

以下数据集均可通过导入sklearn.datasets获得

- 实现鸢尾花数据集kmeans聚类并用散点图展示；使用PCA将鸢尾花数据集降到2维并用散点图展示

```
1  #coding=utf-8
2  import matplotlib.pyplot as plt
3  from sklearn.decomposition import PCA
4  from sklearn.datasets import load_iris
5  from mpl_toolkits.mplot3d import Axes3D
6  from sklearn.cluster import KMeans
7  import numpy as np
8
9
10
11  ##计算欧式距离
12  def distEuclid(x, y):
13      return np.sqrt(np.sum((x - y) ** 2))
14
15
16  ##随机产生n个dim维度的数据 (这里为了展示结果 dim取2或者3)
17  def genDataset(n, dim):
18      data = []
19      while len(data) < n:
20          p = np.around(np.random.rand(dim) * size, decimals=2)
21          data.append(p)
22      return data
23
24
25  ## 初始化簇中心点 一开始随机从样本中选择k个 当做各类簇的中心
26  def initCentroid(data, k):
27      num, dim = data.shape
28      centpoint = np.zeros((k, dim))
29      l = [x for x in range(num)]
30      np.random.shuffle(l)
31      for i in range(k):
32          index = int(l[i])
33          centpoint[i] = data[index]
34      return centpoint
35
36
37  ##进行KMeans分类
38  def KMeans(data, k):
39      ##样本个数
40      num = np.shape(data)[0]
41
42      ##记录各样本 簇信息 0:属于哪个簇 1:距离该簇中心点距离
43      cluster = np.zeros((num, 2))
44      cluster[:, 0] = -1
45
46      ##记录是否有样本改变簇分类
47      change = True
48      ##初始化各簇中心点
49      cp = initCentroid(data, k)
50
```

```

51     while change:
52         change = False
53
54         ##遍历每一个样本
55         for i in range(num):
56             minDist = 9999.9
57             minIndex = -1
58
59             ##计算该样本距离每一个簇中心点的距离 找到距离最近的中心点
60             for j in range(k):
61                 dis = distEuclid(cp[j], data[i])
62                 if dis < minDist:
63                     minDist = dis
64                     minIndex = j
65
66             ##如果找到的簇中心点非当前簇 则改变该样本的簇分类
67             if cluster[i, 0] != minIndex:
68                 change = True
69                 cluster[i, :] = minIndex, minDist
70
71             ## 根据样本重新分类 计算新的簇中心点
72             for j in range(k):
73                 pointincluster = data[[x for x in range(num) if cluster[x, 0] ==
j]]
74                 cp[j] = np.mean(pointincluster, axis=0)
75
76     print("finish!")
77     return cp, cluster
78
79
80     ##展示结果 各类簇使用不同的颜色 中心点使用X表示
81     def Show(data, k, cp, cluster):
82         num, dim = data.shape
83         color = ['b', 'r', 'g', 'c', 'y', 'm', 'k']
84         ##二维图
85         if dim == 2:
86             for i in range(num):
87                 mark = int(cluster[i, 0])
88                 plt.plot(data[i, 0], data[i, 1], color[mark] + 'o')
89
90             for i in range(k):
91                 plt.plot(cp[i, 0], cp[i, 1], color[i] + 'x')
92         ##三维图
93         elif dim == 3:
94             ax = plt.subplot(122, projection='3d')
95             for i in range(num):
96                 mark = int(cluster[i, 0])
97                 ax.scatter(data[i, 0], data[i, 1], data[i, 2], c=color[mark])
98             ax.set_title('k-means result')
99             for i in range(k):
100                 ax.scatter(cp[i, 0], cp[i, 1], cp[i, 2], c=color[i], marker='x')
101
102     plt.show()
103     if __name__ == "__main__":
104         data = load_iris()#以字典形式加载鸢尾花数据集
105         y = data.target #使用y表示数据集中的标签
106         x = data.data #使用x表示数据集中的属性数据
107         #使用PCA 算法, 设置降维后主成分数目为 2

```

```

108     #print(x, '\n', y)
109     #print(x)
110     #print(type(x))
111     size = 20 ##取值范围
112     pca = PCA(n_components=3)
113     #对原始数据进行降维，保存在 reduced_X 中
114     reduced_X = pca.fit_transform(x)
115     # print('降维后的数据为:\n', reduced_X)#降维后的数据
116     print('各主成分方差解释度为: ', pca.explained_variance_ratio_)#方差解释度
117     print('主成分对应的载荷矩阵为', pca.components_)
118     red_x, red_y, red_z = [], [], []
119     blue_x, blue_y, blue_z = [], [], []
120     green_x, green_y, green_z = [], [], []
121
122     for i in range(len(reduced_X)):
123         #标签为0时，3维标签数据保存到列表red_x, red_y, red_z中
124         if y[i] == 0:
125             red_x.append(reduced_X[i][0])
126             red_y.append(reduced_X[i][1])
127             red_z.append(reduced_X[i][2])#
128         elif y[i] == 1:
129             blue_x.append(reduced_X[i][0])
130             blue_y.append(reduced_X[i][1])
131             blue_z.append(reduced_X[i][2])#
132         else:
133             green_x.append(reduced_X[i][0])
134             green_y.append(reduced_X[i][1])
135             green_z.append(reduced_X[i][2])#
136     X = reduced_X[:, :3] ##表示我们取特征空间中的3个维度
137     #print(X.shape)
138     #print(X)
139     #print(type(X))
140     num = 50 ##点个数
141     k=3 ##分类个数
142     data = X
143     cp, cluster = KMeans(data, k)
144     ax=plt.figure().add_subplot(121, projection='3d')
145     ax.scatter(red_x, red_y, red_z, c='r', marker='o')
146     ax.scatter(blue_x, blue_y, blue_z, c='b', marker='o')
147     ax.scatter(green_x, green_y, green_z, c='g', marker='o')#散点图中用s，其余图
    用markersize可调节散点的大小
148     ax.set_title('PCA result')
149     Show(data, k, cp, cluster)

```

```

1  各主成分方差解释度为: [0.92461872 0.05306648 0.01710261]
2  主成分对应的载荷矩阵为 [[ 0.36138659 -0.08452251 0.85667061 0.3582892 ]
3  [ 0.65658877 0.73016143 -0.17337266 -0.07548102]
4  [-0.58202985 0.59791083 0.07623608 0.54583143]]
5  finish!

```

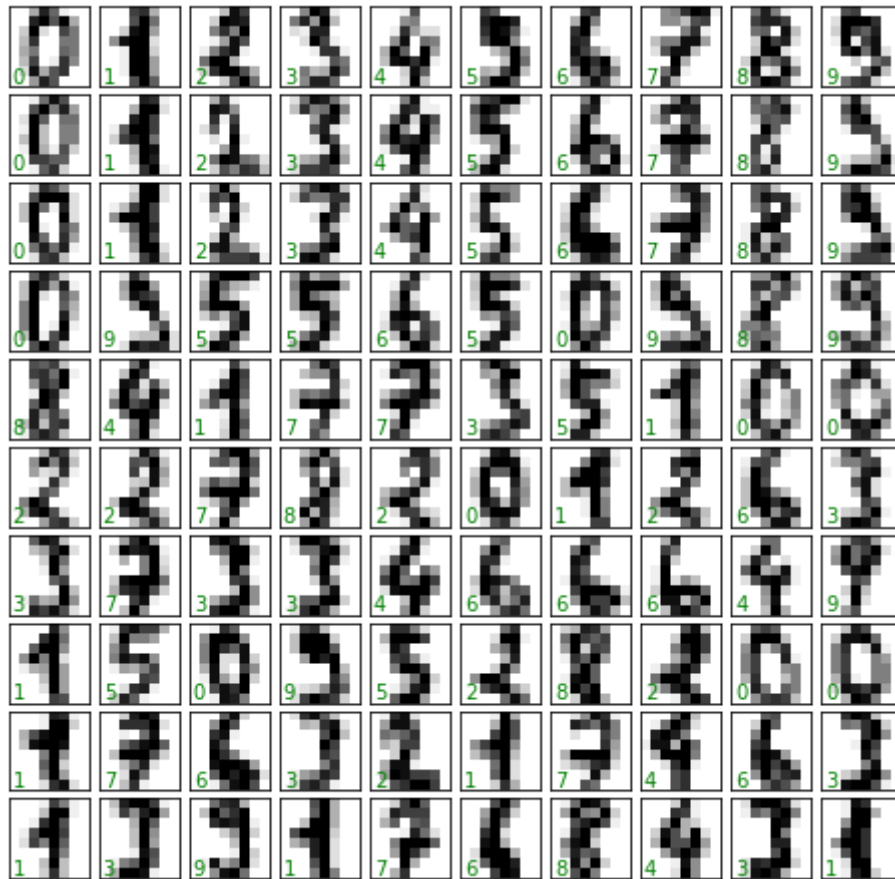

- 实现鸢尾花数据集knn分类，并计算正确率（每类中各取20个测试）

```
1 import pandas as pd
2 from sklearn.datasets import load_iris
3 # (1)加载鸢尾花数据集
4 iris = load_iris()
5 X = iris.data
6 y = iris.target
7 # (2)分割数据
8 from sklearn.model_selection import train_test_split
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3,
10 random_state = 123)
11 # (3)选择模型
12 from sklearn.neighbors import KNeighborsClassifier
13 # (4)生成模型对象
14 knn = KNeighborsClassifier(n_neighbors = 3)
15 # (5)训练模型(数据拟合)
16 knn.fit (X, y)
17 # (6)模型预测
18 # (6)-A单个数据预测
19 knn.predict ([[4,3,5,3]]) #输出 array( [2 ])
20 # (6)-B大集合数据预测
21 y_predict_on_train = knn.predict(X_train)
22 y_predict_on_test = knn.predict(X_test)
23 # (7)模型评估
24 from sklearn.metrics import accuracy_score
25 print("训练集的准确率为: {:.2f}%".format (100 * accuracy_score (y_train[:20],
26 y_predict_on_train[:20])))
27 print("测试集的准确率为: {:.2f}%".format (100 * accuracy_score (y_test[:20],
28 y_predict_on_test[:20])))
```

```
1 训练集的准确率为: 100.00%
2 测试集的准确率为: 95.00%
```

- 对MNIST手写数字数据集使用knn分类，计算正确率；先对MNIST手写数字数据集用PCA降维，选择合适的维数，再使用knn分类，比较两者识别率；

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn import datasets
4
5 digits = datasets.load_digits() # 加载数据
6
7 digits.images.shape
8
9 # 绘制坐标轴
10 fig, axes = plt.subplots(10,10, figsize=(8, 8), subplot_kw={'xticks':[], 'yticks':
11 []},
12 gridspec_kw=dict(hspace=0.1, wspace=0.1))
13
14 # 显示图像
15 for i, ax in enumerate(axes.flat):
16     ax.imshow(digits.images[i], cmap='binary', interpolation='nearest')
17     ax.text(0.05, 0.05, str(digits.target[i]), transform=ax.transAxes,
18 color='green')
```



```

1  X = digits.data # 数据
2  y = digits.target # 标签
3
4  from sklearn.model_selection import train_test_split
5  # 随机划分为训练集和测试集
6  X_train,X_test,y_train,y_test = train_test_split(X,y,random_state=666)
7  X_train.shape
8
9
10 # 不对数据降维，使用knn对手写数字分类并测试性能
11 from sklearn.neighbors import KNeighborsClassifier
12
13 # 实现k近邻投票算法的分类器。
14 knn_clf = KNeighborsClassifier()
15 # 拟合
16 knn_clf.fit(X_train,y_train)
17 # 正确率
18 knn_clf.score(X_test,y_test)
19

```

```

1  0.9866666666666667

```

```

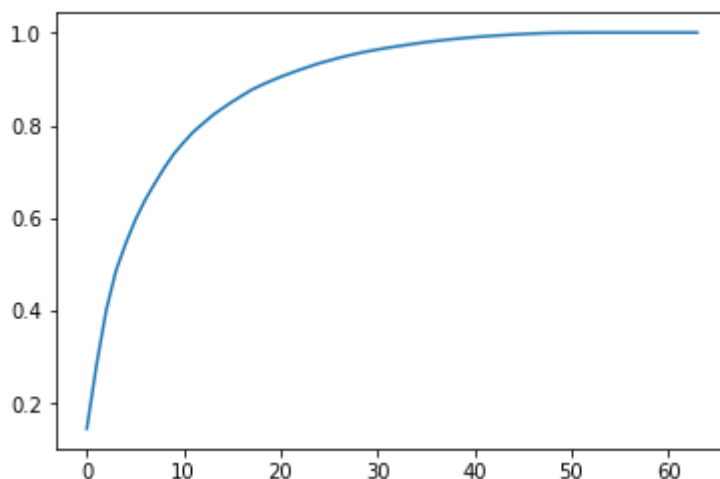
1  '''PCA降维后，使用knn'''
2
3  from sklearn.decomposition import PCA
4

```

```

5  # 降维后的主成分数量为2的PCA
6  pca = PCA(n_components=2)
7  # 拟合
8  pca.fit(X_train)
9  # 降维
10 X_train_reduction = pca.transform(X_train)
11 X_test_reduction = pca.transform(X_test)
12
13 pca = PCA(n_components=X_train.shape[1])
14 pca.fit(X_train)
15 #通过sklearn.PCA.explaine_variance_ration_来查看刚刚的2个纬度的方差爱解释度:
16 pca.explained_variance_ratio_
17
18
19 '''表示取前n个主成分能解释多少百分比的方差'''
20 plt.plot([i for i in range(X_train.shape[1])],\
21          [np.sum(pca.explained_variance_ratio_[:i+1]) for i in
22           range(X_train.shape[1])])
23 plt.show()
24
25 pca = PCA(0.95)
26 pca.fit(X_train)
27 #查看选择特征的数量
28 pca.n_components_ # 28
29
30 X_train_reduction = pca.transform(X_train)
31 X_test_reduction = pca.transform(X_test)
32
33 kcc_clf = KNeighborsClassifier()
34 knn_clf.fit(X_train_reduction, y_train)
35 knn_clf.score(X_test_reduction, y_test)

```



1 0.98

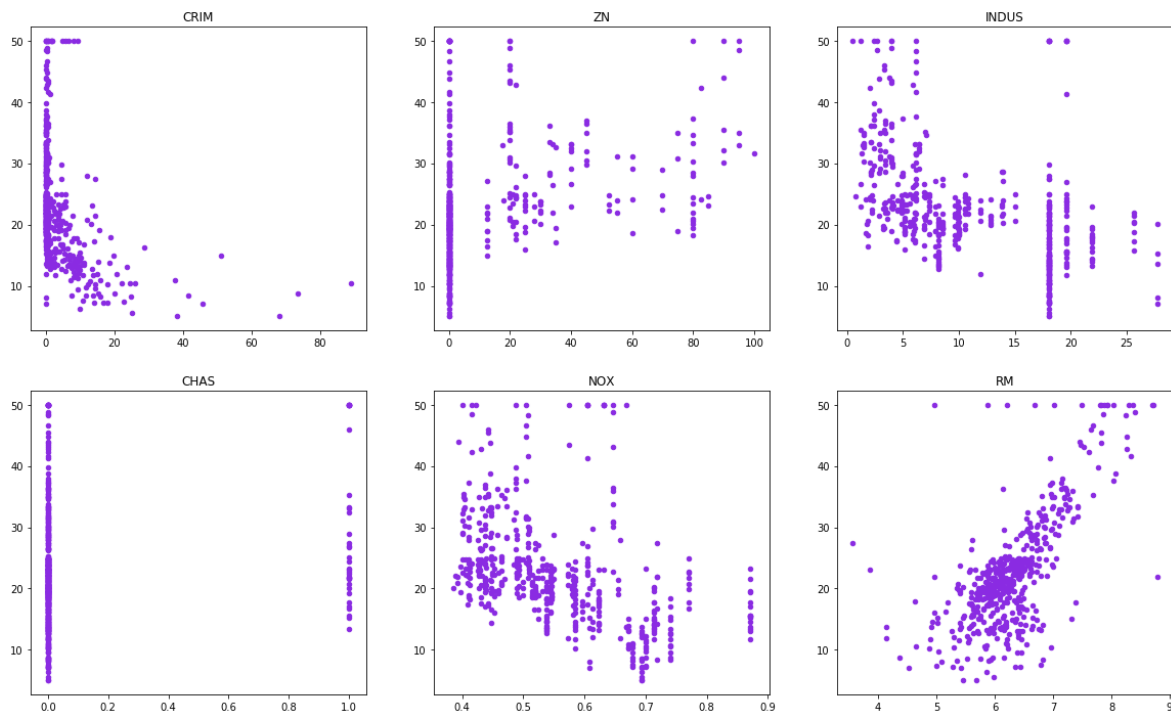
虽然正确率的降低比较小，但是训练时间降低非常多。

- 分析波士顿房价数据集，对其建模并实现预测

```
1  # 导入库
2  from sklearn.datasets import load_boston
3  import pandas as pd
4  from pandas import Series, DataFrame
5  import numpy as np
6  from matplotlib import pyplot as plt
7
8  #导入数据集
9  boston_data=load_boston()
10 x_data = boston_data.data
11 y_data = boston_data.target
12
13 names=boston_data.feature_names
14 FeaturesNums = 6
15 DataNums = len(x_data)
```

观察特征与标签的关系：

```
1  # 每个Feature和target二维关系图
2  plt.subplots(figsize=(20,12))
3  for i in range(FeaturesNums):
4      plt.subplot(231+i)
5      plt.scatter(x_data[:,i],y_data,s=20,color='blueviolet')
6      plt.title(names[i])
7  plt.show()
```

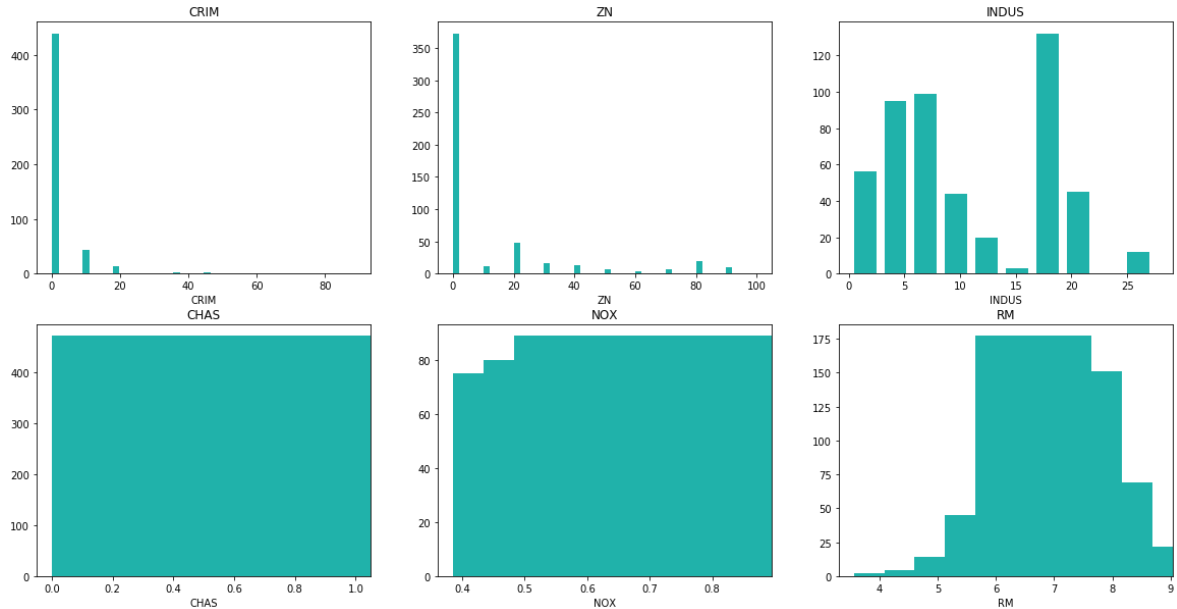


特征数据的分布：


```

1 plt.subplots(figsize=(20,10))
2 for i in range(FeaturesNums):
3
4     plt.subplot(231+i)
5     plt.hist(x_data[:,i],color='lightseagreen',width=2)
6     plt.xlabel(names[i])
7     plt.title(names[i])
8 plt.show()

```



数据处理:

```

1 from sklearn import preprocessing
2
3 # 清除异常特征
4 DelList0=[]
5 for i in range(DataNums):
6     if (y_data[i] >= 49 or y_data[i] <= 1):
7         DelList0.append(i)
8 DataNums -= len(DelList0)
9 x_data = np.delete(x_data, DelList0, axis=0)
10 y_data = np.delete(y_data, DelList0, axis=0)
11
12 # 去除无用特征
13 DelList1=[]
14 for i in range(FeaturesNums):
15     if (names[i] == 'ZN' or
16         names[i] == 'INDUS' or
17         names[i] == 'RAD' or
18         names[i] == 'TAX' or
19         names[i] == 'CHAS' or
20         names[i] == 'NOX' or
21         names[i] == 'B' or
22         names[i] == 'PTRATIO'):
23         DelList1.append(i)
24 x_data = np.delete(x_data, DelList1, axis=1)
25 names = np.delete(names, DelList1)
26 FeaturesNums -= len(DelList1)
27
28 #数据分割
29 from sklearn.model_selection import train_test_split

```

```

30 x_train, x_test, y_train, y_test = train_test_split(x_data, y_data,
31 test_size=0.3)
32 # 归一化
33 from sklearn.preprocessing import MinMaxScaler, scale
34 nms = MinMaxScaler()
35 x_train = nms.fit_transform(x_train)
36 x_test = nms.fit_transform(x_test)
37 y_train = nms.fit_transform(y_train.reshape(-1,1))
38 y_test = nms.fit_transform(y_test.reshape(-1,1))
39

```

训练模型:

```

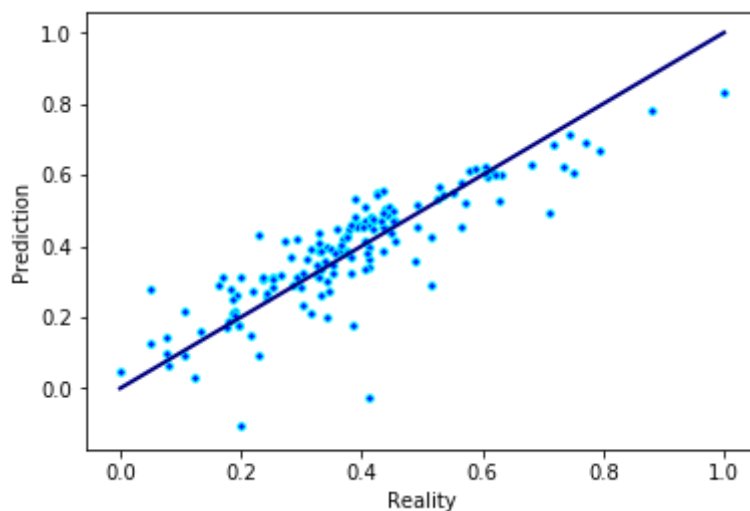
1 # 线性回归
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error, r2_score
4
5 model = LinearRegression()
6 model.fit(x_train, y_train)
7 y_pred = model.predict(x_test)
8 print ("MSE =", mean_squared_error(y_test, y_pred),end='\n\n')
9 print ("R2 =", r2_score(y_test, y_pred),end='\n\n')
10
11 # 画图
12 fig, ax = plt.subplots()
13 ax.scatter(y_test, y_pred, c="blue", edgecolors="aqua",s=13)
14 ax.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k', lw=2,
15 color='navy')
16 ax.set_xlabel('Reality')
17 ax.set_ylabel('Prediction')
18 plt.show()

```

```

1 MSE = 0.007832421992124205
2
3 R2 = 0.7319660538060528

```



```

1 # SVR Linear
2 from sklearn.svm import SVR

```

```

3  from sklearn.model_selection import cross_val_predict, cross_val_score
4  linear_svr = SVR(kernel='linear')
5  # linear_svr.fit(x_train, y_train)
6  # linear_pred = linear_svr.predict(x_test)
7  linear_svr_pred = cross_val_predict(linear_svr, x_train, y_train, cv=5)
8  linear_svr_score = cross_val_score(linear_svr, x_train, y_train, cv=5)
9  linear_svr_meanscore = linear_svr_score.mean()
10 print ("Linear_SVR_Score =",linear_svr_meanscore,end='\n')
11
12
13 from sklearn.svm import SVR
14 from sklearn.model_selection import cross_val_predict, cross_val_score
15 poly_svr = SVR(kernel='poly')
16 poly_svr.fit(x_train, y_train)
17 poly_pred = poly_svr.predict(x_test)
18 poly_svr_pred = cross_val_predict(poly_svr, x_train, y_train, cv=5)
19 poly_svr_score = cross_val_score(poly_svr, x_train, y_train, cv=5)
20 poly_svr_meanscore = poly_svr_score.mean()
21 print ('\n', "Poly_SVR_Score =",poly_svr_meanscore,end='\n')
22

```

```

1  Linear_SVR_Score = 0.7222405483607128
2
3  Poly_SVR_Score = 0.6066134996650142

```

总结：SVR的线性核更好。