

# 2022-2023 学年秋季学期《模式识别》作业二

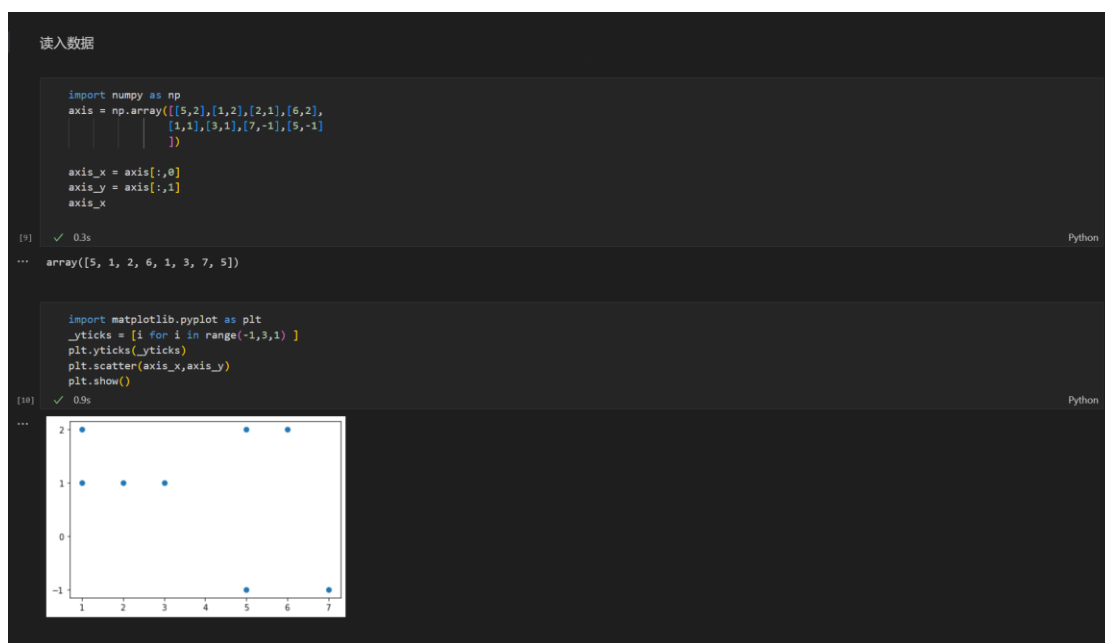
## （聚类分析、特征选择和特征提取）

教师：王玮

姓名：张奕桢 学号：2001111102

1. 分别采用顺序聚类，谱系聚类和 K-均值聚类算法将下列 8 个样本聚成两个类别，并用 Dunn 指数和 Davies-Bouldin 指数评估聚类效果。

$$x_1 = (5, 2)^T, x_2 = (1, 2)^T, x_3 = (2, 1)^T, x_4 = (6, 2)^T$$
$$x_5 = (1, 1)^T, x_6 = (3, 1)^T, x_7 = (7, -1)^T, x_8 = (5, -1)^T$$



顺序聚类是手算的

$x_1 = (5, 2)$  label = 0  
 $x_2 = (1, 2)$   $d(x_2, C_0) = 4 > \sqrt{5}$  label = 1  
 $x_3 = (2, 1)$   $d(x_3, C_0) = \sqrt{5}$  label = 1  
 $x_4 = (6, 2)$   $d(x_4, C_0) = \max(d(x_4, x_1), d(x_4, C_1)) = \max(\sqrt{10}, \sqrt{5}) = \sqrt{10}$   
 $d(x_4, x_3) = \max(\sqrt{20}, \sqrt{5}) = \sqrt{20}$  label = 0  
 $x_5 = (1, 1)$   $d(x_5, C_0) = \max(d(x_5, x_1), d(x_5, x_4)) = \max(\sqrt{10}, \sqrt{5}) = \sqrt{5}$   
 $d(x_5, C_1) = \max(d(x_5, x_3), d(x_5, x_5)) = \max(\sqrt{5}, \sqrt{5}) = \sqrt{5}$   
label = 1  
 $x_6 = (3, 1)$   $d(x_6, C_0) = \max(d(x_6, x_1), d(x_6, x_4)) = \max(\sqrt{10}, \sqrt{5}) = \sqrt{5}$   
 $d(x_6, C_1) = \max(d(x_6, x_3), d(x_6, x_5)) = \max(\sqrt{5}, \sqrt{5}) = \sqrt{5}$   
 $\sqrt{5} < \sqrt{10}$   
label = 1  
 $x_7 = (7, -1)$   $d(x_7, C_0) = \max(d(x_7, x_1), d(x_7, x_4)) = \max(\sqrt{10}, \sqrt{5}) = \sqrt{10}$   
 $d(x_7, C_1) = \max(d(x_7, x_3), d(x_7, x_5)) = \max(\sqrt{50}, \sqrt{5}) = \sqrt{50}$   
 $\sqrt{10} < \sqrt{50}$  label = 0  
 $x_8 = (5, -1)$   $d(x_8, C_0) = \max(d(x_8, x_1), d(x_8, x_4), d(x_8, x_7)) = \max(\sqrt{10}, \sqrt{2}) = \sqrt{10}$   
 $d(x_8, C_1) = \max(d(x_8, x_3), d(x_8, x_5), d(x_8, x_6)) = \max(\sqrt{5}, \sqrt{5}, \sqrt{5}) = \sqrt{5}$   
 $\sqrt{5} < \sqrt{10}$  label = 0

```

axis
✓ 0.4s Python
array([[ 5,  2],
       [ 1,  2],
       [ 2,  1],
       [ 6,  2],
       [ 1,  1],
       [ 3,  1],
       [ 7, -1],
       [ 5, -1]])

顺序聚类
+ 代码 + Markdown
label=[0,1,1,0,1,1,0,0]
label
280 ✓ 0.3s Python
[0, 1, 1, 0, 1, 1, 0, 0]

```

```

Kmeans
+ 代码 + Markdown
from sklearn.cluster import KMeans
clf = KMeans(n_clusters=2)
y_pred_kmeans = clf.fit_predict(axis)
# 输出完整Kmeans函数，包括很多省略参数
print("K均值模型:\n", clf)
# 输出聚类预测结果，20行数据，每个y_pred对应x一行或一个球员，聚成3类，类标为0、1、2
print("聚类结果:\n", y_pred_kmeans)

177 ✓ 0.6s Python
K均值模型:
KMeans(n_clusters=2)
聚类结果:
[0 1 1 0 1 1 0 0]

谱聚类
+ 代码 + Markdown
from sklearn.cluster import SpectralClustering as sc
clf = sc(gamma=1, n_clusters=2)
y_pred_sc = clf.fit_predict(axis)
# 输出完整Kmeans函数，包括很多省略参数
print("谱聚类模型:\n", clf)
# 输出聚类预测结果，20行数据，每个y_pred对应x一行或一个球员，聚成3类，类标为0、1、2
print("聚类结果:\n", y_pred_sc)

181 ✓ 0.5s Python
谱聚类模型:
SpectralClustering(gamma=1, n_clusters=2)
聚类结果:
[0 1 1 0 1 1 0 0]

```

```
计算距离矩阵

n=axis.shape
D = np.zeros([n,n])
for i in range(n):
    for j in range(i+1,n):
        d = axis[i,:]-axis[j,:]
        D[i,j] = np.sqrt(np.dot(d,d))
        D[j,i] = D[i,j]
D

array([[0., 4., 3.16227766, 1., 4.12310563,
        2.23606798, 3.60555128, 3., ],
       [4., 0., 1.41421356, 5., 1.,
        2.23606798, 6.70820393, 5., ],
       [3.16227766, 1.41421356, 0., 4.12310563, 1.,
        1., 5.38516481, 3.60555128],
       [1., 5., 4.12310563, 0., 5.09901951,
        3.16227766, 3.16227766, 3.16227766],
       [4.12310563, 1., 1., 5.09901951, 0.,
        2., 6.32455532, 4.47213595],
       [2.23606798, 2.23606798, 1., 3.16227766, 2.,
        0., 4.47213595, 2.82842712],
       [3.60555128, 6.70820393, 5.38516481, 3.16227766, 6.32455532,
        4.47213595, 0., 2.],
       [3., 5., 3.60555128, 3.16227766, 4.47213595,
        2.82842712, 2., 0.]])

label0=[0,3,6,7]
label1=[1,2,4,5]
min_d = D[0][1]
for i in label0:
    for j in label1:
        min_d = min(D[i][j],min_d)
max_d0=0
for i in label0:
    for j in label1:
        if(i!=j):
            max_d0=max(D[i][j],max_d0)
max_d1=0
for i in label0:
    for j in label1:
        if(i!=j):
            max_d1=max(D[i][j],max_d1)
print("min_d=",min_d)
print("max_d0=",max_d0)
print("min_d1=",max_d1)
max_d = max(max_d0,max_d1)
J_Dunn = min_d/max_d
print("J_Dunn",J_Dunn)

min_d= 2.23606797749979
max_d0= 3.605551275463989
min_d1= 3.605551275463989
J_Dunn 0.6201736729460423

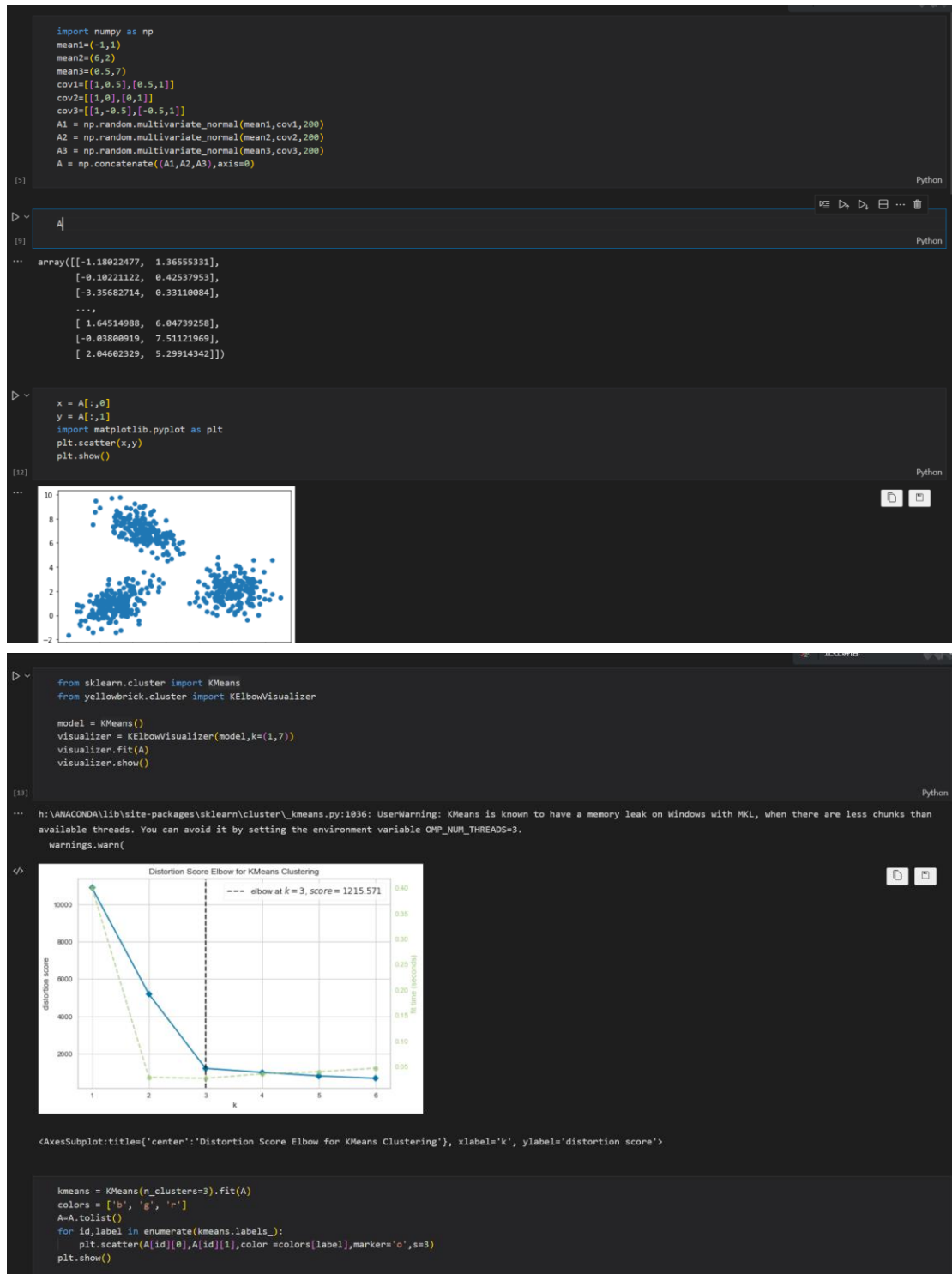
from sklearn import metrics
from pandas.core.frame import DataFrame
df = DataFrame(axis)
score = metrics.davies_bouldin_score(df, label)
print("Davies-Bouldin-score: ", score)

Davies-Bouldin-score: 0.6331473961991908
```

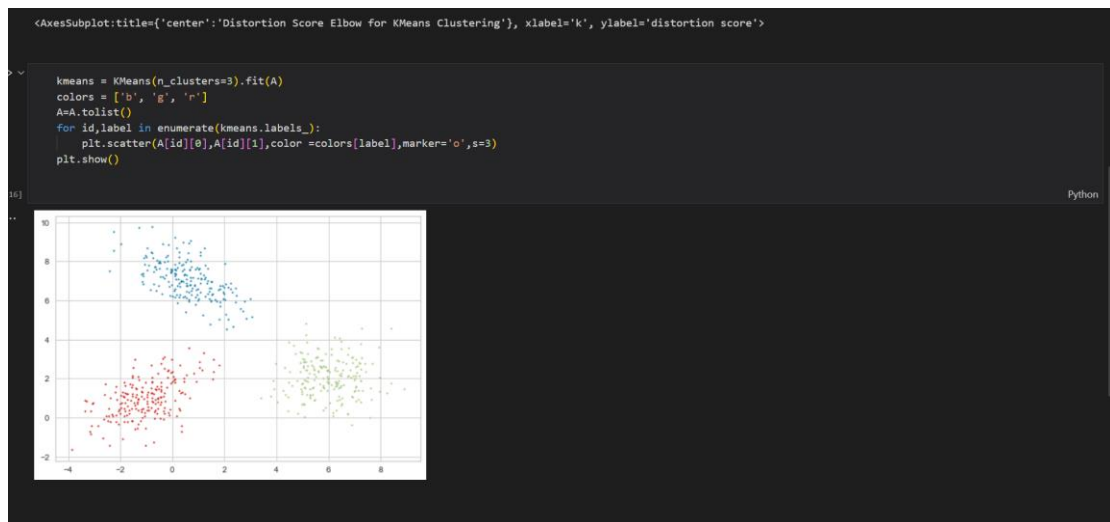
由于三种方法聚类结果是一样的，因此算出的两种评分标准是一样的

2. 根据下列 3 组高斯分布参数分别生成 200 个 2 维随机样本，使用 K-均值算法将所有样本聚成若干类别，采用类内距离准则确定合理聚类数。

$$\mu_1 = (-1, -1)^T, \mu_2 = (6, 2)^T, \mu_3 = (0.5, 7)^T$$
$$\Sigma_1 = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \Sigma_3 = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$



其中 `distortion_score` 是没有求类内距离平均之前的值。（和平均之后是等价的）



3. 有两类样本集合：

- 类别 1 中有 5 个样本  $c_1 = [(1,2), (2,3), (3,3), (4,5), (5,5)]$
- 类别 2 中有 6 个样本  $c_2 = [(1,0), (2,1), (3,1), (3,2), (5,3), (6,5)]$

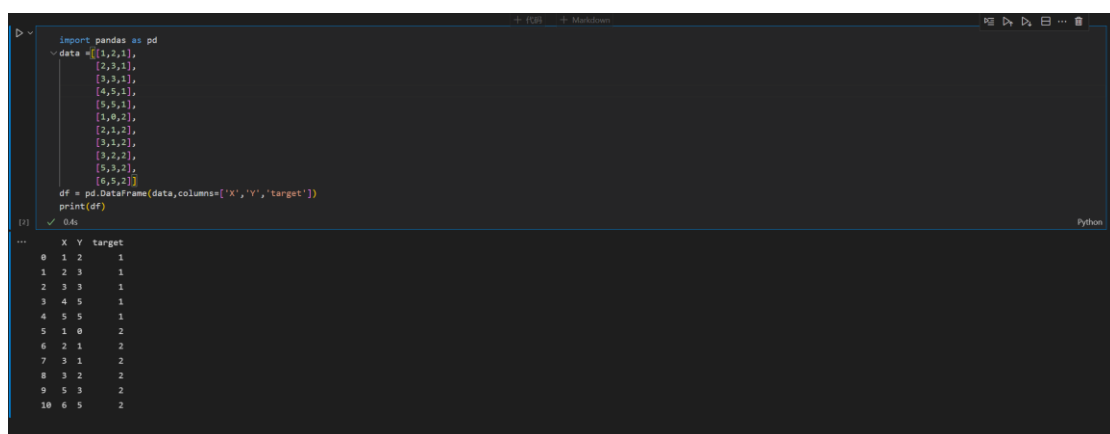
(1) 使用主成分分析将特征降为 1 维，在二维空间中画出样本和投影坐标轴。

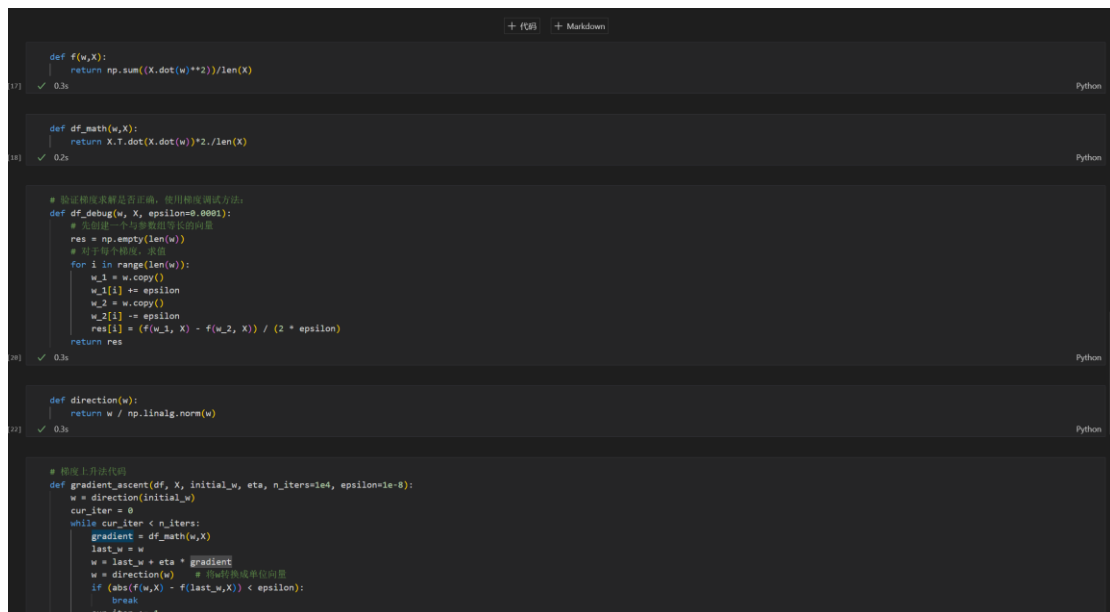
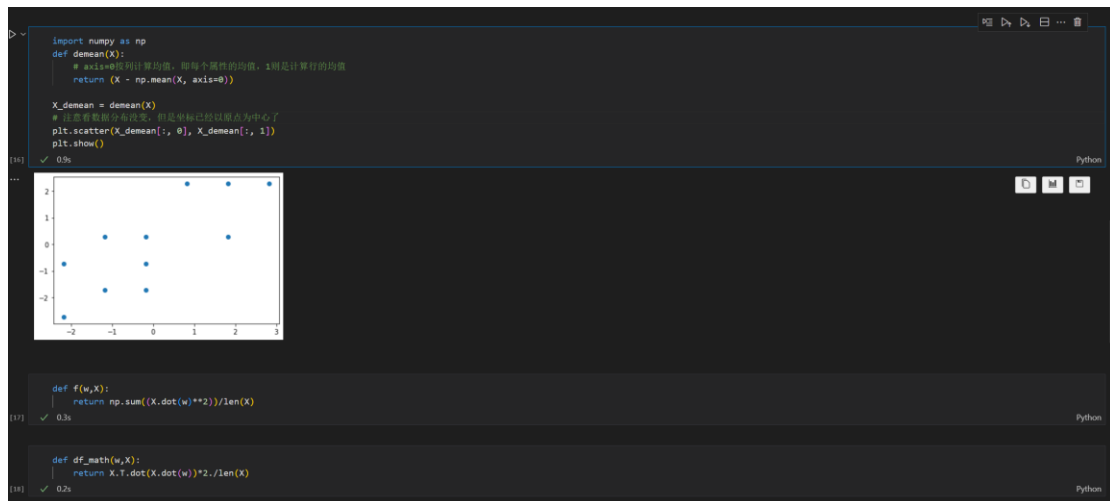
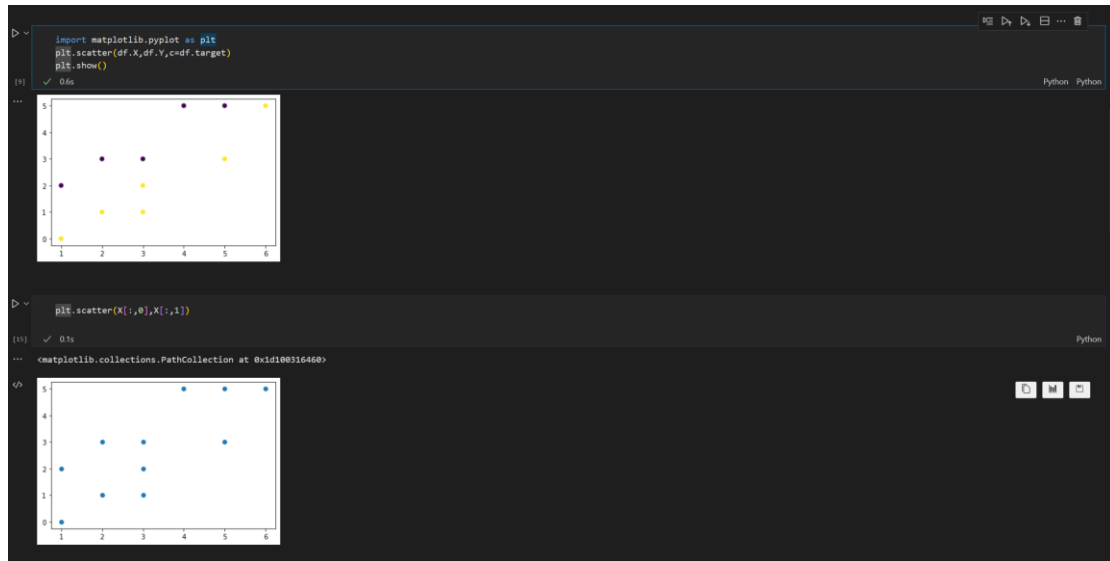
(2) 使用 Fisher 判别分析将特征降为 1 维，在二维空间中画出样本和投影坐标轴。

(3) 在使用 Fisher 判别分析降维后的 1 维空间中构建线性分类器区分两类样本。

(4) 如果在主成分分析降维后的 1 维空间中构建线性分类器，效果是否理想？为什么？

(1)







```

import pandas as pd
import numpy as np
data = [[1,2,1],
        [2,3,1],
        [3,3,1],
        [4,5,1],
        [5,5,1],
        [1,0,2],
        [2,1,2],
        [3,1,2],
        [3,2,2],
        [5,3,2],
        [6,5,2]]
df = pd.DataFrame(data,columns=['X','Y','target'])
df_X=df[['X','Y']]
X=df_X.to_numpy()
df_Y=df[['target']]
Y=df_Y.to_numpy()
c_1=X[0:5][:]
c_2=X[5:-1][:]

```

[17]

✓ 0.3s

Python

```

def cal_cov_avg(samples):
    """
    给定一个类别的数据，计算协方差矩阵和平均向量
    :param samples:
    :return:
    """
    u1 = np.mean(samples, axis=0)
    cov_m = np.zeros((samples.shape[1], samples.shape[1]))
    for s in samples:
        t = s - u1
        cov_m += t * t.reshape(2, 1)
    return cov_m, u1

def fisher(c_1, c_2):
    """
    fisher算法实现(请参考上面推导出来的公式，那个才是精华部分)
    :param c_1:
    :param c_2:
    :return:
    """
    cov_1, u1 = cal_cov_avg(c_1)
    cov_2, u2 = cal_cov_avg(c_2)
    s_w = cov_1 + cov_2
    u, s, v = np.linalg.svd(s_w) # 奇异值分解
    s_w_inv = np.dot(np.dot(v.T, np.linalg.inv(np.diag(s))), u.T)
    return np.dot(s_w_inv, u1 - u2)

```

[18]

✓ 0.3s

Python

```

def judge(sample, w, c_1, c_2):
    """
    true 属于1
    false 属于2
    :param sample:
    """

```



```

def judge_center_1(c_1, c_2):
    :param center_2:
    :return:
    """
    u1 = np.mean(c_1, axis=0)
    u2 = np.mean(c_2, axis=0)
    center_1 = np.dot(w.T, u1)
    center_2 = np.dot(w.T, u2)
    pos = np.dot(w.T, sample)
    return abs(pos - center_1) < abs(pos - center_2)

w = fisher(c_1, c_2) # 调用函数，得到参数w
out = judge(c_1[1], w, c_1, c_2) # 判断所属类别
print(out)

[19] ✓ 0.2s Python
... True

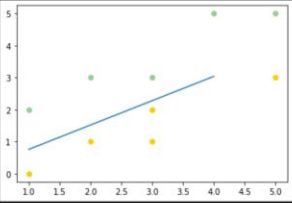
import matplotlib.pyplot as plt

plt.scatter(c_1[:, 0], c_1[:, 1], c='#99CC99')
plt.scatter(c_2[:, 0], c_2[:, 1], c='#FFCC00')
line_x = np.arange(min(np.min(c_1[:, 0]), np.min(c_2[:, 0])),
                    max(np.max(c_1[:, 0]), np.max(c_2[:, 0])),
                    step=1)

line_y = - (w[0] * line_x) / w[1]
plt.plot(line_x, line_y)
plt.show()

[20] ✓ 0.1s Python
...

```



```

第3问
Markdown

```

```

from sklearn import discriminant_analysis
#定义Fisher分类器对象fisher_clf
fisher_clf = discriminant_analysis.LinearDiscriminantAnalysis()
#调用该对象的训练方法
fisher_clf.fit(X,Y)

[4] ✓ 0.8s Python
... h:\ANACONDA\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape

```

```

第3问
Markdown

```

```

>
from sklearn import discriminant_analysis
#定义Fisher分类器对象fisher_clf
fisher_clf = discriminant_analysis.LinearDiscriminantAnalysis()
#调用该对象的训练方法
fisher_clf.fit(X,Y)

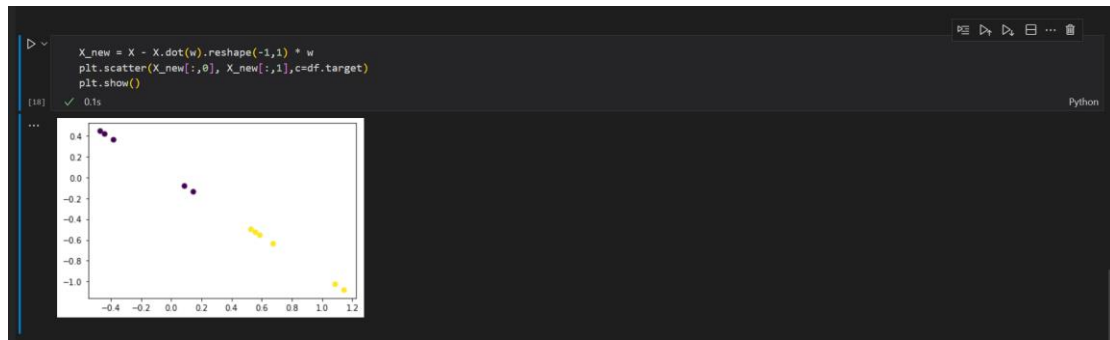
[4] ✓ 0.8s Python
... h:\ANACONDA\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape
of y to (n_samples, ), for example using ravel().
y = column_or_1d(y, warn=True)

LinearDiscriminantAnalysis()

y_pred=fisher_clf.predict(X)
print("测试数据集的正确标签为:",Y)
print("测试数据集的预测标签为:",y_pred)
from sklearn.metrics import accuracy_score
testing_accuracy_score(Y, y_pred)*100
print("Fisher线性分类器测试准确率: {:.2f}%".format(testing_acc))

[5] ✓ 0.3s Python
...
测试数据集的正确标签为: [1]
[1]
[1]
[1]
[1]
[2]
[2]
[2]
[2]
[2]
[2]
测试数据集的预测标签为: [1 1 1 1 2 2 2 2 2]
Fisher线性分类器测试准确率: 100.00%

```



比较理想因为各个点之间没有相互交叉，特征比较明显