

Algorithm of K-means

Given samples $X = \{x_1, x_2, \dots, x_n\}$, where for each x_i represents a sample, the purpose of k means is arranging those samples into k different classes.

It use the squared Euclidean distance as the distance between samples:

$$d(x_i, x_j) = \sum_{k=1}^m (x_{ki} - x_{kj})^2$$

$$= \|x_i - x_j\|^2$$

The loss function is defined as the sum of the distance of different samples to its centroids.

$$W(c) = \sum_{l=1}^k \sum_{C(i)=l} \|x_i - x_l\|^2$$

The purpose of K means is optimise this loss function:

$$C^* = \operatorname{argmin} W(C)$$

In realistic, we use iteration to solve this optimisation, because this may be a NP hard problem.

The iteration steps are:

input: X

output: cluster of X

- steps:
- (1) initialise: random select k samples as initial centroids
 - (2) calculate all distances between X and initial centroids, arranging each samples into its nearest class
 - (3) calculate new centroids, which is the mean value of each feature
 - (4) if the iterations converge or satisfy the condition. Stop and output cluster

The time will spend $O(mnk)$.

Algorithm of PCA

Given samples $X = \{x_1, x_2, \dots, x_n\}$, where for each x_i represents a sample:

$$X = [x_1 \ x_2 \ \dots \ x_n] = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix} \quad (1)$$

Then we can get covariance S:

协方差矩阵

公式: $S = [s_{ij}]_{m \times m}$

实数 $X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}$

对 X 转置

$$(注1) \quad s_{ij} = \frac{1}{n-1} \sum_{k=1}^n (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j), \quad i, j = 1, 2, \dots, m \quad (1)$$

And its Correlation matrix R:

样本相关矩阵 R 为 (1) 1-10

公式: $R = [r_{ij}]_{m \times m}, \quad r_{ij} = \frac{s_{ij}}{\sqrt{s_{ii}s_{jj}}}, \quad i, j = 1, 2, \dots, m$

标准差: $\sigma = \sqrt{\operatorname{var}(x)}$

Our purpose is finding the eigenvalues and eigenvectors of R , because these eigenvectors are **orthogonal**, so these values and vectors can be seen as the most important part of getting matrix's effective informations. Which means they can reduce redundancy information or noise.

Arranging the eigenvalues from the largest to smallest.

$$\lambda_1 \geq \lambda_2 \geq \dots$$

Then the correlated eigenvectors becomes the main component of X :

$$V = (v_1, v_2, \dots, v_k)$$

Thus we can get $Y(k, n) = V.T(k, m) * X(m, n)$
The Y is the target matrix after dimension reduced.

However, in real world, we cannot guarantee these eigenvectors are orthogonal, so we offer use **SVD**.

This is my study note (please ignore the Chinese language part):

图 16.3 因子负荷量的分布图

16.2.3 数据矩阵的奇异值分解算法

给定样本矩阵 X , 利用数据矩阵奇异值分解进行主成分分析。具体过程如下。这里假设有 k 个主成分。

$A^T A$ eigenvector V
 $A A^T$ eigenvector U

$A = U \Sigma V^T = (\text{rotate})(\text{stretch})(\text{rotate})$
 $= \begin{bmatrix} | & | \\ u_1 & u_2 \\ | & | \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix} \begin{bmatrix} -v_1^T \\ -v_2^T \end{bmatrix}$

$A^T A = (V \Sigma^T U^T) U \Sigma V^T = V (\underbrace{\Sigma^T \Sigma}_{\text{diagonalization}}) V^T$
 $\lambda \text{ for } A^T A = \sigma^2 \text{ for } A$

We can know that, for our X , we can use $X.T * X$ to get the eigenvectors.

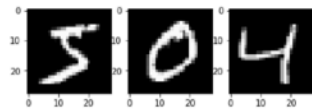
Suppose $X' = 1/\sqrt{n-1} * X.T$
Then $X'.T.X' = 1/(n-1) * X * X.T = R(\text{of } X)$
So the problem change into finding the SVD of X' , where $X' = U * \sigma * V.T$
Then the first k rows of $V.T$ compose the k numbers of principal components.
Thus we can get $Y(k, n) = V.T(k, m) * X(m, n)$

Experiments:

Preprocessing the data:

One thing I want to discuss here is the image performance after using standard scaler:

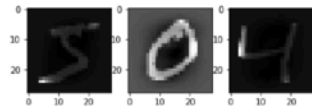
```
ax0.imshow(train_features[0].reshape(28,28), cmap='gray')
ax1.imshow(train_features[1].reshape(28,28), cmap='gray')
ax2.imshow(train_features[2].reshape(28,28), cmap='gray')
plt.show()
```



standard scaler

```
In [11]: from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
train_features=scaler.fit_transform(train_features)
test_features=scaler.fit_transform(test_features)
```

```
In [12]: fig,axs = plt.subplots(1,3)
ax0,ax1,ax2 = axs.ravel()
ax0.imshow(train_features[0].reshape(28,28), cmap='gray')
ax1.imshow(train_features[1].reshape(28,28), cmap='gray')
ax2.imshow(train_features[2].reshape(28,28), cmap='gray')
plt.show()
```



From the picture we can see that the structure of each number still remains, it seems that the only change is the shade of colour, which mainly because standard scaler does not change the distribution of data, it only narrow the data into -1 and +1.

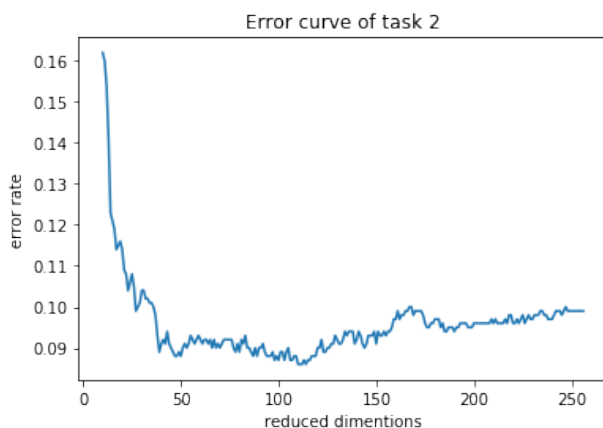
Task 1:

In task 1 I build 2 models, one is the traditional PCA and the other using SVD. Because I want to see whether they have the similar trend and their performance.

Task2:

This picture is the PCA using SVD, we can see some informations from it.

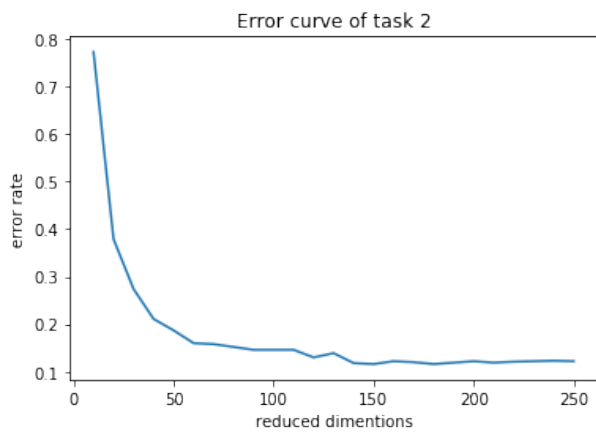
- (1) The error curve sharp decreased from 10 to 50, and it becomes steady flat after 50, which can be explained that the front 50 eigenvectors are very important that can not be ignored. The eigenvalues that after 50 are much more smaller and we cannot get too much informations from them.
- (2) The curve slightly arise during 100 to 250 and then becomes stable (because I checked dimensions in [512, 784], and their error is [0.102, 0.103])



```
In [16]: error2 = []
for i in [512,784]:
    temp = PCA(train_features.T,i)
    train_x = temp.fit(train_features.T)
    test_x = temp.fit(test_features.T)
    error = KNN(train_x,train_label,test_x).error(test_label)
    error2.append((error,i))
print(error2)

[(0.102, 512), (0.103, 784)]
```

Compared to SVD, the traditional PCA performs similar:

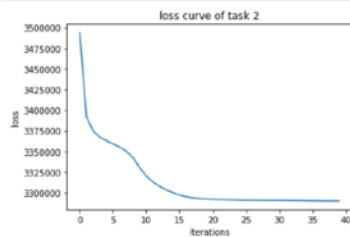


This picture's dimensions interval is 10, so it seems more smooth, but both PCA and SVD have similar trend and all of their errors get narrowed around 0.1.

Task 3:

My iteration sets to 100, and the loss curve looks like:

```
plt.figure()
plt.title('loss curve of task 2')
plt.xlabel('iterations')
plt.ylabel('loss')
plt.plot(x_list,y_list)
plt.show()
```



task 4 -random create centroids

As we can see, the loss value decreases dramatically before the first 20 iterations, and then becomes stable. This implies that the curve has ability to converge. But if updating the centroids can not reduce the loss anymore, then theoretically, all the samples are almost classified correctly, then the accuracy should high.

Task4:

Without using PCA, K means performs like:

```
In [43]: k_object2 = KMeans_task4(train_features,train_label,new_cent)
k_object2.fit(train_features)
k_object2.calculate_percentage(train_label)
k_object2.report()
k_object2.average_precision()

20
TP: True Positive
FP: False Positive
Precision = TP/(TP+FP)
class original predict true TP FP precision
0 502 401 364 364 37 0.9077306733167082
1 671 1070 657 657 413 0.6140186915887851
2 581 335 191 191 144 0.5701492537313433
3 608 451 155 155 296 0.343680709343681
4 623 1105 262 262 843 0.23710407239819004
5 514 520 237 237 283 0.45576923076923076
6 608 560 470 470 59 0.8392857142857143
7 651 178 154 154 24 0.8651605393258427
8 551 885 194 194 691 0.2192098395480226
9 601 495 177 177 318 0.3575757575757576
```

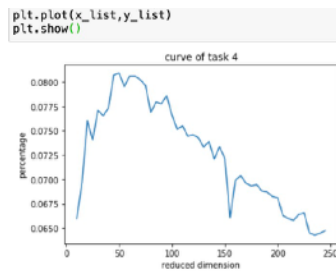
Out[43]: 0.5409691682073963

To (601): average stat list longer off

With PCA, K means performs like:

Why the curve arise, and then decrease?

I guess... when the dimensions are smaller than 50, the important informations are not sufficient,
So the curve increases when dimensions arise.

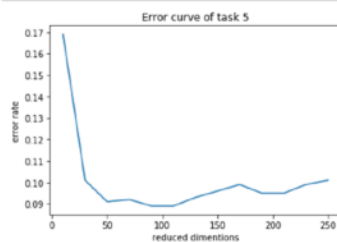


Task 5:

The original image is composed with 0 and 1, so I use generate 256 0 and 1 randomly and add them to the original data. Lets call it new data

Then I use the standard scaler to the new data for PCA,
Then, I get PCA picture:

```
In [38]: y_list = [i[0] for i in error_plot_list]
x_list = [i[1] for i in error_plot_list]
plt.figure()
plt.title('Error curve of task 5')
plt.xlabel('reduced dimention')
plt.ylabel('error rate')
plt.plot(x_list,y_list)
plt.show()
```



We can find that this curve is extremely similar to the task 2- curve. They all have a sharp decreasing before 50 and then becomes stable around 0.1.

The reason of that is PCA can decrease noise. The main idea of PCA is getting the unique(the most important) informations, replacing the redundant informations and abandoning rubbish informations. Thus PCA can do quite well in reducing rubbish features.

Besides, the SVM looks like:

```
>>> predict test result
y_pred_with_scale = clf.predict(X_test)
print("Prediction report with scale:")
print(classification_report(y_test, y_pred_with_scale))
```

Prediction report with scale:				
	precision	recall	f1-score	support
0.0	0.97	1.00	0.99	99
1.0	0.95	0.98	0.97	113
2.0	0.98	0.94	0.96	94
3.0	0.96	0.88	0.92	108
4.0	0.96	0.95	0.95	93
5.0	0.92	0.94	0.93	96
6.0	0.96	0.97	0.97	101
7.0	0.96	0.92	0.94	103
8.0	0.95	0.92	0.93	99
9.0	0.88	0.98	0.92	94
accuracy			0.95	1000
macro avg	0.95	0.95	0.95	1000
weighted avg	0.95	0.95	0.95	1000

We can see that the accuracy is 0.95, which is quite high. It implies that SVM can also performs well when dealing with noise, which is also called **outliers**. As we have discussed in the former assignment, it is the points that near the SVM boundary important, the outliers can hardly influence the boundary, thus SVM can also performs well even with noise.