

ECE657A ASSIGNMENT № 2

He Bing, 20848700 , b29he@uwaterloo.ca

Feb 2020

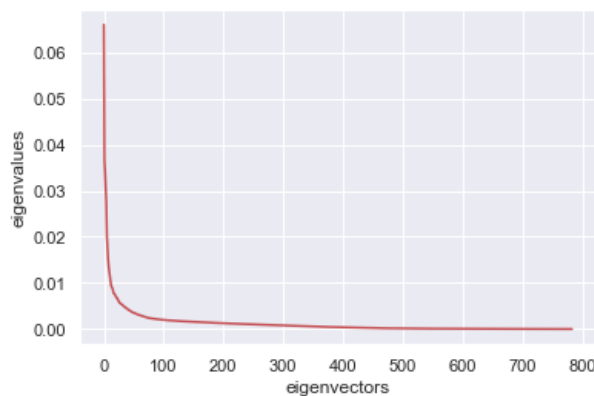
PLEASE SEE THE **JUPYTER NOTEBOOK PDF**
AND **CODE** IN THE END OF THE REPORT!!!!

Problem 2

Principal Component Analysis (PCA)

Plot the scree plot and visually discuss which cut-off is good

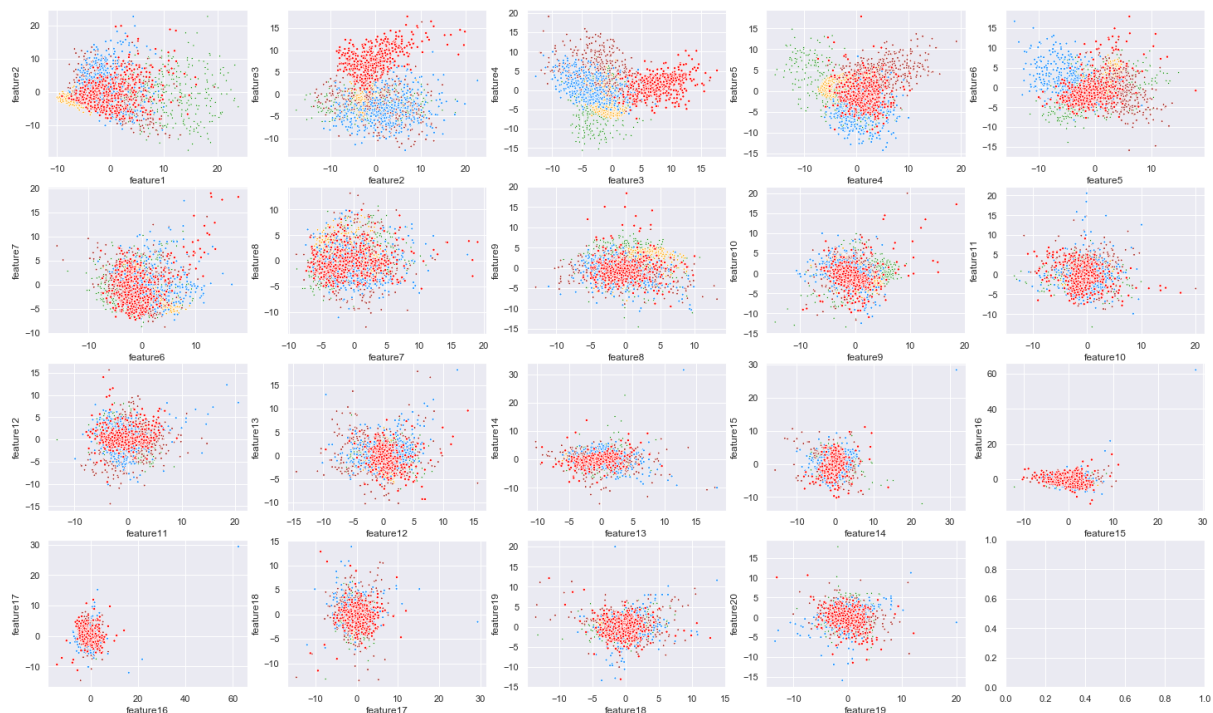
Figure 1: Scree plot of PCA



The cut-off can be around at 10-20. We can see the eigenvalues drop sharply when the index of eigenvector grows. Generally speaking, the first K components take up most part of all the features under PCA. And the eigenvector that is after 50's eigenvalue is very near to 0. So actually, they contribute a very little part to the classification.

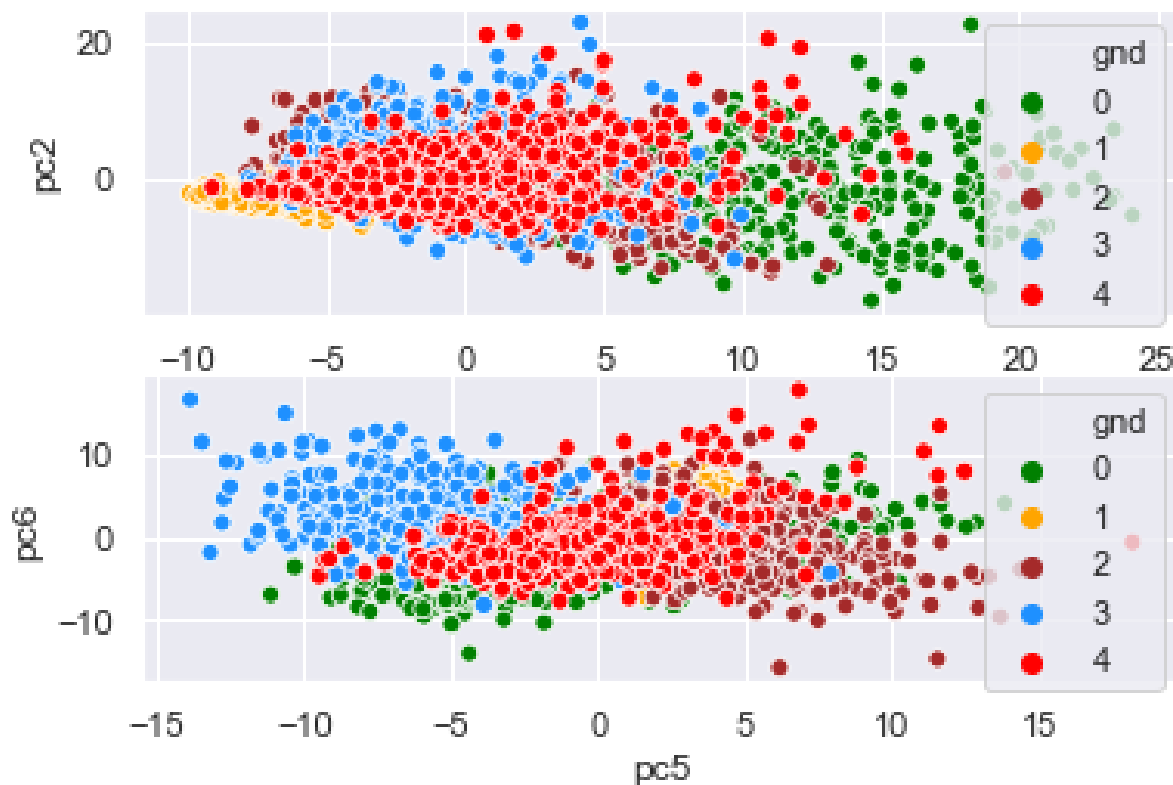
Scatter plot of the projected data with the top 20 eigenvalues

Figure 2: scatter plot of PCA



After observation, we found that the classification after fea.6 vs fea.7 is very vague. So the feature after 7 can be cut off. The result we draw from this picture is very close to that of above. The more clear the subplot's classification is, the greater eigenvalue is. But we cannot have a accurately have a conclusion of which point is best to cut in the scree plot.

fea.1 and fea.2 vs fea.5 and fea.6



We cannot tell which picture has a better classification towards the points, but 0 can be split from other points in the first plot, 3 can be separated in the second picture. Although the principle component in PCA takes up the most part, we cannot do a great job on classifying just according to the first feature and the second feature. More features are beneficial to classification. So we can use PCA and set `n_components` as 7 to get a good performance as good as all features and save a lot of time.

implement

PCA

Listing 1: PCA

```
1 time_start = time.perf_counter()
2 eigValueU , eigVectorU = linalg.eigh(np.dot(data.T, data))
3 indexU = eigValueU.argsort()[::-1]
4 eigValueU = eigValueU[indexU]
5 eigVectorU = eigVectorU[:,indexU ]
6 X_pca = eigVectorU.T.dot(data.T).T
7 time_over = time.perf_counter()
8 print("running time: " + str(time_over-time_start))
9 print(X_pca)
```

dual PCA using SVD

Listing 2: svd

```
1 def my_svd(X):
2     n, m = X.shape
3     eigValueU , eigVectorU = linalg.eigh(np.dot(X, X.T))
4     eigValueV , eigVectorV = linalg.eigh(np.dot(X.T, X))
5     indexU = eigValueU.argsort()[::-1]
6     indexV = eigValueV.argsort()[::-1]
7     eigValueU = eigValueU[indexU]
8     eigVectorU = eigVectorU[:,indexU ]
9     eigValueV = eigValueV[indexV]
10    eigVectorV = eigVectorV[:,indexV]
11    if n>m:
12        sigma=np.sqrt(eigValueU)
13    else:
14        sigma=np.sqrt(eigValueV)
15    return eigVectorU, sigma ,eigVectorV.T
```

Listing 3: dual PCA via SVD

```
1 time_start = time.perf_counter()
2 U,s,VT = my_svd(data.T)
3 time_over = time.perf_counter()
4 print("svd running time: "+str(time_over-time_start))
5 time_start = time.perf_counter()
6 X_dual_PCA = np.diag(s).dot(VT)[:784].T
7 time_over = time.perf_counter()
8 print("dual pca running time: "+str(time_over-time_start))
9 print(X_dual_PCA)
```

Time comparing

The running time of PCA is 0.35 s, svd is 2.63 s and dual PCA is 0.39 s. For this data set, the number of dimension is 784 (28 * 28), in the mean time, the number of samples is 2066. Only when the number of samples far greater than that of dimension, dual PCA will be faster than PCA. Because the running time of dual PCA depends on the number of the samples.

Prove that PCA is the best linear method for reconstruction

Suppose the line that make classification represented as:

$$L = b + \alpha v \quad (1)$$

assume $\|v\| = 1$ for convenience.

And each instance x_i is associated with a point on the line $\hat{x}_i = b + \alpha_i v$.

And reconstruction error can be written to:

$$R = \sum_{i=1}^m \|x_i - \hat{x}_i\|^2 \quad (2)$$

Our goal is to minimize the reconstruction error. We can rewrite equation (2) into:

$$R = \sum_{i=1}^m \|x_i - (b + \alpha_i v)\|^2 \quad (3)$$

We write the gradient of R wrt. α_i and set it to 0.

$$\frac{\partial R}{\partial \alpha_i} = 2\|v\|^2 \alpha_i + 2v x_i + 2bv = 0 \quad (4)$$

We draw from equation(4):

$$\alpha_i = v(x_i - b) \quad (5)$$

We write the gradient of R wrt. b and set it to 0.

$$\frac{\partial R}{\partial b} = 2mb - 2 \sum_{i=1}^m x_i + 2 \left(\sum_{i=1}^m \alpha_i \right) v = 0 \quad (6)$$

$$\sum_{i=1}^m \alpha_i = v^T \left(\sum_{i=1}^m x_i - mb \right) \quad (7)$$

By plugging (7) into (6), we get:

$$v^T \left(\sum_{i=1}^m x_i - mb \right) v = \left(\sum_{i=1}^m x_i - mb \right) \quad (8)$$

This is satisfied when $\left(\sum_{i=1}^m x_i - mb \right) = 0$, which means:

$$b = \frac{1}{m} \sum_{i=1}^m x_i \quad (9)$$

Substituting equation (5) into the optimization problem, we get a new optimization problem:

$$\max_v \sum_{i=1}^m v^T (x_i - b)(x_i - b)^T v \text{ (s.t. } \|v\|^2 = 1) \quad (10)$$

The Lagrangian is:

$$L(v, \lambda) = \sum_{i=1}^m v^T (x_i - b)(x_i - b)^T v + \lambda + \lambda \|v\|^2 \quad (11)$$

Let $S = \sum_{i=1}^m (x_i - b)(x_i - b)^T$, which we call it scatter matrix.

$$\frac{\partial L}{\partial v} = 2Sv - 2\lambda v = 0 \quad (12)$$

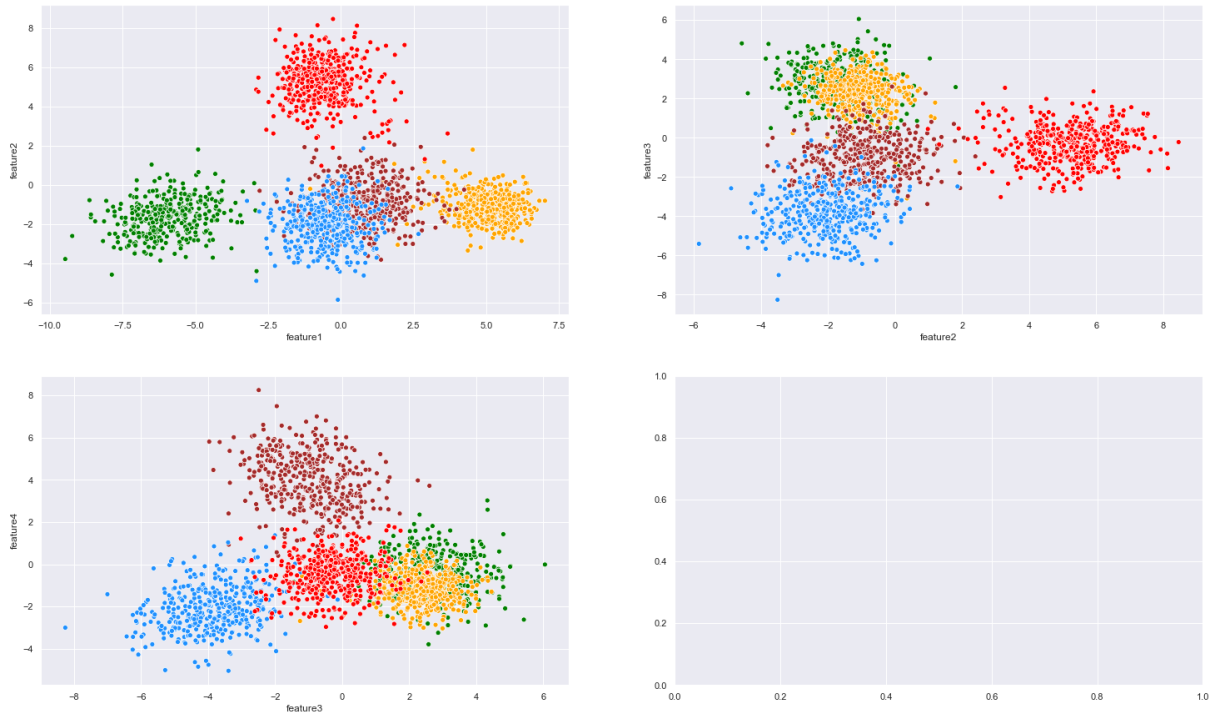
$$Sv = \lambda v \quad (13)$$

We can know from equation (13) that when v is the eigenvector of S , we can get the minimized reconstruction error. And in the mean time, the v in PCA is the eigenvector of $X^T X$, it is the same as the S matrix we said above. So PCA is the best linear method.

FDA

The plots

Figure 3: LDA



Analysis

Plot 1 (feature 1 vs. feature 2) can separate 4, 0 and 1 perfectly. Plot 2 (feature 2 vs. feature 3) can separate 3 and 4 perfectly. Plot 3 (feature 3 vs. feature 4) can separate 2, 3 and 4 perfectly. But each of them is indispensable for the classifying.

Compare the results of the LDA with the results obtained by using PCA.

Firstly, LDA can only take $n-1$ (n is the number of categories) features, so there are only 4 features extracted. The pictures of LDA look more clear than those of PCA. PCA is not a classifying method, but a dimension-reducing method. But in the mean time, LDA is a supervised classifying method. This is why LDA has a more clear classification than PCA.

Theoretical Question

The optimization of PCA is to maximize $\text{tr}(U^T S U)$ subject to $U^T U = I$, but in the mean time, LDA is to maximize $\frac{\text{tr}(U^T S_B U)}{\text{tr}(U^T S_W U)}$ subject to $U^T S_W U = I$. The matrix of U in PCA is the eigenvector of $S = X^T X$, but that of LDA is the eigenvector of $S_W^{-1} S_B$. And from the pictures, we can see that the scatter plots of LDA is far clearer than those of PCA, this is because PCA is actually not a classifying method, but a method to reduce the dimensions. But LDA will consider about the label of all the data, which is a supervised classification method.