

Introduction

This project requires the use of learned statistics to process real world data. It requires first performing principal component analysis on the dataset, then estimating the parameters of the reduced dimensionality samples, and finally calculating the error rate using the estimated parameters.

The dataset used in this project comes from images (with modifications) from the Fashion MNIST dataset. This project only uses images for class “T-shirt” and class “Sneaker”. There are 6000 T-shirt images and 6000 Sneaker images in the training dataset and 1000 T-shirt and 1000 Sneaker images in the testing dataset. T-shirt is labeled in “0” and Sneaker is labeled in “1”. Each image is 28*28 pixels.

Method

After I loaded the data, I reshaped each sample from the original 28*28 matrix into a 784-dimensional vector and calculated the sample mean and variance vectors for the training set. Then I normalize the sample set.

Next, I used **np.cov()** to get the dataset covariance matrix and used **np.linalg.eig()** to get the eigenvectors. I used the first and second eigenvectors as the first and second principal components. After that I used **np.linalg.norm()** to standardize principal components.

I then projected the normalized samples in the direction of the principal components and obtained their position in the new coordinate system. I visualized the training samples in this 2-d space in scatterplot and each of their classes looks like a normal distribution.

So I used MLE to estimate mean and variance for each class. I directly used the conclusion that the estimated mean for data in normal distribution is equal to its average and estimated variance is equal to data’s covariance matrix.

Finally I use the obtained parameters to calculate the error rate. For each of the data in the training and test sets, after normalizing them and projecting them in the direction of the principal components, I used the resulting mean and variance to calculate the likelihood $p(x|w_i)$. Since the prior probability is the same for both categories, this is done directly by comparing the likelihood values. The judgment is then compared to ground truth, and if it is wrong then loss +1. Finally, I calculate the loss/total sample number to get the minimum error rate.

Result and Observation

Part 1:

```
Part 1, mean vector = [6.6666667e-04 9.7500000e-03 7.3833333e-02 1.7041667e-01
2.2733333e-01 3.4683333e-01 7.9266667e-01 2.7800833e+00
7.7342500e+00 1.6573250e+01 2.9193667e+01 3.8230667e+01
3.2550000e+01 2.5306500e+01 2.4452333e+01 2.8843000e+01
3.9167750e+01 3.7331583e+01 2.3788916e+01 1.1966833e+01
4.8364167e+00 1.4420000e+00 4.5866667e-01 2.8900000e-01
1.6650000e-01 6.7250000e-02 7.4166667e-03 1.1000000e-02
3.6666667e-03 1.3000000e-02 1.5041667e-01 4.6391667e-01
7.4941667e-01 3.2333333e+00 1.3301333e+01 3.0809833e+01
4.8755333e+01 6.2815333e+01 7.1993250e+01 8.2118000e+01
8.6069250e+01 8.2520667e+01 8.1206667e+01 8.6787250e+01
8.7659167e+01 7.7689333e+01 6.8840500e+01 5.8079416e+01
4.1448000e+01 2.1757000e+01 6.9380000e+00 1.6586667e+00
6.7525000e-01 2.6650000e-01 6.0583333e-02 3.9833333e-02
7.4166667e-03 2.6333333e-02 2.2150000e-01 8.1475000e-01
```

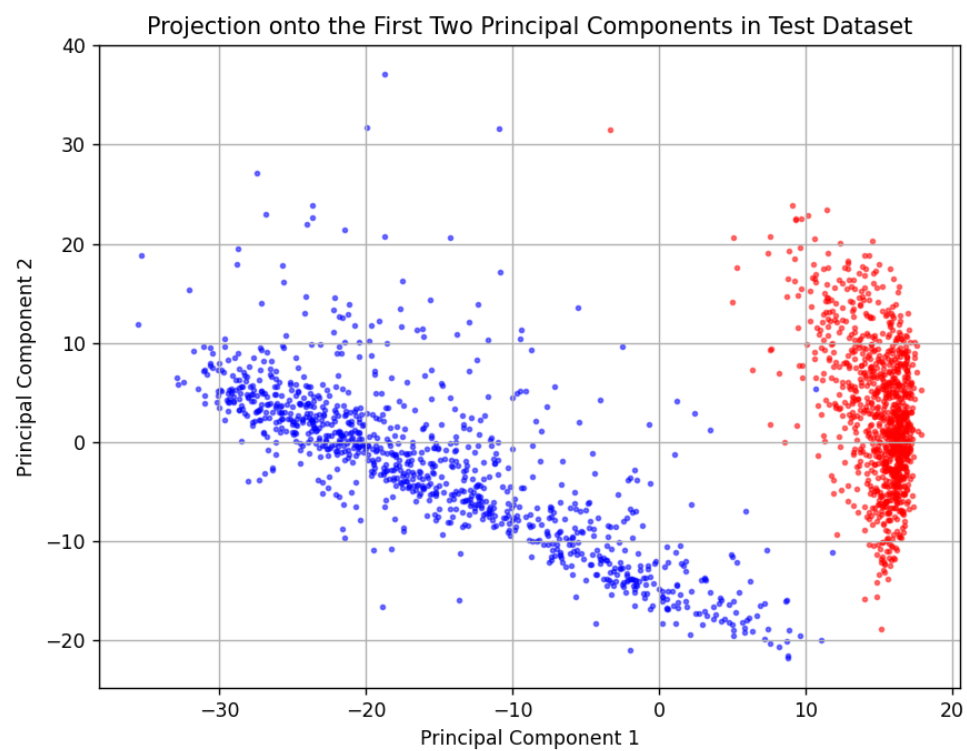
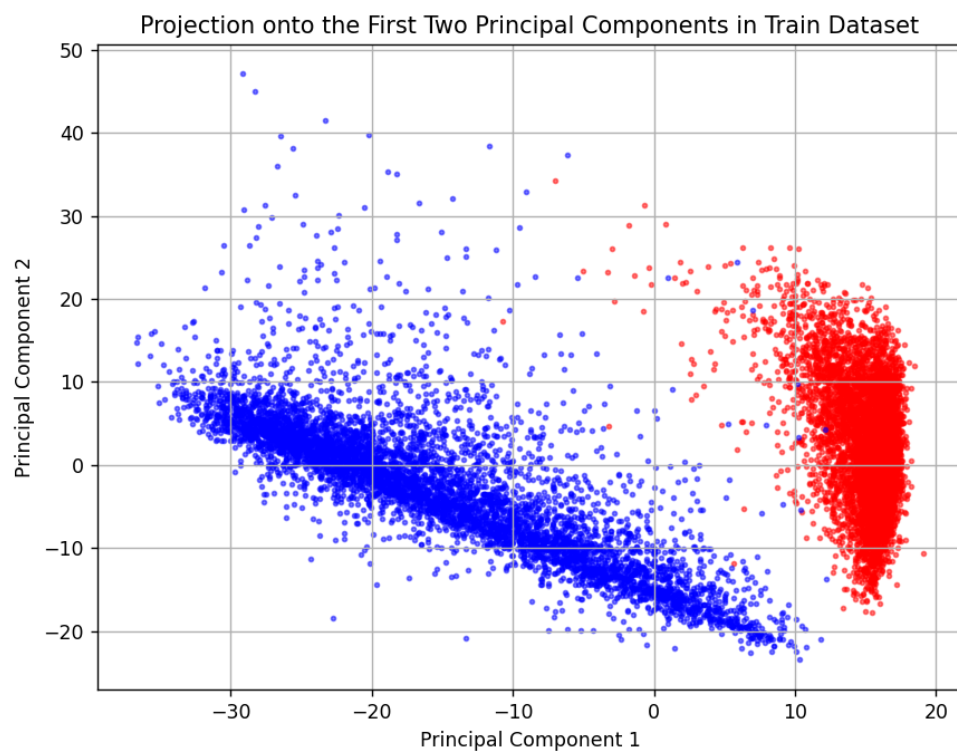
```
Part 1, std vector = [6.4546279e-02 2.8836597e-01 1.2368839e+00 2.7593552e+00
3.5117877e+00 4.4833997e+00 7.2684486e+00 1.5236268e+01
2.7662639e+01 4.0782107e+01 5.4184850e+01 6.3006898e+01
5.8434734e+01 5.0800653e+01 5.0221914e+01 5.4928308e+01
6.4862297e+01 6.2724745e+01 4.9657599e+01 3.4635392e+01
2.1066624e+01 1.0905553e+01 5.2230538e+00 3.9800727e+00
3.6979649e+00 1.9801247e+00 2.3985063e-01 1.1068178e+00
1.9617990e-01 3.1357774e-01 3.8934292e+00 7.4410257e+00
9.9758186e+00 1.8840940e+01 3.9466123e+01 6.0664590e+01
7.6688138e+01 8.7843063e+01 9.3277520e+01 9.8712786e+01
1.0328751e+02 1.0120898e+02 1.0050277e+02 1.0442457e+02
1.0365555e+02 9.6538796e+01 9.1795839e+01 8.4323692e+01
7.0730960e+01 5.1390428e+01 2.8754585e+01 1.4703689e+01
9.8380615e+00 5.8630178e+00 2.6958758e+00 2.5063412e+00
2.1062207e-01 9.6668155e-01 5.3238398e+00 1.0144108e+01
2.0894584e+01 5.0880357e+01 7.5722869e+01 8.4488315e+01
8.7259147e+01 8.8333532e+01 9.0284841e+01 9.2260337e+01
9.5176143e+01 9.5365969e+01 9.5508286e+01 9.6702207e+01
9.3570770e+01 9.0746362e+01 8.9131630e+01 8.7948537e+01
8.6608891e+01 8.2280766e+01 6.5383683e+01 3.4971909e+01
```

Part 2:

```
=====
First PCA: [-6.64461141e-04 -2.10239358e-03 -3.72519317e-03 -3.29025897e-03
-4.01885380e-03 -4.96809009e-03 -6.35460152e-03 -1.01534045e-02
-1.64850021e-02 -2.41267167e-02 -3.22589368e-02 -3.61794719e-02
-3.29855480e-02 -2.98896816e-02 -2.92143258e-02 -3.12336583e-02
-3.51197480e-02 -3.49700847e-02 -2.83077474e-02 -2.01043561e-02
-1.26913609e-02 -7.07754442e-03 -5.00744050e-03 -3.90573083e-03
-2.65044469e-03 -1.53462408e-03 -1.49134011e-03 -2.75162426e-04
-5.43098285e-04 -2.22896564e-03 -2.28417575e-03 -4.04821449e-03
-5.10328585e-03 -1.12836135e-02 -2.21666963e-02 -3.33053800e-02
-4.20871667e-02 -4.72122783e-02 -4.99110649e-02 -5.17274921e-02
-5.02451414e-02 -4.79064868e-02 -4.68806977e-02 -4.90809960e-02
-5.08006516e-02 -5.06108897e-02 -4.88629329e-02 -4.53594349e-02
```

```
Second PCA: [ 2.65302063e-04 5.24337961e-03 5.53663175e-03 4.33290888e-03
6.64708440e-03 5.65805224e-03 4.76100149e-03 2.31190846e-03
2.18227454e-04 -9.48351797e-03 -2.08669002e-02 -2.68049906e-02
-2.37244946e-02 -1.94199903e-02 -1.85653162e-02 -2.14317786e-02
-2.85795673e-02 -2.56994764e-02 -1.45396953e-02 -4.30601053e-03
1.19669468e-03 3.65954115e-03 4.25778360e-03 4.21387901e-03
6.52520594e-03 4.08876080e-03 3.38069453e-04 4.89777315e-03
1.59128192e-03 5.14883157e-03 7.08527307e-03 9.24331382e-03
1.06499799e-02 6.39207411e-03 9.62553250e-04 -3.65532805e-03
-6.22977022e-03 -9.47465608e-03 -1.40809685e-02 -2.19096287e-02
-2.74796979e-02 -3.07504609e-02 -3.20491306e-02 -3.03713013e-02
-2.80041767e-02 -1.80838263e-02 -1.05648573e-02 -7.56323561e-03
-5.94273633e-03 -2.90920151e-03 3.83932834e-03 8.57972462e-03
```

Part 3:



Part 4 and 5:

```
=====
Estimated mean of T shirt: [-14.96013579 -2.53759713]
Estimated cov of T shirt: [[103.38209995 -69.15661464]
[-69.15661464 73.43418283]]
Estimated mean of Sneaker: [14.96013579 2.53759713]
Estimated cov of Sneaker: [[ 4.94656978 -6.78163691]
[-6.78163691 54.67473161]]
=====
Error Rate in Training Dataset: 0.0025
Error Rate in Testing Dataset: 0.001
PS E:\workplace\CSE-569> 
```

Conclusion

This project was very successful in performing principal component analysis on the training set and estimating the parameter of the projection in the direction of the principal components. And the error rate of the final speculation on the samples of the test set is very small.