

第四次上机作业

18029100040

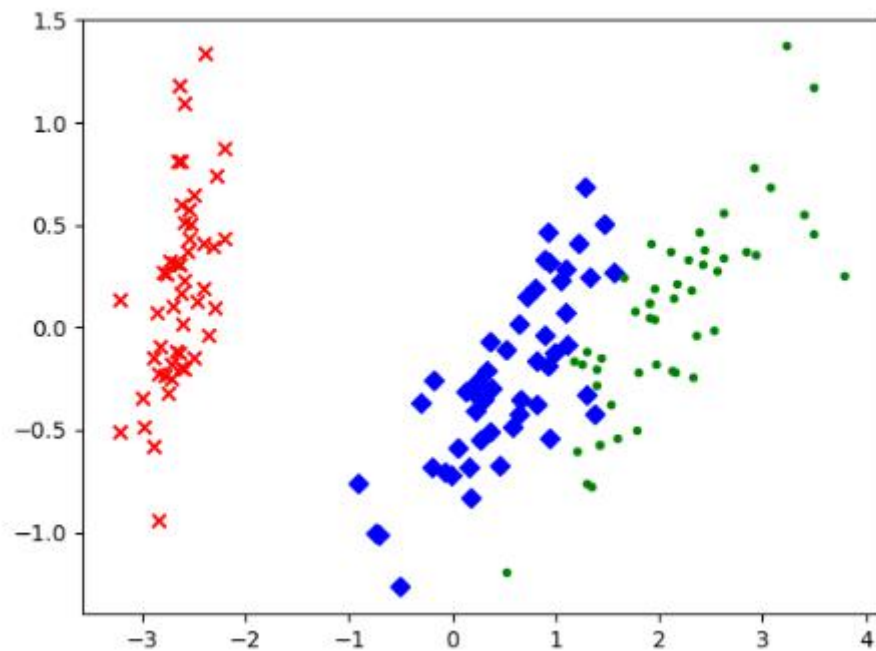
吴程锴

一、作业 26：PCA 实现高维数据可视化：

1.1 代码

```
1. import matplotlib.pyplot as plt
2. from sklearn.decomposition import PCA
3. from sklearn.datasets import load_iris
4. data = load_iris()
5. y = data.target
6. x = data.data
7. pca = PCA(n_components = 2)
8. reduced_x = pca.fit_transform(x)
9. red_x, red_y = [], []
10. blue_x, blue_y = [], []
11. green_x, green_y = [], []
12. for i in range(len(reduced_x)):
13.     if y[i] == 0:
14.         red_x.append(reduced_x[i][0])
15.         red_y.append(reduced_x[i][1])
16.     elif y[i] == 1:
17.         blue_x.append(reduced_x[i][0])
18.         blue_y.append(reduced_x[i][1])
19.     else:
20.         green_x.append(reduced_x[i][0])
21.         green_y.append(reduced_x[i][1])
22. plt.scatter(red_x, red_y, c = 'r', marker = 'x')
23. plt.scatter(blue_x, blue_y, c = 'b', marker = 'D')
24. plt.scatter(green_x, green_y, c = 'g', marker = '.')
25. plt.show()
```

1.2 结果



二、作业 27：降维之 NMF

2.1 代码

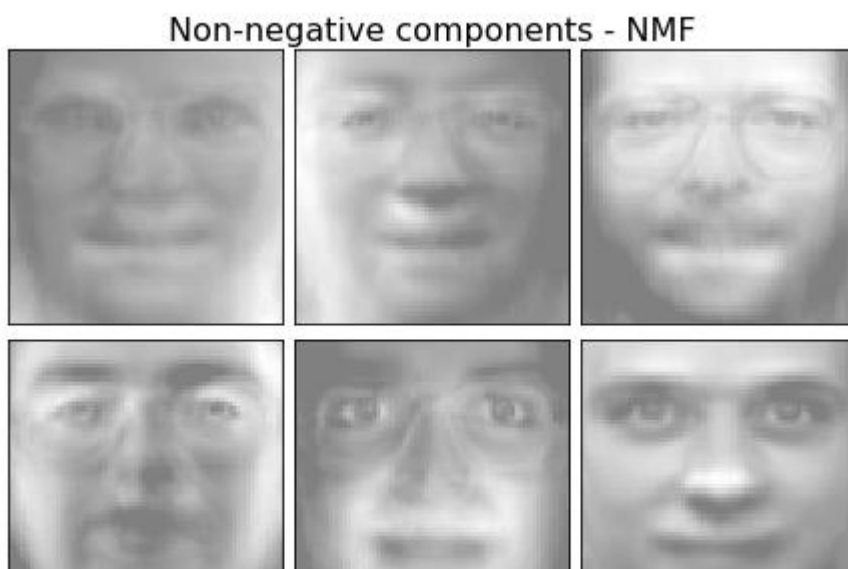
```
1. import matplotlib.pyplot as plt
2. from sklearn import decomposition
3. from sklearn.datasets import fetch_olivetti_faces
4. from numpy.random import RandomState
5.
6.
7.
8. n_row, n_col = 2, 3
9.
10. n_components = n_row = n_col
11.
12. image_shape = (64, 64)
13. dataset = fetch_olivetti_faces(shuffle=True, random_state=RandomState(0))
14. faces = dataset.data
15. def plot_gallery(title, images, n_col=n_col, n_row=n_row):
16.     plt.figure(figsize=(2. * n_col, 2.26 * n_row))
17.     plt.suptitle(title, size=16)
18.     for i, comp in enumerate(images):
```

```

19.     plt.subplot(n_row, n_col, i + 1)
20.     vmax = max(comp.max(), -comp.min())
21.     plt.imshow(comp.reshape(image_shape), cmap=plt.cm.gray, interpolation='nearest', vmin=-vmax,
22.                 vmax=vmax)
23.     plt.xticks(())
24.     plt.yticks(())
25.     plt.subplots_adjust(0.01, 0.05, 0.99, 0.93, 0.04, 0.)
26.
27.
28. # 创建特征提取的对象 NMF，使用 PCA 作为对
29. estimators = [('Eigenfaces - PCA using randomized SVD', decomposition.PCA(n_components=6, whiten=True)),
30.               ('Non-negative components - NMF', decomposition.NMF(n_components=6, init='nndsvda', tol=5e-3))]
31. # NMF 和 PCA 实例，将它们放在一个列表中
32.
33. # 降维后数据点的可视化
34. for name, estimator in estimators:
35.     estimator.fit(faces)
36.     components_ = estimator.components_
37.     plot_gallery(name, components_[:n_components])
38.     plt.show() # 可视化

```

2.2 结果

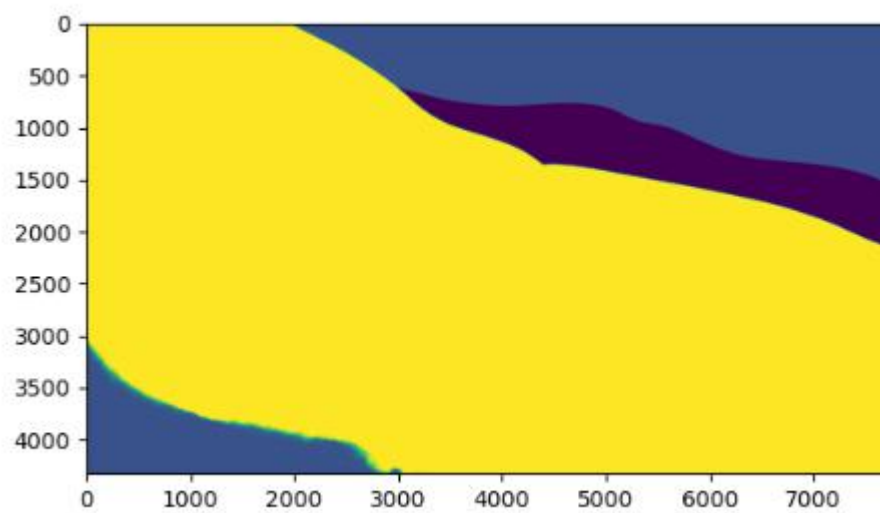


三、作业 28：图像分割

3.1 代码

```
1. import numpy as np
2. from PIL import Image #加载 PIL 包，用于加载创建图片
3. from sklearn.cluster import KMeans #加载 Kmeans 算法
4. import matplotlib.pyplot as plt #绘制图像
5.
6. def loadData(filePath):
7.     f = open(filePath, 'rb') #以二进制形式打开文件
8.     data = []
9.     img = Image.open(f) #以列表的形式返回图片像素值
10.    m, n = img.size #获取图片的大小
11.    for i in range(m): #将每个像素点的 RGB 颜色处理到 0-1
12.        for j in range(n):
13.            x,y,z = img.getpixel((i,j))
14.            data.append([x/256.0, y/256.0, z/256.0]) #范围内并存入 data
15.    f.close()
16.    return np.mat(data), m, n #以矩阵的形式返回 data，以及图片大小
17.
18. imgData,row,col = loadData('bull.jpg')
19. #加载 Kmeans 聚类算法
20. km = KMeans(n_clusters= 3) #其中 n_clusters 属性指定了聚类中心的个数为 3
21.
22. #聚类获取每个像素所属的类别
23. label = km.fit_predict(imgData)
24. label = label.reshape([row, col])
25. #创建一张新的灰度图保存聚类后的结果
26. pic_new = Image.new('L', (row, col))
27.
28. #根据所属类别向图片中添加灰度值
29. # 最终利用聚类中心点的 RGB 值替换原图中每一个像素点的值，便得到了最终的分割后的图片
30. for i in range(row):
31.     for j in range(col):
32.         pic_new.putpixel((i, j), int(256 / (label[i][j] + 1)))
33.
34. #以 JPEG 格式保存图片
35. pic_new.save("bull_r.jpg", "JPEG")
36. plt.imshow(pic_new)
37. plt.show()
```

3.2 结果



四、作业 29：人体运动状态信息评级

4.1 代码

```
1. import pandas as pd
2. import numpy as np
3.
4. from sklearn.preprocessing import Imputer
5. from sklearn.cross_validation import train_test_split
6. from sklearn.metrics import classification_report
7.
8. from sklearn.neighbors import KNeighborsClassifier
9. from sklearn.tree import DecisionTreeClassifier
10. from sklearn.naive_bayes import GaussianNB
```

```

11.
12. def load_datasets(feature_paths, label_paths):
13.     feature = np.ndarray(shape=(0,41))
14.     label = np.ndarray(shape=(0,1))
15.     for file in feature_paths:
16.         df = pd.read_table(file, delimiter=',', na_values='?', header=None)
17.         imp = Imputer(missing_values='NaN', strategy='mean', axis=0)
18.         imp.fit(df)
19.         df = imp.transform(df)
20.         feature = np.concatenate((feature, df))
21.
22.     for file in label_paths:
23.         df = pd.read_table(file, header=None)
24.         label = np.concatenate((label, df))
25.
26.     label = np.ravel(label)
27.     return feature, label
28.
29. if __name__ == '__main__':
30.     featurePaths = ['A/A.feature', 'B/B.feature', 'C/C.feature', 'D/D.feature', 'E/E.feature']
31.     labelPaths = ['A/A.label', 'B/B.label', 'C/C.label', 'D/D.label', 'E/E.label']
32.     x_train, y_train = load_datasets(featurePaths[:4], labelPaths[:4])
33.     x_test, y_test = load_datasets(featurePaths[4:], labelPaths[4:])
34.     x_train, x_, y_train, y_ = train_test_split(x_train, y_train, test_size = 0.0)
35.
36.     print('Start training knn')
37.     knn = KNeighborsClassifier().fit(x_train, y_train)
38.     print('Training done')
39.     answer_knn = knn.predict(x_test)
40.     print('Prediction done')
41.
42.     print('Start training DT')
43.     dt = DecisionTreeClassifier().fit(x_train, y_train)
44.     print('Training done')
45.     answer_dt = dt.predict(x_test)
46.     print('Prediction done')
47.
48.     print('Start training Bayes')
49.     gnb = GaussianNB().fit(x_train, y_train)
50.     print('Training done')
51.     answer_gnb = gnb.predict(x_test)
52.     print('Prediction done')
53.
54.     print('\n\nThe classification report for knn:')

```

```

55.     print(classification_report(y_test, answer_knn))
56.     print('\n\nThe classification report for DT:')
57.     print(classification_report(y_test, answer_dt))
58.     print('\n\nThe classification report for Bayes:')
59.         print(classification_report(y_test, answer_gnb))

```

4.2 结果

	precision	recall	f1-score	support
0.0	0.55	0.83	0.66	102341
1.0	0.79	0.96	0.87	23699
2.0	0.91	0.84	0.87	26864
3.0	0.93	0.73	0.82	22132
4.0	0.63	0.95	0.76	32033
5.0	0.73	0.50	0.59	24646
6.0	0.05	0.01	0.02	24577
7.0	0.32	0.14	0.20	26271
12.0	0.60	0.66	0.63	14281
13.0	0.67	0.48	0.56	12727
16.0	0.57	0.07	0.13	24445
17.0	0.86	0.85	0.86	33034
24.0	0.37	0.30	0.33	7733

五、作业 30：上证指数涨跌预测实例

5.1 代码

```

1. import pandas as pd
2. # : 支持高级大量的维度数组与矩阵运算，此外也针对数组运算提供大量的数学函数库
3. import numpy as np
4. # sklearn 下 svm: SVM 算法
5. from sklearn import svm
6. # sklearn 下 cross_validation: 交叉验证
7. from sklearn import model_selection
8.
9. # parse_dates=第 0 列解析为日期， index_col= 用作行索引的列编号
10. data = pd.read_csv(r'000777.csv', encoding='gbk', parse_dates=[0],
11.                    index_col=0)
12. # DataFrame.sort_index(axis=0 (按 0 列排), ascending=True (升序),
13. # inplace=False (排序后是否覆盖原数据)) data 按照时间升序排列
14. data.sort_index(0, ascending=True, inplace=True)
15.
16. # 选取 5 列数据作为特征: 收盘价 最高价 最低价 开盘价 成交量
17. # dayfeature: 选取 150 天的数据
18. # featurenum: 选取的 5 个特征*天数

```

```

19. # x: 记录 150 天的 5 个特征值 y: 记录涨或者跌
20. dayfeature = 150
21. featurenum = 5 * dayfeature
22. # data.shape[0]-dayfeature 意思是因为我们要用 150 天数据做训练,
23. # 对于条目为 200 条的数据, 只有 50 条数据是有前 150 天的数据来训练的,
24. # 所以测试集的大小就是 200-150, 对于每一条数据, 他的特征是前 150 天的所有特征数据,
25. # 即 150*5, +1 是将当天的开盘价引入作为一条特征数据
26. x = np.zeros((data.shape[0] - dayfeature, featurenum + 1))
27. y = np.zeros((data.shape[0] - dayfeature))
28.
29. for i in range(0, data.shape[0] - dayfeature):
30.     # /将数据中的“收盘价”“最高价”“开盘价”“成交量”存入 x 数组中
31.     # u:unicode 编码 reshape:转换成 1 行, featurenum 列
32.     x[i, 0:featurenum] = np.array(data[i:i + dayfeature] \
33.                                     [[u'收盘价', u'最高价',
34.                                         u'最低价', u'开盘价', u'成交量
35.                                     ]]).reshape((1, featurenum))
36.     x[i, featurenum] = data.ix[i + dayfeature][u'开盘价']
37.     # 最后一列记录当日的开盘价 ix :索引
38.     for i in range(0, data.shape[0] - dayfeature):
39.         if data.ix[i + dayfeature][u'收盘价'] >= data.ix[i + dayfeature][u'开盘价']:
40.             y[i] = 1
41.         else:
42.             y[i] = 0
43.         # 如果当天收盘价高于开盘价, y[i]=1 代表涨, 0 代表跌
44. # 创建 SVM 并进行交叉验证
45. clf = svm.SVC(kernel='rbf')
46. # 调用 svm 函数, 并设置 kernel 参数, 默认是 rbf, 其它: 'linear' 'poly' 'sigmoid'
47. result = []
48. for i in range(5):
49.     # x 和 y 的验证集和测试集, 切分 80 - 20 % 的测试集
50.     x_train, x_test, y_train, y_test = \
51.         model_selection.train_test_split(x, y, test_size=0.2)
52.     # 训练数据进行训练
53.     clf.fit(x_train, y_train)
54.     # 将预测数据和测试集的验证数据比对
55.     result.append(np.mean(y_test == clf.predict(x_test)))
56. print("svm classifier accuracy:")
57.     print(result)

```

5.2 结果

```

svm classifier accuracy:
[0.5635179153094463, 0.5754614549402823, 0.5266015200868621, 0.5450597176981542, 0.5407166123778502]

```


六、作业 31：线性回归+房价与房屋尺寸关系的线性拟合

6.1 代码

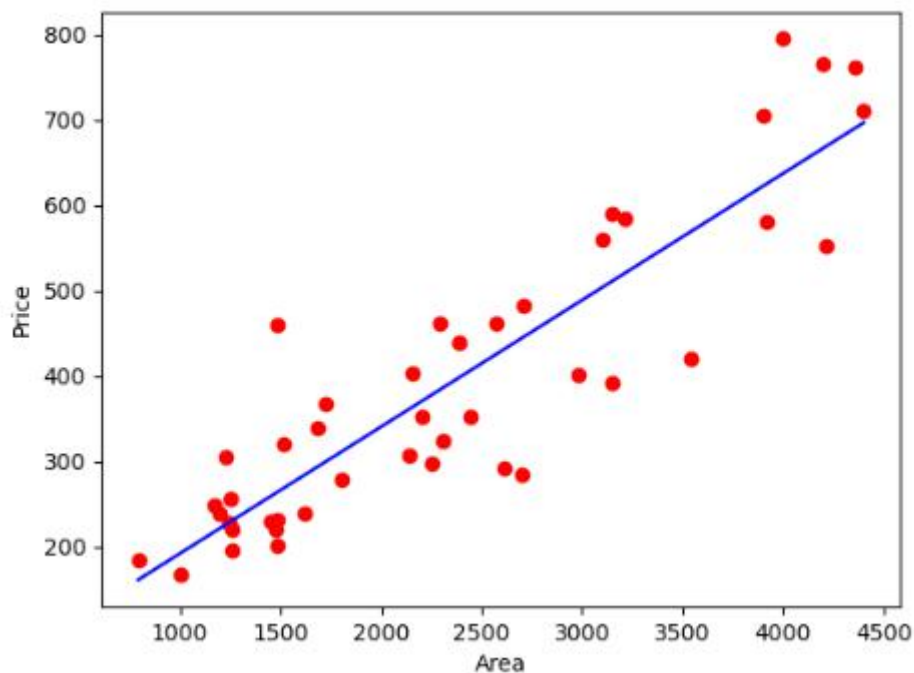
```
1.      #表示matplotlib的pyplot子库，它提供了和matlab类似的绘图API
2. import matplotlib.pyplot as plt
3. #表示可以调用sklearn中的linear_model模块进行线性回归。
4. from sklearn import linear_model
5. import numpy as np
6. #建立datasets_X和datasets_Y用来存储数据中的房屋尺寸和房屋成交价格。
7. datasets_X=[]
8. datasets_Y=[]
9. fr=open('prices.txt','r')
10. #一次读取整个文件。
11. lines=fr.readlines()
12. #逐行进行操作，循环遍历所有数据
13. for line in lines:
14.     #去除数据文件中的逗号
15.     items=line.strip().split(',')
16.     #将读取的数据转换为int型，并分别写入datasets_X和datasets_Y。
17.     datasets_X.append(int(items[0]))
18.     datasets_Y.append(int(items[1]))
19. #求得datasets_X的长度，即为数据的总数。
20. length=len(datasets_X)
21. #将datasets_X转化为数组，并变为二维，以符合线性回归拟合函数输入参数要求
22. datasets_X=np.array(datasets_X).reshape([length,1])
23. #将datasets_Y转化为数组
24. datasets_Y=np.array(datasets_Y)
25.
26. minX=min(datasets_X)
27. maxX=max(datasets_X)
28. #以数据datasets_X的最大值和最小值为范围，建立等差数列，方便后续画图。
29. X=np.arange(minX,maxX).reshape([-1,1])
30.
31. linear=linear_model.LinearRegression()
32. linear.fit(datasets_X,datasets_Y)#调用线性回归模块，建立回归方程，拟合数据
33. #查看回归方程系数
34. print('Coefficients:',linear.coef_)
35. #查看回归方程截距
36. print('intercept',linear.intercept_)
37.
38. #3.可视化处理
39. #scatter函数用于绘制数据点，这里表示用红色绘制数据点；
40. plt.scatter(datasets_X,datasets_Y,color='red')
```

```

41. #plot 函数用来绘制直线，这 里表示用蓝色绘制回归线；
42. #xlabel 和 ylabel 用来指定横纵坐标的名称
43. plt.plot(X,linear.predict(X),color='blue')
44. plt.xlabel('Area')
45. plt.ylabel('Price')
46. plt.show()

```

6.2 结果



七、作业 32：多项式回归+房价与房屋尺寸的非线性拟合

7.1 代码

```

1. #房价与房屋尺寸关系的非线性拟合（多项式回归）
2. import numpy as np
3. import matplotlib.pyplot as plt
4. from sklearn import linear_model
5. from sklearn.preprocessing import PolynomialFeatures as PF
6.
7. plt.rcParams['font.sans-serif']='SimHei'
8. plt.rcParams['axes.unicode_minus']=False
9. data_x=[] #设房屋的尺寸（面积）为 data_x
10. data_y=[] #房价为 data_y

```

```
11. f=open('prices.txt','r')
12.
13. lines=f.readlines()
14. for line in lines:
15.     items=line.strip().split(',')
16.     data_x.append(int(items[0]))
17.     data_y.append(int(items[1]))
18. length=len(data_x)
19. data_x=np.array(data_x).reshape([length,1])
20. data_y=np.array(data_y)
21. minX=min(data_x)
22. maxX=max(data_x)
23. x=np.arange(minX,maxX).reshape([-1,1])
24.
25. poly_reg=PF(degree=2) #degree=2 表示建立 data_x 的二 次多项式特征 x_poly
26. x_poly=poly_reg.fit_transform(data_x)
27. linear=linear_model.LinearRegression()
28. linear.fit(x_poly,data_y) #拟合 x, y
29.
30. plt.scatter(data_x,data_y,color='red')
31. plt.plot(x,linear.predict(poly_reg.fit_transform(x)),color='blue')
32. plt.xlabel('Area') #x 轴标签
33. plt.ylabel('Price')
34. plt.title('房价与房屋尺寸关系的线性关系图')
35.     plt.show()
```

7.2 结果

