

哈尔滨工业大学计算机科学与技术学院

## 实验报告

课程名称： 机器学习

课程类型： 选修

实验题目： PCA

学号： 1190202107

姓名： 姚舜宇

## 一、实验目的

实现一个 PCA 模型，能够对给定数据进行降维。

## 二、实验要求及实验环境

1. 人工生成一些数据，让它们主要分布在低维空间中，然后使用 PCA 模型对数据进行主成分提取。
2. 找一个人脸数据，使用 PCA 方法进行降维，找出一些主成分，然后用这些主成分对每一副人脸图像进行重建，并比较它们与原图的差别。

环境：Windows10, python3.9, pycharm

## 三、设计思想（本程序中的用到的主要算法及数据结构）

定义数据集  $X = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_N]^T_{N \times p} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Np} \end{bmatrix}_{N \times p}$  ,  $\mathbf{x}_i \in \mathcal{R}^p, i=1, 2, \cdots, N$ 。

可以计算样本均值  $\bar{\mathbf{x}}_{p \times 1} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \frac{1}{N} [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_N] \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \frac{1}{N} X^T \mathbf{1}_N$

样本协方差矩阵  $S = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T = \frac{1}{N} X^T H X$  , 其中  $H = I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T$

定义  $p$  维向量  $\mathbf{x}$  到  $p$  维向量  $\mathbf{y}$  的线性变换

$$\mathbf{y} = A\mathbf{x}$$

其中  $A = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_p]^T = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ a_{21} & a_{22} & \cdots & a_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ a_{p1} & a_{p2} & \cdots & a_{pp} \end{bmatrix}$  ,  $\mathbf{a}_i = \begin{bmatrix} a_{i1} \\ a_{i2} \\ \vdots \\ a_{ip} \end{bmatrix}$  ,  $\mathbf{a}_i^T \mathbf{a}_i = 1, i=1, 2, \cdots, p$ 。

现计算最大投影方差  $J$ 。假设某一个单位方向向量  $\mathbf{a}_i$ ，投影方差为：

$$\begin{aligned} J &= \frac{1}{N} \sum_{i=1}^N [(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{a}]^2 \\ &= \frac{1}{N} \sum_{i=1}^N \mathbf{a}^T (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{a} \\ &= \mathbf{a}^T S \mathbf{a} \end{aligned}$$

现在问题就转化为优化问题：

$$\begin{cases} \hat{\mathbf{a}} = \arg \max_{\mathbf{a}} \mathbf{a}^T S \mathbf{a} \\ s.t. \ \mathbf{a}^T \mathbf{a} = 1 \end{cases}$$

使用拉格朗日乘子法，定义  $L(\mathbf{a}, \lambda) = \mathbf{a}^T S \mathbf{a} + \lambda(1 - \mathbf{a}^T \mathbf{a})$ 。对  $\mathbf{a}$  求导并令导数为零，可以得到  $S\mathbf{a} = \lambda\mathbf{a}$ 。可以看出  $\mathbf{a}$  即协方差矩阵的特征向量。又因为逐一选取方差最大方向等价于直接取最大的一些特征值，所以问题就转化为求协方差矩阵的特征向量。

下面简要推导特征值分解。令  $\Lambda = \begin{bmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \lambda_p \end{bmatrix}$  为特征值矩阵，其中  $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ ，

$P = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_p]$  为特征向量矩阵，顺序与特征值矩阵对应。由  $S\mathbf{a} = \lambda\mathbf{a}$ ，扩展到矩阵形式，即为  $SP = P\Lambda$ ， $S = PAP^{-1}$ 。

计算出特征向量  $\mathbf{a}_i$  后，即可得到变换矩阵  $A = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_p]^T$ 。对于当前数据集中的某一个线性相关的样本  $\mathbf{x}_i$ ，经过变换得到线性无关的样本  $\mathbf{y}_i = A(\mathbf{x}_i - \bar{\mathbf{x}})$ 。转化为矩阵形式，当前数据集为  $X = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]^T$ ，设变换之后的数据集为  $Y = [\mathbf{y}_1 \ \mathbf{y}_2 \ \dots \ \mathbf{y}_N]^T$ 。由线性变换公式可以得到  $Y^T = A(X^T - \bar{\mathbf{x}}\mathbf{1}_N^T)$ ，即  $Y = (X - \mathbf{1}_N \bar{\mathbf{x}}^T)A^T$ 。

需要降维时，一般选择  $q$  个主成分（线性无关变量）来代替  $p$  个原有的线性相关的变量，使问题得以简化，并能保留原有变量的大部分信息。对于选择  $q$  个主成分的情况，变换公式为  $\mathbf{y}_i = B(\mathbf{x}_i - \bar{\mathbf{x}})$ ，代入数据集，即为  $Y = (X - \mathbf{1}_N \bar{\mathbf{x}}^T)B^T$ ，其中  $B = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_q]^T$ ，为  $q \times p$  矩阵。

为了选择合适的主成分个数  $q$ ，需要定义第  $k$  主成分的方差贡献率  $\eta_k = \frac{\lambda_k}{\sum_{i=1}^p \lambda_i}$ ， $p$  个主成分

的累计方差贡献率  $\sum_{i=1}^q \eta_i = \frac{\sum_{i=1}^q \lambda_i}{\sum_{i=1}^p \lambda_i}$ 。累计方差贡献率反映了主成分保留信息的比例。通常取

$q$  使得累计方差贡献率达到规定的百分比以上，如 70%~80% 以上。

若要重建图像，由于  $BB^T = I_q$ ，所以根据式  $Y = (X - \mathbf{1}_N \bar{\mathbf{x}}^T)B^T$ ， $X = YB + \mathbf{1}_N \bar{\mathbf{x}}^T$ ，每一行即为一副图像经过降维之后的特征。将其重新修改尺寸为原来的大小，即可和原图计算信噪比。均方误差计算公式为： $MSE = \frac{1}{mn} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} \|I(i,j) - J(i,j)\|^2$ 。信噪比计算公式为：

$$PSNR = 10 \log_{10} \left( \frac{MAX_I^2}{MSE} \right) = 20 \log_{10} \left( \frac{MAX_I}{\sqrt{MSE}} \right)$$

在此问题中为 255。一般情况下，当 PSNR 大于 40 时，说明图像质量好，即与原图非常接近；大于 30 小于 40 时，说明图像质量较好，有失真但可以接受；大于 20 小于 30 时，说明图像质量差；小于 20 时，图像质量不可接受。

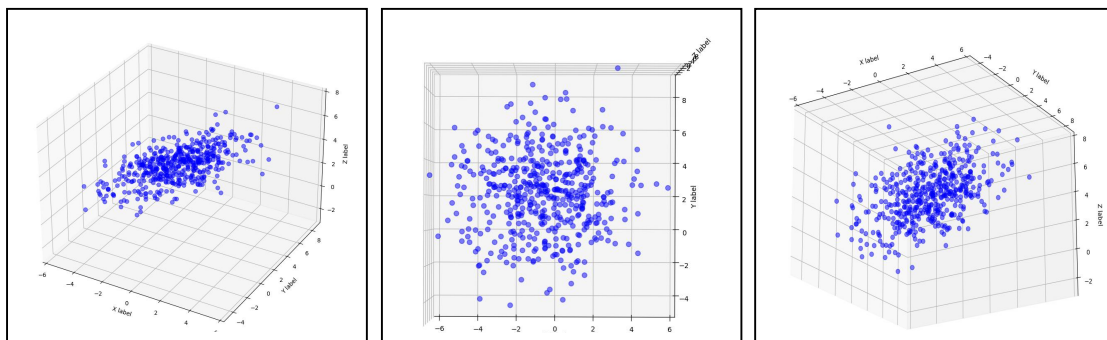
## 四、实验结果与分析

### 1. 使用自己生成的数据

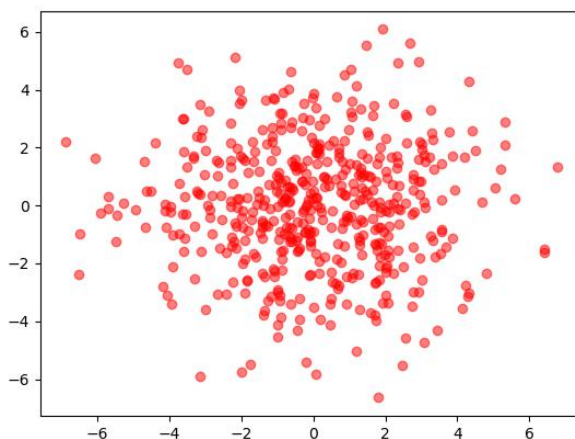
#### 1) 使用三维高斯分布的数据

假设样本量为 500，三个维度的均值分别为[1,2,3]，协方差矩阵为  $\begin{bmatrix} 5 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 0.5 \end{bmatrix}$ ，关

于 y 轴旋转  $30^\circ$ ，生成样本效果见下图。



可以看到某一方向上的方差较小，即旋转之前的 z 轴方向，因为协方差矩阵中 z 轴方向的样本方差较小。使用 PCA 进行降到 2 维，生成的结果如下：



从结果图可以看出，当降到 2 维时，两个维度分布较为均衡，原数据中方差最小的方向上的数据在降维过程中得到了一定程度上的忽略。根据输出结果，上述情况下，由三维降到二维，数据信息保留了 95.95%。

使用不同的协方差矩阵进行测试。

当协方差矩阵为  $\begin{bmatrix} 5 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 5 \end{bmatrix}$ ，即三个维度方差相等时，定性来说三个维度信息量相

近，经过降维效果应当表现较差。经测试，由三维降到二维，数据信息保留了 70.34%。

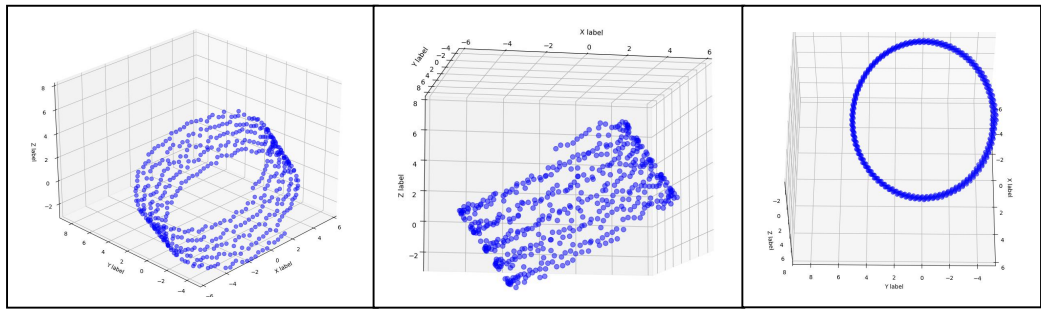
当协方差矩阵为  $\begin{bmatrix} 5 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 2.5 \end{bmatrix}$ ，同理分析并经测试，由三维降到二维，数据信息保

留了 81.20%。

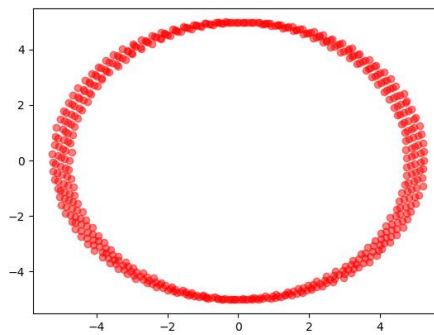
#### 2) 使用三维螺旋线数据

假设样本量为 500，z 的范围为[-3,3]。在竖坐标上加方差为 0.2 的高斯噪声，关于 y

轴旋转  $30^\circ$ ，生成样本效果见下图。

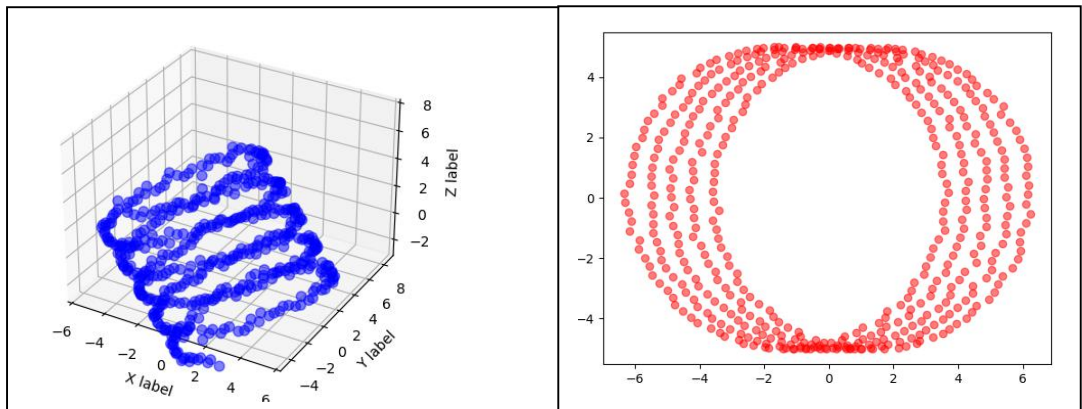


使用 PCA 降到 2 维，效果如下：



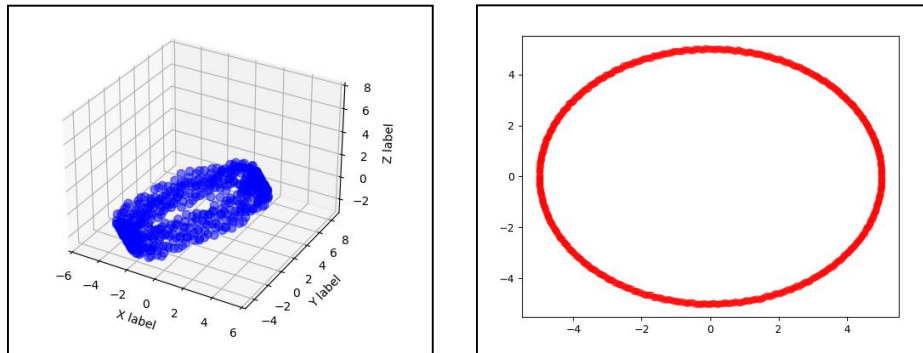
输出信息保留率为 89.29%。如果在生成数据时将旋转之前  $z$  的范围调整到  $[-5,5]$  或  $[-1,1]$ ，根据理论，保留率应该分别降低和升高。实验结果如下。

$z \in [-5, 5]$ ：



信息保留率为 76.50%；

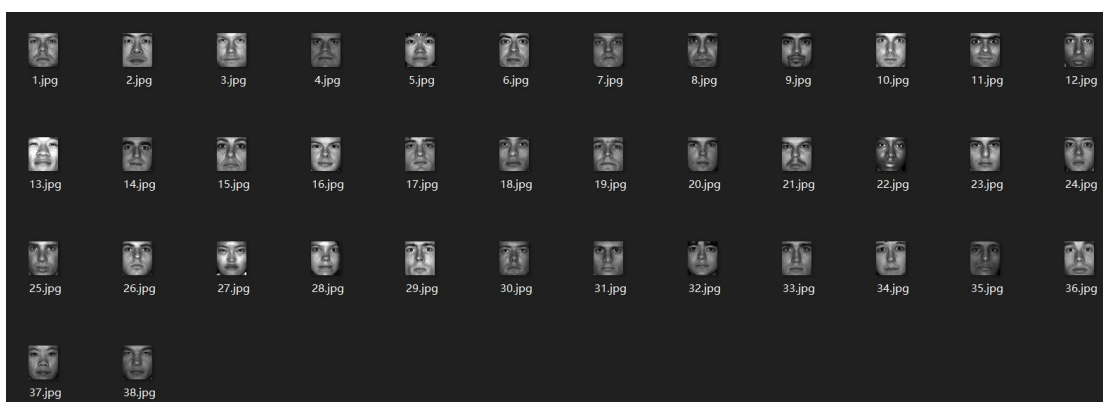
$z \in [-1, 1]$ ：



信息保留率为 98.55%。

## 2. 使用人脸数据

在网站上下载了灰度人脸数据集，选择了光线和视角较好的 38 张图片。由于计算机内存限制，难以对太大的矩阵进行操作，所以使用 `opencv` 将图片压缩到 42\*48，效果如下。



分别使用至少 95%、90%、85%和 60%的信息保留率对图像进行降维压缩，并和原始图像计算信噪比。

考虑到图像较多，所以选择少量人脸进行展示。从左到右依次为源图像、95%、90%、85%和 60%。



明显看到从左到右依次变得模糊。

下面简要分析计算出的信噪比情况。

信噪比计算公式为： $PSNR = 10\log_{10}\left(\frac{MAX_I^2}{MSE}\right) = 20\log_{10}\left(\frac{MAX_I}{\sqrt{MSE}}\right)$ ， $MAX_I$  表示图像

点颜色的最大值， $MSE$  表示均方误差，公式为： $MSE = \frac{1}{mn} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} \|I(i,j) - J(i,j)\|^2$ 。

在此问题中为 255。一般情况下，当 PSNR 大于 40 时，说明图像质量好，即与原图非常接近；大于 30 小于 40 时，说明图像质量较好，有失真但可以接受；大于 20 小于 30 时，说明图像质量差；小于 20 时，图像质量不可接受。

源图像的维度为 42\*48=2016。对于保留率达到 95%的图像，降维之后的维度为 25，PSNR

都在 30 以上 40 以下，均值为 33.3。对于保留率达到 90% 的图像，降维之后的维度为 18，PSNR 都在 29 以上 37 以下，均值为 31.1。对于保留率达到 85% 的图像，降维之后的维度为 14，PSNR 都在 30 左右，均值为 30.3。对于保留率达到 65% 的图像，降维之后的维度为 4，PSNR 都在 28 到 30，均值为 29.0。

在实验过程中，有时特征值会出现虚部，处理方式保留实部，丢弃虚部，不影响结果。

## 五、结论

主成分分析是一种常用的无监督学习方法，首先利用正交变换把由线性相关变量表示的观测数据转换为少数几个由线性无关变量表示的数据。计算中样本协方差矩阵进行了特征值分解，按照要求取前几个较大的特征值，其余舍去。一般来说，前几个特征值就占据了所有特征值之和的大部分，所以 PCA 能够在仅损失少量信息的情况下对数据进行较大程度的降维。

## 六、参考文献

《统计学习方法》 李航

《机器学习》 周志华

## 七、附录：源代码（带注释）

pca.py

```
1. import numpy as np
2.
3.
4. def cal_cov_matrix(X):
5.     """
6.     计算 X 的协方差矩阵
7.     :param X: 样本矩阵 N*M, N 个样本, 每个样本 M 维
8.     :return: X 的协方差矩阵
9.     """
10.    n = X.shape[0]
11.    H = np.mat(np.identity(n)) - np.multiply(1 / n, np.mat(np.ones(n)).T @ np.mat(np.ones(n)))
12.    return np.multiply(1 / n, X.T @ H @ X)
13.
14.
15. def PCA_solve(X, eta=float(1), k=0):
16.     """
17.     PCA 求解
18.     :param X: 样本矩阵
19.     :param eta: 保留信息比例
20.     :param k: 保留特征值数
21.     :return: 降维后中心化的样本 (N*降过之后的维度), 变换矩阵, 原样本均值
22.     """
23.    S = cal_cov_matrix(X)
```

```

24.     mean = np.multiply(1 / X.shape[0], X.T @ np.mat(np.ones(X.shape[0])).T)
25.     eigval, eigvec = np.linalg.eig(S)
26.     sorted_indices = np.argsort(eigval)
27.     eigvec = eigvec[:, sorted_indices[::-X.shape[1] - 1:-1]]
28.     eigval = sorted(eigval, reverse=True)
29.     cnt = 0
30.     sum_eigval = np.sum(np.array(eigval))
31.     pre_sum = 0
32.     if 1 > eta > 0 and k == 0:
33.         for i in range(len(eigval)):
34.             pre_sum += eigval[i]
35.             if pre_sum / sum_eigval > eta:
36.                 cnt = i + 1
37.                 break
38.     pre_sum = 0
39.     if k > 0:
40.         cnt = k
41.         for i in range(cnt):
42.             pre_sum += eigval[i]
43.         print('维度: '+str(cnt))
44.         print('信息保留百分比: ' + str(np.real(pre_sum / sum_eigval)))
45.     B = eigvec.T[0]
46.     for i in range(1, cnt):
47.         B = np.vstack((B, eigvec.T[i]))
48.     B = np.real(B)
49.     return (X - mean.T) @ B.T, B, mean

```

mydata.py

```

1. import numpy as np
2. from pca import PCA_solve
3. import matplotlib.pyplot as plt
4.
5.
6. def Generate_Gaussian_data():
7.     """
8.     生成一组高斯分布的数据
9.     :return: 高斯数据矩阵
10.    """
11.    num = 500
12.    mean = [1, 2, 3]
13.    cov = [[5, 0, 0], [0, 5, 0], [0, 0, 2.5]]
14.    data = np.random.multivariate_normal(mean, cov, num)
15.    return np.mat(data)

```



```

16.
17.
18. def Generate_Spiral_data():
19.     """
20.     生成一组螺旋线数据
21.     :return: 螺旋线数据矩阵
22.     """
23.     num = 500
24.     theta = np.linspace(-5*np.pi, 5*np.pi, num)
25.     z = np.linspace(-1, 1, num)
26.     sigma = np.random.normal(loc=0, scale=0.2, size=num)
27.     z += sigma
28.     r = 5
29.     x = r * np.sin(theta)
30.     y = r * np.cos(theta)
31.     X = np.mat(x)
32.     X = np.vstack((X, np.mat(y)))
33.     X = np.vstack((X, np.mat(z)))
34.     return X.T
35.
36.
37. def rotate(X, axis, theta):
38.     """
39.     对数据 X 沿着某一个轴进行旋转一定的角度
40.     :param X: 数据矩阵
41.     :param axis: 轴
42.     :param theta: 旋转角度
43.     :return: 旋转之后的数据矩阵
44.     """
45.     if axis == 'x':
46.         rotate_matrix = np.mat([[1, 0, 0], [0, np.cos(theta), np.sin(theta)],
47.             , [0, -np.sin(theta), np.cos(theta)]])
48.         return X @ rotate_matrix.T
49.     elif axis == 'y':
50.         rotate_matrix = np.mat([[np.cos(theta), 0, -np.sin(theta)], [0, 1, 0],
51.             , [np.sin(theta), 0, np.cos(theta)]])
52.         return X @ rotate_matrix.T
53.     elif axis == 'z':
54.         rotate_matrix = np.mat([[np.cos(theta), np.sin(theta), 0], [-np.sin(
55.             theta), np.cos(theta), 0], [0, 0, 1]])
56.         return X @ rotate_matrix.T
57.     else:
58.         print('wrong')
59.         exit(0)

```

```

57.
58.
59. def show3D(X):
60.     """
61.     将数据在三维空间中绘制出来
62.     :param X: 数据
63.     :return: 无
64.     """
65.     plt.figure(1)
66.     ax1 = plt.axes(projection='3d')
67.     ax1.scatter3D(np.array(X)[: , 0], np.array(X)[: , 1], np.array(X)[: , 2], c
        = 'b', marker='o', s=49, alpha=0.5)
68.     ax1.set_xlim(-6, 6)
69.     ax1.set_ylim(-5, 9)
70.     ax1.set_zlim(-3, 8)
71.     ax1.set_xlabel('X label')
72.     ax1.set_ylabel('Y label')
73.     ax1.set_zlabel('Z label')
74.
75.
76. def show2D(Y):
77.     """
78.     将数据在二维空间中绘制出来
79.     :param Y: 数据
80.     :return: 无
81.     """
82.     plt.figure(2)
83.     plt.scatter(np.array(Y)[: , 0], np.array(Y)[: , 1], c='r', marker='o', s=3
        6, alpha=0.5)
84.
85.
86. if __name__ == '__main__':
87.     # X = Generate_Gaussian_data()
88.     X = Generate_Spiral_data() # 生成三维的数据
89.     X = rotate(X, 'y', np.pi / 6)
90.     Y, B, m = PCA_solve(X, k=2) # 把三维的数据降到二维
91.     plt.ion()
92.     show3D(X)
93.     show2D(Y)
94.
95.     plt.ioff()
96.     plt.show()

```

```

1. import numpy as np
2. from pca import PCA_solve
3. import cv2 as cv
4.
5.
6. def compression(eta):
7.     """
8.     对人脸数据集进行压缩
9.     :param eta: 保留信息比例
10.    :return: 无
11.    """
12.    img0 = cv.imread('dataset/1.jpg', 0)
13.    X = np.mat(np.array(img0).flatten())
14.    for i in range(1, 38):
15.        file_path = 'dataset/' + str(i + 1) + '.jpg'
16.        img = cv.imread(file_path, 0)
17.        X = np.vstack((X, np.mat(np.array(img).flatten())))
18.    Y, B, m = PCA_solve(X, eta)
19.    x = Y @ B + m.T
20.    for i in range(38):
21.        img = x[i].reshape((48, 42))
22.        cv.imwrite('dataset_' + str(eta) + '/' + str(i + 1) + '.jpg', img)
23.
24.
25. def psnr(img1, img2):
26.     """
27.     计算峰值信噪比
28.     :param img1: 原图
29.     :param img2: 目的图像
30.     :return: 峰值信噪比
31.     """
32.     delta = img1 - img2
33.     mse = np.mean(np.multiply(delta, delta))
34.     ret = 20 * np.log10(255 / np.sqrt(mse))
35.     return ret
36.
37.
38. def cal_psnr():
39.     str2 = ''
40.     str3 = ''
41.     str4 = ''
42.     str5 = ''
43.     for i in range(0, 38):
44.         img1 = cv.imread('dataset/' + str(i + 1) + '.jpg', 0)

```

```
45.     img2 = cv.imread('dataset_0.95/' + str(i + 1) + '.jpg', 0)
46.     img3 = cv.imread('dataset_0.9/' + str(i + 1) + '.jpg', 0)
47.     img4 = cv.imread('dataset_0.85/' + str(i + 1) + '.jpg', 0)
48.     img5 = cv.imread('dataset_0.6/' + str(i + 1) + '.jpg', 0)
49.     ret2 = round(psnr(img1, img2), 4)
50.     ret3 = round(psnr(img1, img3), 4)
51.     ret4 = round(psnr(img1, img4), 4)
52.     ret5 = round(psnr(img1, img5), 4)
53.     str2 += str(i + 1) + '.jpg: ' + str(ret2) + '\n'
54.     str3 += str(i + 1) + '.jpg: ' + str(ret3) + '\n'
55.     str4 += str(i + 1) + '.jpg: ' + str(ret4) + '\n'
56.     str5 += str(i + 1) + '.jpg: ' + str(ret5) + '\n'
57.
58.     f2 = open('dataset_0.95/psnr.txt', 'w')
59.     f2.write(str2)
60.     f3 = open('dataset_0.9/psnr.txt', 'w')
61.     f3.write(str3)
62.     f4 = open('dataset_0.85/psnr.txt', 'w')
63.     f4.write(str4)
64.     f5 = open('dataset_0.6/psnr.txt', 'w')
65.     f5.write(str5)
66.
67.
68. if __name__ == '__main__':
69.     compression(0.6)
70.     cal_psnr()
```