are in the same area because they are both samples from nigeria.

#### 4.4 d

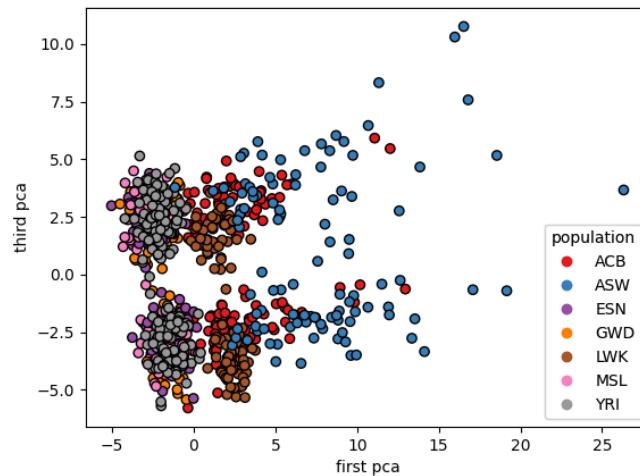


Figure 17: first and third PCAs

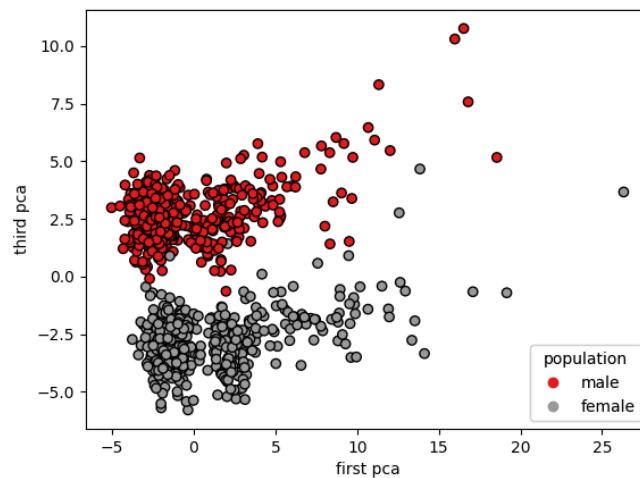


Figure 18: first and third PCAs

we can see that these components can separate clusters based on gender.

#### 4.5 e

based on the plots in last section we can see that third PCA captures the gender features.

#### 4.6 f

```

1 #computing Y matrix
2 Y = df.drop([0,1,2] , axis = 1)
3 Y = Y.apply(lambda x: x!=x.mode()[0]).to_numpy()
4 Y = Y.astype('int')
5 #changing mean to zero
6 Y = Y - Y.mean(axis = 0).reshape(10101,1).T
7

```



```
8 #performing pca on Y
9 pca = PCA(n_components=3)
10 Y_reduced = pca.fit_transform(Y)
11
12 #plotting nucleobase index versus the absolute value of the third principal component
13 fig, ax = plt.subplots()
14 scatter = ax.scatter(range(10101) , np.abs(pca.components_[2]) )
15 plt.xlabel('index')
16 plt.ylabel('abs of third pca')
17 plt.savefig('4f')
```

Listing 8: plotting nucleobase index versus the absolute value of the third principal component

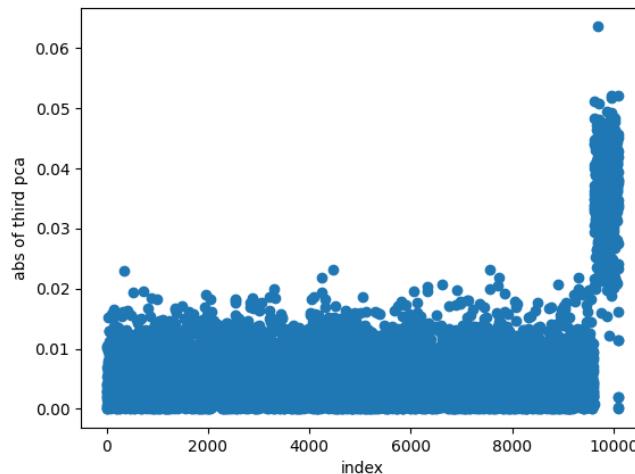


Figure 19: nucleobase index versus the absolute value of the third principal component

in part D we saw that third PCA can separate samples based on genders. in this section we can see third PCA absolute values are greater for last nucleobases. a possible explanation is maybe theses last nucleobase and last chromosomes are different for men and women.

## 5 How Twitter Reacts to a Crisis?

### 5.1 a

some functions that i implemented are these with their descriptions:

```
1 def jaccard_dist(a , b):
2     """calculate jaccard dist between 2 sets
3
4     Args:
5         a ([set])
6         b ([set])
7
8     Returns:
9         [int]: [jaccard dist between set a and set b]
10    """
11    jaccard_index = len(a.intersection(b)) / len(a.union(b))
12    return 1 - jaccard_index
13
14 def assign_cluster(data , centroids):
15     """for point data calculate which centroid it belongs to
16
17     Args:
18         data ([set]): [set of words of a tweet]
19         centroids ([dict]): [centroids id's and their set of words]
20
```



```
21     Returns:
22         [int]: [it's the id of centroid which data belongs to]
23     """
24     min_dist = 1
25     for key in centroids.keys():
26         dist = jaccard_dist(data , centroids[key])
27         if dist < min_dist:
28             min_dist = dist
29             assign = key
30     return assign
31
32 def df_to_dict(df):
33     """ a function for changing data to a dict with this format:
34         every key is an id of a tweet and for every key we have
35         a set which it has unique words of it's tweet
36
37     Args:
38         df ([dataframe]): [ids and tweets]
39
40     Returns:
41         [dict]
42     """
43     df.set_index("id", drop=True, inplace=True)
44     df = df.to_dict(orient="index")
45     for i in df.keys():
46         df[i] = set(df[i]['text'].split())
47     return df
48
49 def new_centroid(cluster):
50     """for a cluster calculates it's centriod
51
52     Args:
53         cluster ([dict]): [it's a dict of id's as keys and their set of words of
54                             corresponding tweet
55                                         which these id's belongs to a single cluster]
56
57     Returns:
58         [int]: [new centroid id]
59     """
60     min_dist = 10**4
61     for i in cluster.keys():
62         dist = 0
63         for j in cluster.keys():
64             dist+=jaccard_dist(cluster[i] , cluster[j])
65         if dist < min_dist:
66             min_dist = dist
67             centroid = i
68     return centroid
```

Listing 9: usefull functions

```
1 #reading data and centroids
2 df = pd.read_json('tweets.json', lines=True)
3 data = df[['text' , 'id']]
4 data = df_to_dict(data)
5
6 centroids = pd.read_csv('initial_centroids.txt' , header=None)[0]
7 centroids = list(centroids)
```

Listing 10: reading data

and clustering steps are these:

```
1 converge = False
2 while(not converge):
3     result = {}
4     #assign each tweet to it's nearest centroid and saving them in result
5     for i in data.keys():
6         cluster = assign_cluster(data[i] , {centroid:data[centroid] for centroid in
6                                     centroids})
7
```



## Assignment IV

```
8         if cluster not in result.keys():
9             result[cluster] = []
10
11         result[cluster].append(i)
12
13     centroids_new = []
14
15     #calculating new centroid for each cluster
16     for i in result.keys():
17         cluster_ids = result[i]
18         cluster_ids.append(i)
19
20         centroids_new.append(new_centroid({j:data[j] for j in cluster_ids}))
21
22     #checking converge
23     if set(centroids_new) == set(centroids):
24         converge = True
25     centroids = centroids_new
```

Listing 11: clustrings

and in this section results with wanted format will save on a file with name result-a.txt

```
1 with open("result_a.txt", "w") as file:
2     for i in result.keys():
3         file.write(str(i) + ', : ')
4         for j in result[i]:
5             file.write(str(j) + ', ')
6         file.write('\n')
```

Listing 12: saving result in a file

Figure 20: result

5.2 b

in this section i used kmeans++ idea for initialization seeds , and the function for this work and it's description is shown here:

```
1 def initial_seeds(data , k):
2     """this function works based on initialization of kmeans++
3         The first seed is randomly selected from the set of tweets.
4         Then , for each tweet t , the distance between t and the nearest seed that has
5             already
6                 been chosen is computed. Another seed is then chosen with probability of
7                     its distance squared among the sum of all distances squared.
8                         These steps are computed until k seeds have been selected.
9
10    Args:
11        data ([dict]): [ a dict with tweet id's as keys and their words as values]
12        k ([int]): [number of centroids that we want]
```



```
12     Returns:  
13         [list]: [id's of tweets that the function choose to be centroid]  
14     """  
15     centroids = []  
16     centroids.append(np.random.choice(list(data.keys())))  
17  
18     while(len(centroids) < k):  
19         dists_from_near_seed = {}  
20         probs = {}  
21         for key in data.keys():  
22             dists_from_near_seed[key] = min([jaccard_dist(data[key], data[centroid])  
23                 for centroid in centroids])  
24         for key in data.keys():  
25             probs[key] = dists_from_near_seed[key]**2 / (sum([dists_from_near_seed[i]  
26                 **2 for i in dists_from_near_seed.keys()]))  
27  
28     centroids.append(np.random.choice(a = list(probs.keys()), p = [probs[key]  
for key in probs.keys()]))  
29     return centroids
```

Listing 13: initilizaliza

the result of this section is saved in a file name result-b.txt

## 6 A Tribute to Stephen Hawking

### 6.1 a

```
1 #reading data  
2 data = pd.read_csv('ALS_train.csv')  
3 data = data.set_index('ID')  
4  
5 #covariance matrix of data  
6 cov = data.corr()  
7 #selecting 10 columns with biggest value of covariance with ALSFRS_slope  
8 cov = cov['ALSFRS_slope'].abs()  
9 cov = cov.sort_values()  
10 columns = cov[-11:].index  
11 data = data[columns]
```

Listing 14: selecting 10 features

```
1 #plotting heatmap of cov matrix  
2 sns_plot = sns.heatmap(data.corr())  
3 sns_plot.get_figure().savefig("6a1.png")  
4  
5 #plotting data in ALSFRS_slope vs other featurs spaces  
6 fig, axs = plt.subplots(2, 5, figsize=(12, 5))  
7 temp = 0  
8 for i in range(2):  
9     for j in range(5):  
10         axs[i,j].scatter(data['ALSFRS_slope'], data[columns[temp]], marker='.')  
11         axs[i,j].set_title(columns[temp])  
12     temp+=1  
13  
14 for j in range(5,10):  
15     axs.flat[j].set(xlabel='ALSFRS_slope')  
16  
17 for ax in fig.get_axes():  
18     ax.label_outer()  
19  
20 plt.savefig('6a2')
```

Listing 15: visualization

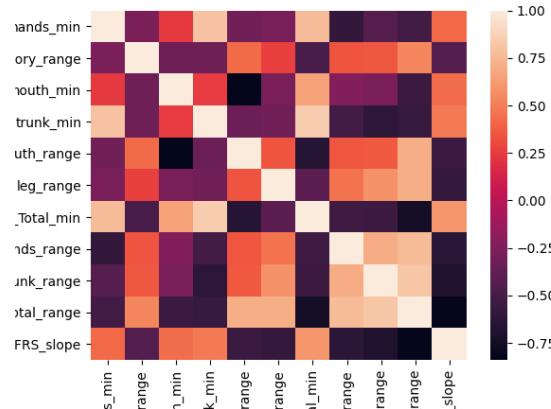


Figure 21: heatmap for covariance matrix of new data, the last column is for ALSFRS\_slope column

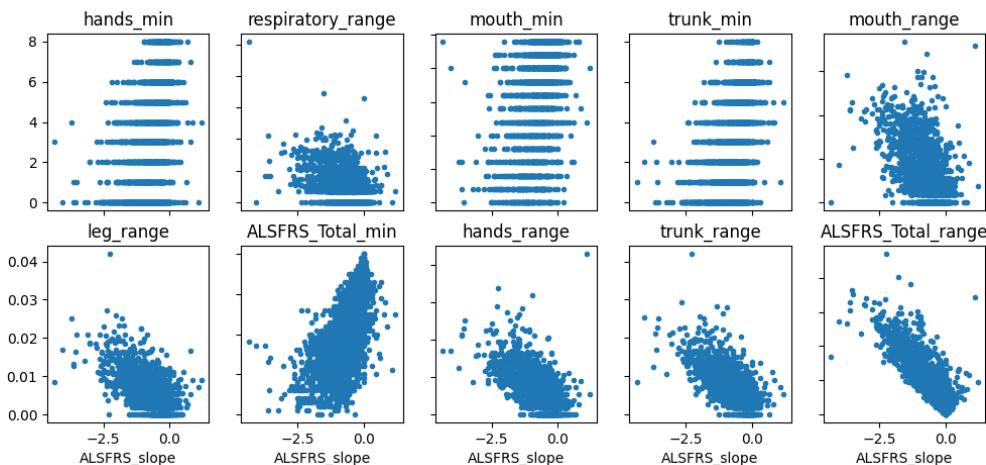


Figure 22: plotting data in different spaces with ALSFRS\_slope as X axis and other features as Y axis

## 6.2 b

for selecting K i used elbow method with SSE criteria and silhouette method.

```

1 train_data = data.drop('ALSFRS_slope', axis = 1)
2
3 #elbow method
4 sse = {}
5 for k in range(1, 15):
6     kmeans = KMeans(n_clusters=k, init = 'random').fit(train_data)
7     sse[k] = kmeans.inertia_ # Inertia: Sum of distances of samples to their closest
    cluster center
8 plt.figure()
9 plt.plot(list(sse.keys()), list(sse.values()))
10 plt.xlabel("Number of cluster")
11 plt.ylabel("SSE")

```

Listing 16: elbow method

```

1 sil = []
2 # dissimilarity would not be defined for a single cluster, thus, minimum number of
    clusters should be 2
3 for k in range(2, 15):
4     kmeans = KMeans(n_clusters = k, init = 'random').fit(train_data)

```



```
5     sil.append(silhouette_score(train_data, kmeans.labels_, metric = 'euclidean'))
6
7 plt.plot(range(2,15) , sil)
8 plt.xlabel("Number of cluster")
9 plt.ylabel("silhouette score")
```

Listing 17: silhouette method

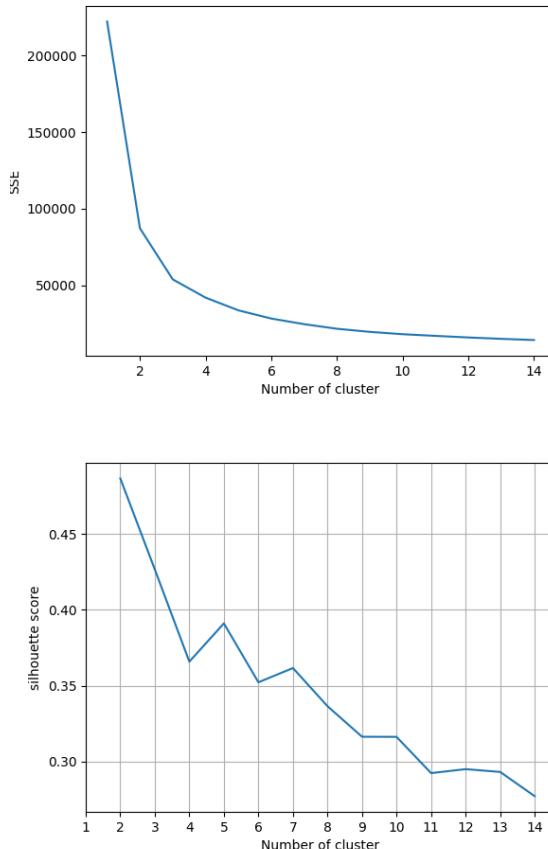


Figure 23: elbow and silhouette method plots for selecting K

based on elbow method it seems in  $K = 2$  and  $K = 3$  we have elbow. based on silhouette method the optimal  $K$  is 2.. i choose 3 for  $K$ .

### 6.3 c

```
1 #reading data
2 data = pd.read_csv('ALS_train.csv')
3 data = data.set_index('ID')
4
5 #covariance matrix of data
6 cov = data.corr()
7 #selecting 10 columns with biggest value of covariance with ALSFRS_slope
8 cov = cov['ALSFRS_slope'].abs()
9 cov = cov.sort_values()
10 columns = cov[-11:].index
11 data = data[columns]
12 train_data = data.drop('ALSFRS_slope' , axis = 1)
13
14 kmeans = KMeans(n_clusters = 3,init = 'random' ).fit(train_data)
15
16 print('centers of clusters:')
```



```
17 print(kmeans.cluster_centers_)
18 print()
19 print('silhouette score: ', silhouette_score(train_data, kmeans.labels_, metric =
    'euclidean'))
20
21 #calculating number of samples in each cluster
22 cluster1_num = np.count_nonzero(kmeans.labels_ == 0)
23 cluster2_num = np.count_nonzero(kmeans.labels_ == 1)
24 cluster3_num = np.count_nonzero(kmeans.labels_ == 2)
25
26 #plotting bar chart
27 fig, ax = plt.subplots()
28 plt.bar([0,1 , 2], [cluster1_num , cluster2_num , cluster3_num] , width = 0.3)
29 plt.xticks([0,1,2], ('cluster 0' , 'cluster 1' , 'cluster 3'))
30 plt.ylabel('number of samples')
31 plt.savefig('6c')
```

Listing 18: Evaluating model performance

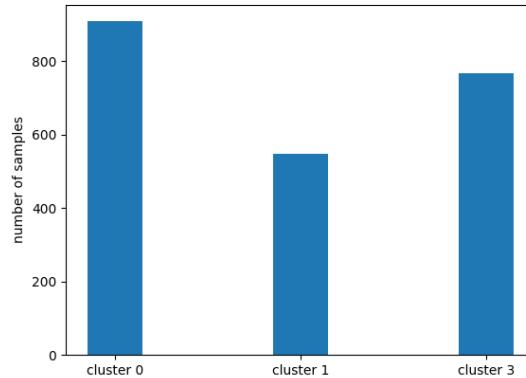


Figure 24: number of samples in each cluster

centers of clusters: [[2.32857143e+00 2.57275631e-03 8.22967033e+00 2.31868132e+00 5.85175584e-03  
6.49584660e-03 1.86230769e+01 7.92365925e-03 7.89464384e-03 2.71946639e-02] [7.57299270e-01 4.14316504e-  
03 3.71532847e+00 6.66058394e-01 1.32495969e-02 8.34827940e-03 8.69160584e+00 9.64347110e-03 1.00976663e-  
02 4.23871766e-02] [5.54248366e+00 1.27555133e-03 1.01516340e+01 5.35424837e+00 2.71081487e-03 4.20080686e-  
03 2.93816993e+01 3.66684084e-03 4.11118900e-03 1.29428526e-02]]  
silhouette score: 0.4266824551341931

## 6.4 d

```
1 kmeans = KMeans(n_clusters = 3,init = 'k-means++').fit(train_data)
2
3 #calculating number of samples in each cluster
4 cluster1_num = np.count_nonzero(kmeans.labels_ == 0)
5 cluster2_num = np.count_nonzero(kmeans.labels_ == 1)
6 cluster3_num = np.count_nonzero(kmeans.labels_ == 2)
7
8 #plotting bar chart
9 fig, ax = plt.subplots()
10 plt.bar([0,1 , 2], [cluster1_num , cluster2_num , cluster3_num] , width = 0.3)
11 plt.xticks([0,1,2], ('cluster 0' , 'cluster 1' , 'cluster 3'))
12 plt.ylabel('number of samples')
13 plt.savefig('6d')
```

Listing 19: kmeans++

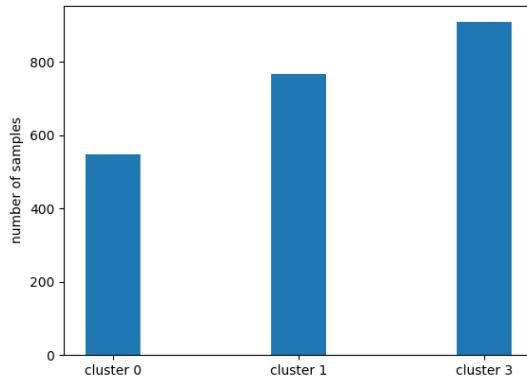


Figure 25: number of samples in each cluster for kmeans++

## 6.5 e

```
1 KMeans(n_clusters = 3, init = 'k-means++')
```

Listing 20: model with optimal parameters

## 6.6 f

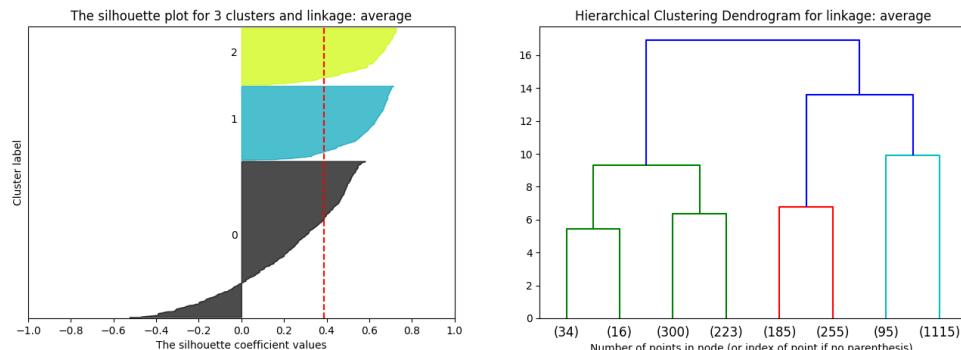


Figure 26: silhouette plot and dendrogram for average linkage

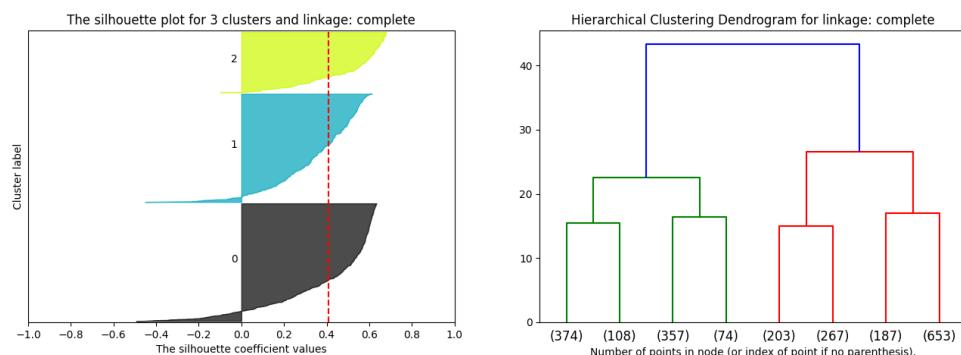


Figure 27: silhouette plot and dendrogram for complete linkage

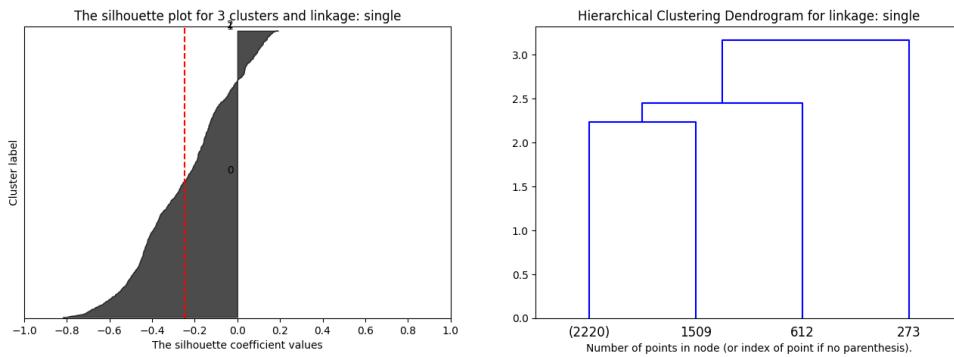


Figure 28: silhouette plot and dendrogram for single linkage

## 6.7 g

we can see that single linkage does not work well comparing to complete and average linkage.

## 7 Are You Into Fashion?

### 7.1 a

```

1 #reading data
2 images = idx2numpy.convert_from_file('t10k-images-idx3-ubyte')
3 labels = idx2numpy.convert_from_file('t10k-labels-idx1-ubyte')
4 labels = labels.reshape(10000,1)
5 images = images.reshape(10000,28*28)
6
7 #performing PCA
8 pca = PCA(n_components=20)
9 pca_reduced = pca.fit_transform(images)
10
11 print('eigenvalues: ', pca.explained_variance_)

```

Listing 21: performing PCA

top 20 eigenvalues: eigenvalues: [1288319.52477778 779197.62253773 265730.43854769 218669.76933454  
169257.23458071 152452.76424582 104674.41864859 83982.28146169 58343.40691272 57195.68410513 43687.9739565  
40723.99085159 33555.72005452 28600.95680321 27416.98190461 25850.84607944 24624.94870779 23563.26086334  
21252.81703952 19652.47616396]

```

1 fig, ax = plt.subplots()
2 scatter = ax.scatter(pca_reduced[:, 0], pca_reduced[:, 1], c=labels, cmap=plt.cm.tab10, edgecolor='k', s = 30)
3 handles = scatter.legend_elements() [0]
4 labels_unique = np.unique(labels)
5 legend1 = ax.legend(handles, labels_unique, loc="upper right", title="category",
6 bbox_to_anchor=(1.15, 1))
7 plt.xlabel('first pca')
8 plt.ylabel('second pca')
9 ax.add_artist(legend1)
10 plt.savefig('a1.png')
11 plt.clf()
12
13 fig = plt.figure()
14 ax = fig.add_subplot(111, projection='3d')
15 ax.scatter(pca_reduced[:, 0], pca_reduced[:, 1], pca_reduced[:, 2], marker='o', c=labels, cmap=plt.cm.tab10, edgecolor='k', s = 20)
16 handles = scatter.legend_elements() [0]
17 labels_unique = np.unique(labels)
18 legend1 = ax.legend(handles, labels_unique, loc="upper right", title="category",
19 bbox_to_anchor=(1.35, 1))
20 plt.xlabel('first pca')
21 plt.ylabel('second pca')

```



```
20 ax.set_zlabel('third pca')
21 ax.add_artist(legend1)
22 plt.savefig('a2.png')
23 plt.clf()
```

Listing 22: plotting data on first 2 and 3 PCs

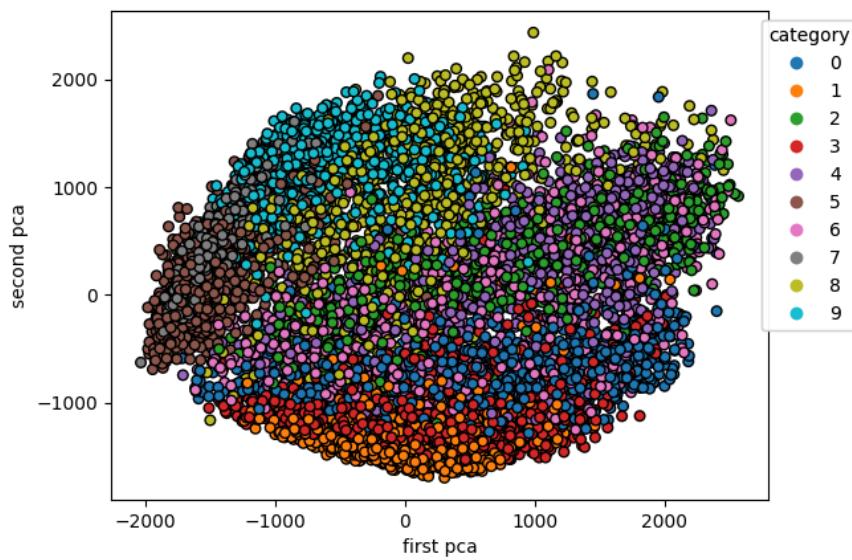


Figure 29: data projection on first 2 PCs

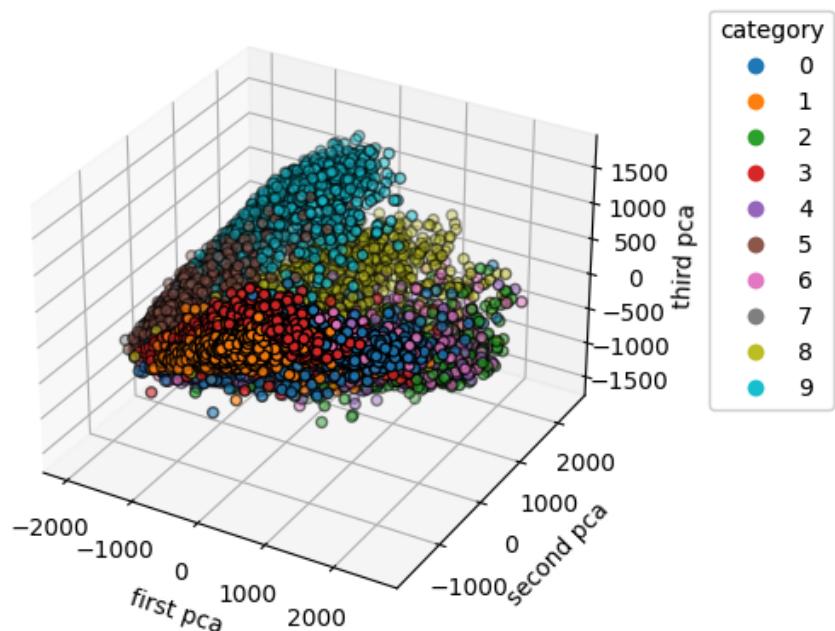


Figure 30: data projection on first 2 PCs



## 7.2 b

```
1 #performing LDA
2 lda = LinearDiscriminantAnalysis(n_components=2)
3 lda_reduced = lda.fit(images, np.ravel(labels)).transform(images)
```

Listing 23: applying LDA on data

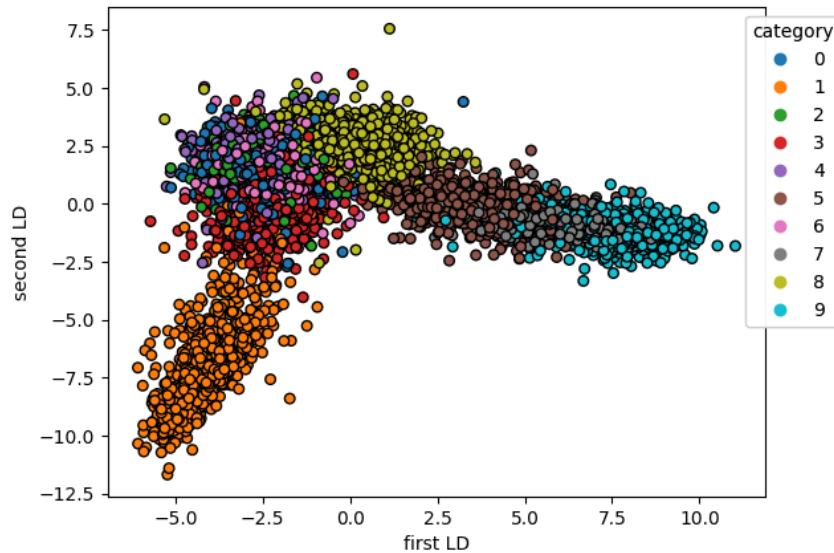


Figure 31: data projection on first 2 LDs

## 7.3 c

```
1 #performing kmeans and plotting clusters
2 for k in [4,7,10]:
3     kmeans = KMeans(n_clusters=k, init = 'random').fit(pca_reduced)
4     fig, ax = plt.subplots()
5     scatter = ax.scatter(pca_reduced[:, 0], pca_reduced[:, 1] ,c=kmeans.labels_ , cmap
6     =plt.cm.tab10, edgecolor='k' , s = 30)
7     handles = scatter.legend_elements() [0]
8     plt.xlabel('first pca')
9     plt.ylabel('second pca')
10    plt.savefig('7c,K= '+str(k))
11    plt.clf()
```

Listing 24: performing kmeans

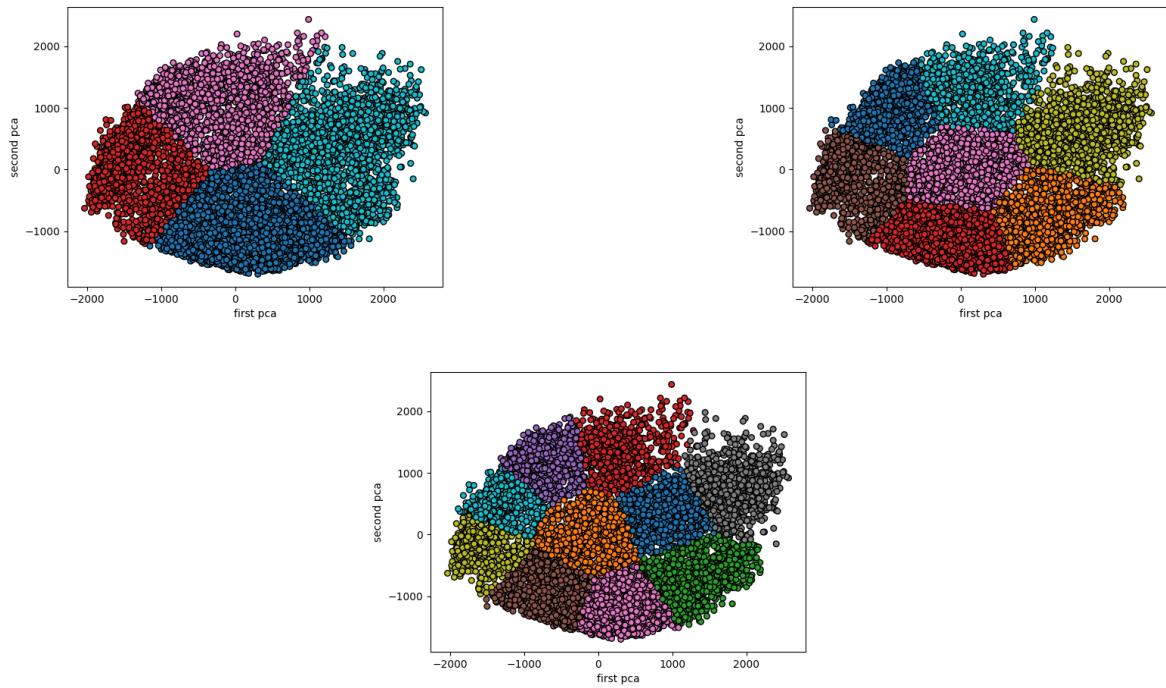


Figure 32: performing kmeans on data using first 2 PCs with different Ks

kmeans try to cluster data with linear boundaries so based on part a plot it cant cluster data using these 2 features in a good way because classes are not linearly separable in this feature space.

#### 7.4 d

```

1 sets = [k_4sets , k_7sets , k_10sets]
2 #performing kmeans and plotting clusters
3 #a loop for different values of k
4 for set in sets:
5     k = len(set)
6     #computing mean of each set
7     means = [pca_reduced[np.isin(labels , i)].mean(axis=0) for i in set]
8     means = np.array(means)
9     #performing kmeans
10    kmeans = KMeans(n_clusters=k, init = means).fit(pca_reduced)

```

Listing 25: computing means and performing kmeans for different Ks

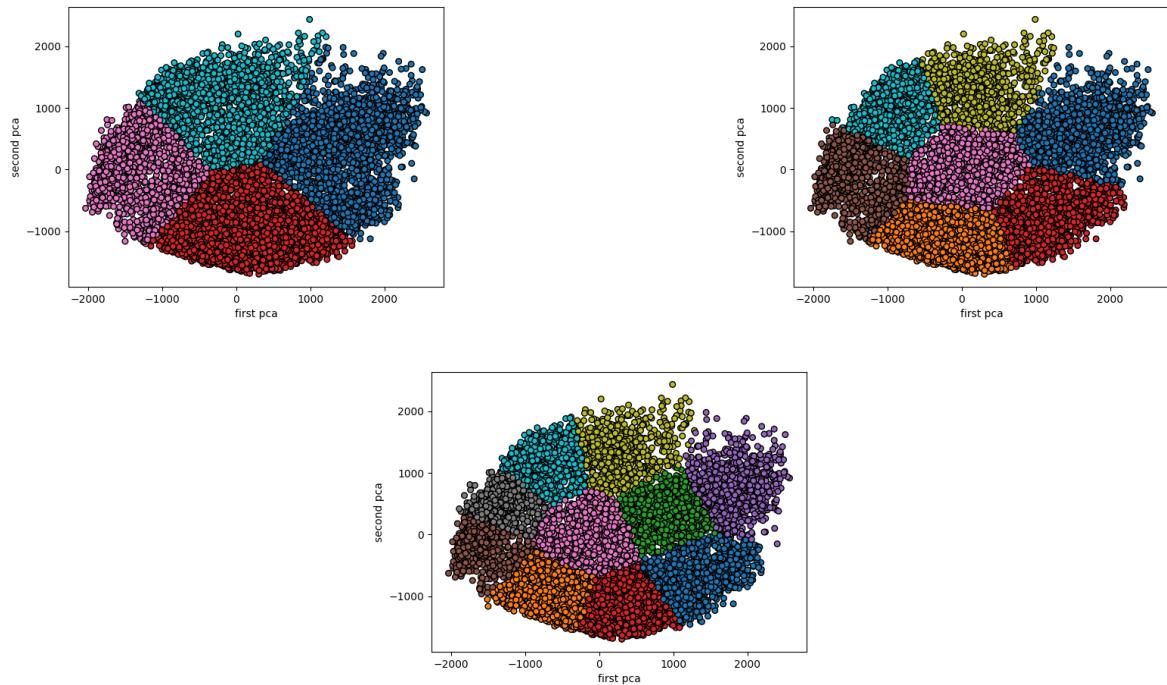


Figure 33: performing kmeans on data using first 2 PCs with different Ks and means as centroids

we can see with centroids as means the result dose not change from random centroids very much

## 7.5 e

185 first PCs are enough

```

1 #performing PCA
2 pca = PCA(n_components=200)
3 pca_reduced = pca.fit_transform(images)
4
5 #calculating min number of PCs for 95% information
6 var = pca.explained_variance_ratio_[0]
7 i = 1
8 while(var < .95):
9     var += pca.explained_variance_ratio_[i]
10    i+=1
11
12 print('enough numbers of PCs for capturing 95% of information: ',i+1)

```

Listing 26: finding enough number of PCs for capturing 0.95 information from data

```

1 #performing pca with i component
2 pca = PCA(n_components=i+1)
3 pca_reduced = pca.fit_transform(images)
4 #reversing reduced data to images
5 inverse_images = pca.inverse_transform(pca_reduced)
6
7 #saving 3 random image and their reduced versions
8 fig, axs = plt.subplots(3, 2 , figsize=(9, 9))
9 for i in range(3):
10     #select 10 indexes randomly from cluster i
11     random_idx = np.random.choice(range(10000))
12     axs[i,0].imshow(images[random_idx].reshape(28,28), cmap='gray', vmin=0, vmax =255 )
13     axs[i,1].imshow(inverse_images[random_idx].reshape(28,28), cmap='gray', vmin=0, vmax=255 )
14 for j in range(3):
15     axs[j,0].set( xlabel='original images')

```



```
16     axs[j,1].set( xlabel='reduced images')
17 for ax in fig.get_axes():
18     ax.label_outer()
19 plt.savefig('7e.png')
```

Listing 27: performing inverse transform on reduced data and showing 3 random pics

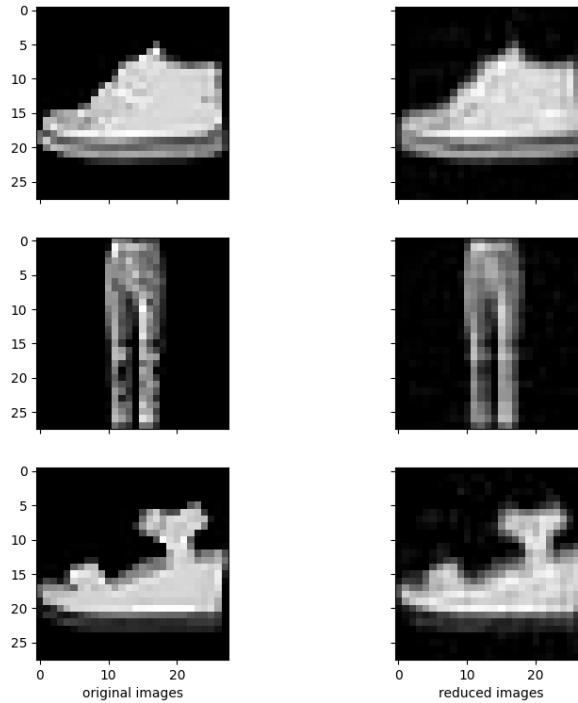


Figure 34: 3 random photos and their reduced versions

## 7.6 f

```
1 #performing PCA
2 pca = PCA(n_components=185)
3 pca_reduced = pca.fit_transform(images)
4 #performing kmeans
5 kmeans = KMeans(n_clusters=10, init = 'k-means++').fit(images)
6
7 #plotting 10 samples from each cluster randomly
8 fig, axs = plt.subplots(10, 10 , figsize=(14, 14))
9
10 for i in range(10):
11     #select 10 indexes randomly from cluster i
12     random_idx = np.random.choice(np.arange(10000)[kmeans.labels_==i], 10)
13
14     for j in range(10):
15         axs[i,j].imshow(images[random_idx[j]].reshape(28,28), cmap='gray', vmin=0,
16                         vmax=255 )
17     for j in range(0,100,10):
18         axs.flat[j].set( ylabel='C'+str(j//10))
19
20 for ax in fig.get_axes():
21     ax.label_outer()
```

```
23 plt.savefig('7f.png')
```

Listing 28: performing kmeans on first 185 PCs

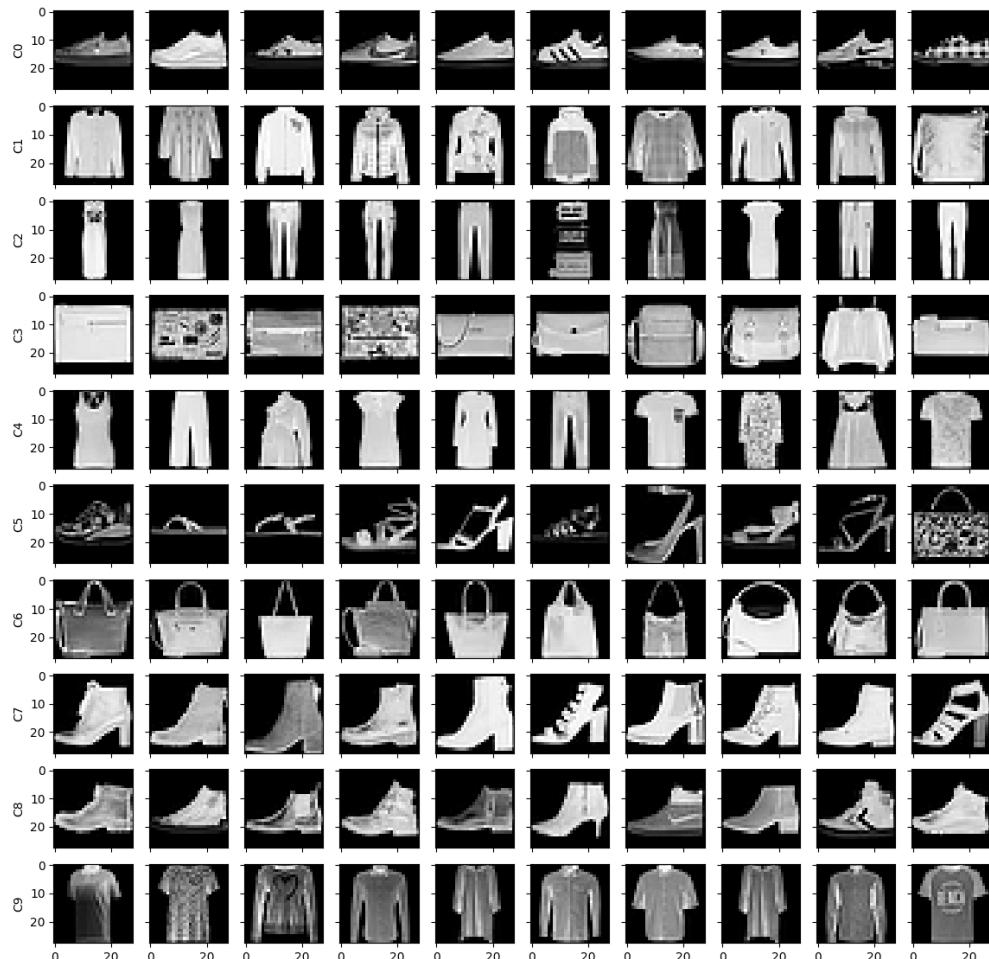


Figure 35: 10 random pics from each cluster

each row represent a cluster. we can see that kmeans on first 185 PCs can divide data in a meaningful way based on fashion categories and can cluster data with a good accuracy based on classes.

## 7.7 g

```
1 #performing PCA
2 pca = PCA(n_components=185)
3 pca_reduced = pca.fit_transform(images)
4
5 #performing kmeans
6 kmeans = KMeans(n_clusters=10, init = 'k-means++').fit(pca_reduced)
7
8 #caculating a dataframe that each row represent a cluster and each column
```



```
9 #reperepresent a class and each cell tell us the percentage of class in cluster
10 percentage = {}
11 for i in range(10):
12     percentage['cluster' + str(i)] = [int(np.count_nonzero(labels[kmeans.labels_ == i]
13 ==k])*100 / np.count_nonzero(kmeans.labels_ == i)) for k in range(10)]
13 df = pd.DataFrame.from_dict(percentage, orient='index')
14
15 #plotting bar plot
16 ax = df.plot(kind='bar', figsize=(20,5), width=0.8, edgecolor=None)
17
18 for p in ax.patches:
19     width = p.get_width()
20     height = p.get_height()
21     x, y = p.get_xy()
22     if height >0:
23         ax.annotate(f'{height}', (x + width/2, y + height*1.02), ha='center')
24
25 plt.legend(labels=df.columns, title = 'classes')
26 plt.savefig('7g')
```

Listing 29: drawing bar grapha

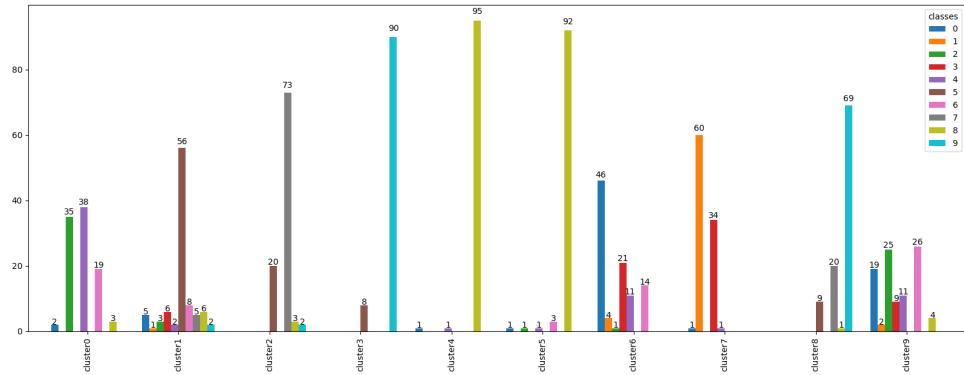


Figure 36: representing the distribution of the samples of different classes in each cluster  
i have rounded percentages so they are not exact.we can see cluster 4 and 5 are like each other.and  
clusters 3 and 8 are like each other.

## 7.8 h

in this section i selected first 3 PCs as features for plotting.

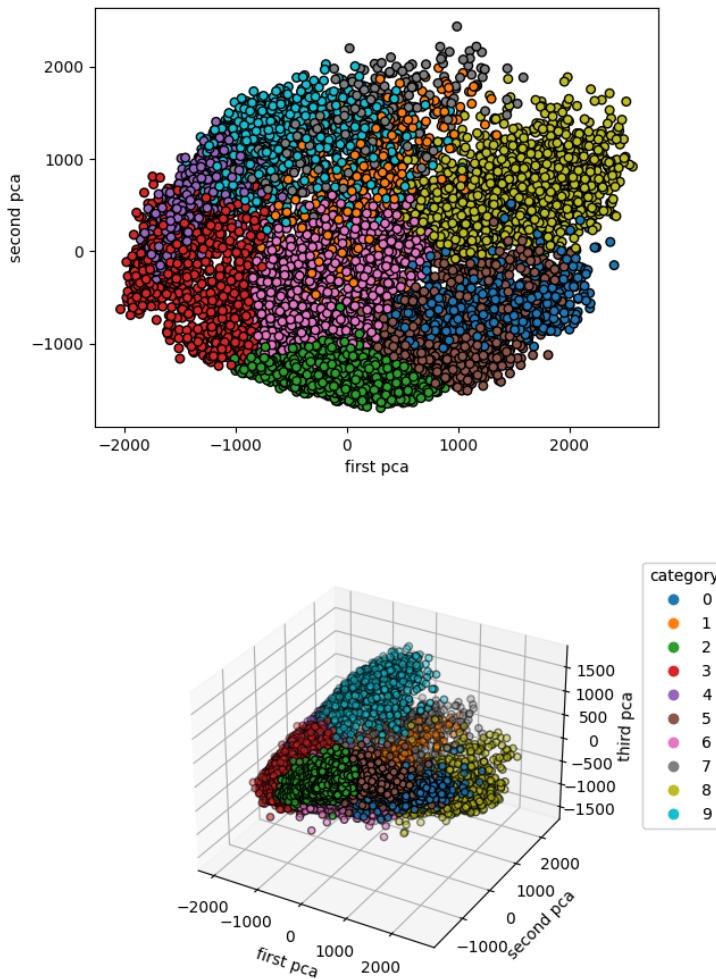


Figure 37: clusters represented on 2d and 3d space based on first 3 PCs for K = 10

## 8 Some Explanatory Questions

### 8.1 a

because it's based on eigen vectors of scatter matrix of data which has  $n * n$  dimension. and each eigen vector give us one of the PCs.so maximum number of PCs that we can have is n because it's the maximum number of eigen vectors of scatter matrix.

### 8.2 b

an image is a matrix of pixels represented by RGB color values. Thus, principal component analysis can be used to reduce the dimensions of the matrix (image) and project those new dimensions to reform the image that retains its qualities but is smaller.

### 8.3 c

yes, basic idea is to replace the euclidean distance computations by kernelized version.

$$\|\Phi(x_n) - \Phi(\mu_k)\|^2 = \|\Phi(x_n)\|^2 + \|\Phi(\mu_n)\|^2 - 2\Phi(x_n)^T \Phi(\mu_k) = K(x_n, x_n) + K(\mu_n, \mu_n) - 2K(x_n, \mu_n) \quad (1)$$



which  $K$  is kernel function and  $\Phi$  is it's feature map. the advantage of this method is it works for nonlinear samples too.

#### 8.4 d

it will decrease. when we increase  $K$ , points that their cluster changes will have smaller error because their new centroids are closer to them, otherwise the cluster which they belongs to would not change. so we would have smaller error and smaller variance.

#### 8.5 e

when  $K$  is equal to  $n$  this happens where  $n$  is the number of samples, because each sample will become a centroid and each cluster is just a sample which the sample itself is centroid so the variance of each cluster becomes zero.