

Experiment 2

孟昶-2021214431

November 9, 2021

Task:

Use the sklearn classification models to finish the experiment.

Goal:

We hope you know better about the classification models and try to use the models better on the given dataset.

Data:

For CS background students, you should use the attached dataset file.

For Non CS background students, you can use the sklearn package to load the MNIST dataset.

Divide the data into train set and test set.

Experiment Step:

- Use the SVM to classify the data. Try to change kernel and other related parameters to see the change of metric you use. Record the result and try to explain it.
- Use the Decision Tree to classify the data, use the API to visualize the tree. Try to change the related parameters to see the change of metric you use. Record the result and try to explain it.
- Use the KMeans to cluster the data and visualize it. (The PCA may not be covered in the lecture currently, just use it.)

Answer:

1.SVM

代码:

```
1 from sklearn.svm import SVC
2 import numpy as np
3 from numpy.lib.function_base import select
4 import pandas as pd
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.model_selection import KFold
9 from sklearn.metrics import accuracy_score
10
11 if __name__ == '__main__':
12     data = pd.read_csv('Data.csv')
13     X_train, X_test, y_train, y_test = train_test_split(data.iloc[:, :-1], data['fetal_health'],
14                                                         test_size=0.1, random_state=666)
15     ss = StandardScaler()
16     ss.fit(X_train)
17     X_train = ss.transform(X_train)
18     X_test = ss.transform(X_test)
19     kernels = ['linear', 'poly', 'rbf', 'sigmoid']
20     Cs = [0.01, 1, 5]
21     for i in kernels:
22         for j in Cs:
23             p1 = SVC(C = j, kernel= i )
24             p1.fit(X_train, y_train)
25             print(i, 'C = ', j, 'score = ', p1.score(X_test, y_test))
26
27     X_train = np.array(X_train)
28     y_train = np.array(y_train)
29     X_test = np.array(X_test)
30     y_test = np.array(y_test)
31     y_train = y_train.ravel()
32     y_test = y_test.ravel()
33     kf = KFold(n_splits=5)
34
35     num = 0
36     for j in kernels:
37         for i in Cs:
38             cls = SVC(C = i, kernel=j)
39             t = 1
40             for train_index, test_index in kf.split(X_train, y_train):
41                 X_train_1, X_test_1 = X_train[train_index], X_train[test_index]
42                 Y_train_1, Y_test_1 = y_train[train_index], y_train[test_index]
43                 cls.fit(X_train_1, Y_train_1)
44                 y_pred = cls.predict(X_test_1)
45                 Y_test = Y_test_1
46                 print('第', t, '次:', j, i)
47                 print(accuracy_score(Y_test, y_pred))
48                 # print ("test_Accuracy:", accuracy)
49                 # print ("test_error_rate:", 1-accuracy)
50                 t += 1
```

Table 1: SVM 测试集上的得分和交叉验证得分表

kernel type	num of C	test score	cross-validation score				
			1st	2nd	3rd	4th	5th
linear	0.01	0.901	0.877	0.896	0.846	0.874	0.898
	1	0.915	0.896	0.888	0.872	0.901	0.921
	5	0.920	0.901	0.877	0.872	0.903	0.911
polynomial	0.01	0.850	0.781	0.836	0.809	0.801	0.825
	1	0.939	0.906	0.901	0.867	0.893	0.906
	5	0.944	0.924	0.906	0.880	0.901	0.927
Gaussian	0.01	0.826	0.742	0.804	0.781	0.775	0.764
	1	0.925	0.903	0.901	0.885	0.898	0.919
	5	0.948	0.916	0.911	0.901	0.919	0.935
sigmoid	0.01	0.826	0.742	0.804	0.781	0.775	0.764
	1	0.817	0.789	0.794	0.7395	0.764	0.764
	5	0.803	0.765	0.778	0.744	0.770	0.772

惩罚系数 C 是最重要的参数： C 越大则对错误分类样本的惩罚力度越大，因此训练集准确率更高，但是泛化能力会有所降低即测试集分类准确率降低； C 越小则对错误分类样本的惩罚力度越小，泛化能力有所增强。因此若出现过拟合现象，则应该适当减小 C ，反之出现欠拟合现象则应该适当增大 C 。

从表中可以看出， C 太小的话会导致训练出现欠拟合现象，导致精度下降，从而使得训练集与测试集错误率都较高；而 C 偏大的话可能会出现过拟合的情况，从 Gaussian 和 polynomial 中的结果就可以看出。从使用的 kernel 上可以看出，Gaussian 与 polynomial 核函数在训练集上错误率较低，而 Linear 在测试集上表现最好，其测试集错误率最低，具有较好的泛化能力。

2.Decision Tree

代码:

```
1 from sklearn.tree import DecisionTreeClassifier
2 import numpy as np
3 from numpy.lib.function_base import select
4 import pandas as pd
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.model_selection import KFold
9 from sklearn.metrics import accuracy_score
10
11 if __name__ == '__main__':
12     data = pd.read_csv('Data.csv')
13     X_train, X_test, y_train, y_test = train_test_split(data.iloc[:, :-1], data['fetal_health'],
14                                                         test_size=0.1, random_state=666)
15     feature_names = X_train.keys()
16     class_names = ['1.0', '2.0', '3.0']
17     ss = StandardScaler()
18     ss.fit(X_train)
19     X_train = ss.transform(X_train)
20     X_test = ss.transform(X_test)
21     kernels = ["gini", "entropy"]
22     Cs = [0, 0.015, 0.030, 0.045]
23     for i in kernels:
24         for j in Cs:
25             p1 = DecisionTreeClassifier(criterion = i, ccp_alpha = j)
26             p1.fit(X_train, y_train)
27             print(i, 'C = ', j, 'score = ', p1.score(X_test, y_test))
28
29     X_train = np.array(X_train)
30     y_train = np.array(y_train)
31     X_test = np.array(X_test)
32     y_test = np.array(y_test)
33     y_train = y_train.ravel()
34     y_test = y_test.ravel()
35     kf = KFold(n_splits=3)
36
37     num = 0
38     for j in kernels:
39         for i in Cs:
40             cls = DecisionTreeClassifier(criterion = j, ccp_alpha = i)
41             t = 1
42             for train_index, test_index in kf.split(X_train, y_train):
43
44                 X_train_1, X_test_1 = X_train[train_index], X_train[test_index]
45                 Y_train_1, Y_test_1 = y_train[train_index], y_train[test_index]
46                 cls.fit(X_train_1, Y_train_1)
47                 y_pred = cls.predict(X_test_1)
48                 Y_test = Y_test_1
49                 print('第', t, '次:', j, i)
50                 print(accuracy_score(Y_test, y_pred))
51
52                 t += 1
```

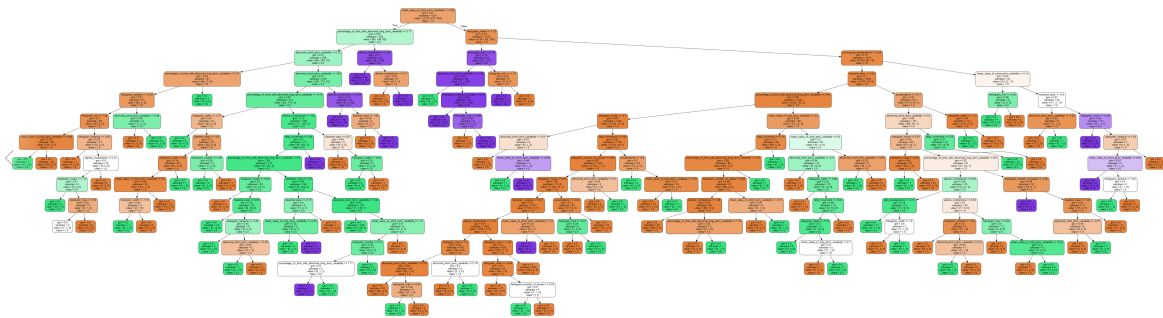


Figure 1: 决策树可视化结果图

Table 2: 决策树测试集上的得分和交叉验证得分表

criterion	ccp_alpha	test score	cross-validation score		
			1st	2nd	3rd
gini	0	0.892	0.917	0.918	0.922
	0.015	0.883	0.901	0.901	0.926
	0.030	0.836	0.862	0.856	0.879
	0.045	0.836	0.862	0.856	0.879
entropy	0	0.915	0.936	0.923	0.928
	0.015	0.906	0.928	0.925	0.917
	0.030	0.887	0.917	0.909	0.900
	0.045	0.883	0.879	0.900	0.895

决策树学习是一种采用树状结构的有监督机器学习方法。决策树是一个预测模型，表示对象特征和对象值之间的一种映射。其原理通俗的说，即给定一个输入值，从树节点不断往下走，直至走到叶节点，这个叶节点就是对输入值的一个预测或者分类。Sklearn 中实现了其三种算法 ID3、C4.5、CART，最常用的 CART 算法其核心公式就是基尼指数的计算，基尼指数越大不确定越大，基尼指数的计算公式为：

$$Gini(D) = 1 - \sum p_i^2$$

其中 p_i 是 D 中元组 C_i 类的概率。

criterion: "gini", "entropy", default="gini": 特征选择标准。

gini 原理：当一个 node p 被分为 k 部分，连续型的属性分割，先把数据有序排列，然后从左到右每个都是一遍，找 gini 最小的那一个。

entropy 原理：log 以 2 为底。当所有分类都一样是 $1/k$ 的时候，熵值最大，是 $\log(nc)$ ；所有值都属于同一类的时候，熵最小，是 0

ccp_alpha : non-negative float, default=0.0: 复杂度参数。用于最小成本复杂度修剪。

随着剪枝的程度越大，会有效抑制过拟合现象，但是当剪枝的程度过大时，模型无法有效拟合，从而使得精度下降。

3.KMeans

代码:

```
1 from sklearn.cluster import KMeans
2 from sklearn.decomposition import PCA
3 import numpy as np
4 from numpy.lib.function_base import select
5 import pandas as pd
6 from sklearn.model_selection import train_test_split
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.model_selection import KFold
9 from sklearn.metrics import accuracy_score
10 from time import time
11 from sklearn import metrics
12 from sklearn.pipeline import make_pipeline
13 import matplotlib.pyplot as plt
14
15 def bench_k_means(kmeans, name, data, labels):
16     t0 = time()
17     estimator = make_pipeline(StandardScaler(), kmeans).fit(data)
18     fit_time = time() - t0
19     results = [name, fit_time, estimator[-1].inertia_]
20
21     # Define the metrics which require only the true labels and estimator
22     # labels
23     clustering_metrics = [
24         metrics.homogeneity_score,
25         metrics.completeness_score,
26         metrics.v_measure_score,
27         metrics.adjusted_rand_score,
28         metrics.adjusted_mutual_info_score,
29     ]
30     results += [m(labels, estimator[-1].labels_) for m in clustering_metrics]
31
32     # The silhouette score requires the full dataset
33     results += [
34         metrics.silhouette_score(
35             data,
36             estimator[-1].labels_,
37             metric="euclidean",
38             sample_size=300,
39         )
40     ]
41
42     # Show the results
43     formatter_result = (
44         "{:9s}\t{:.3f}s\t{:.0f}\t{:.3f}\t{:.3f}\t{:.3f}\t{:.3f}\t{:.3f}\t{:.3f}"
45     )
46     print(formatter_result.format(*results))
47
48
49
50
51 if __name__ == '__main__':
52     data = pd.read_csv('Data.csv')
53     X_train, X_test, y_train, y_test = train_test_split(data.iloc[:, :-1], data['fetal_health'],
54                                                         test_size=0.1, random_state=666)
55     feature_names = X_train.keys()
```

```

55 class_names = ['1.0', '2.0', '3.0']
56 ss = StandardScaler()
57 ss.fit(X_train)
58 X_train = ss.transform(X_train)
59 X_test = ss.transform(X_test)
60
61 (n_samples, n_features), n_digits = X_train.shape, np.unique(y_train).size
62 print(82 * "_")
63 print("init\t\ttime\tinertia\thomo\tcompl\tv-meas\tARI\tAMI\tsilhouette")
64 kmeans = KMeans(init="k-means++", n_clusters=n_digits, n_init=4, random_state=0)
65 bench_k_means(kmeans=kmeans, name="k-means++", data=X_train, labels=y_train)
66
67 kmeans = KMeans(init="random", n_clusters=n_digits, n_init=4, random_state=0)
68 bench_k_means(kmeans=kmeans, name="random", data=X_train, labels=y_train)
69
70 pca = PCA(n_components=n_digits).fit(X_train)
71 kmeans = KMeans(init=pca.components_, n_clusters=n_digits, n_init=1)
72 bench_k_means(kmeans=kmeans, name="PCA-based", data=X_train, labels=y_train)
73 print(82 * "_")
74
75
76
77 # 画图
78 reduced_data = PCA(n_components=2).fit_transform(X_train)
79 kmeans = KMeans(init="k-means++", n_clusters=n_digits, n_init=4)
80 kmeans.fit(reduced_data)
81
82 # Step size of the mesh. Decrease to increase the quality of the VQ.
83 h = 0.02 # point in the mesh [x_min, x_max]x[y_min, y_max].
84
85 # Plot the decision boundary. For that, we will assign a color to each
86 x_min, x_max = reduced_data[:, 0].min() - 1, reduced_data[:, 0].max() + 1
87 y_min, y_max = reduced_data[:, 1].min() - 1, reduced_data[:, 1].max() + 1
88 xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
89
90 # Obtain labels for each point in mesh. Use last trained model.
91 Z = kmeans.predict(np.c_[xx.ravel(), yy.ravel()])
92
93 # Put the result into a color plot
94 Z = Z.reshape(xx.shape)
95 plt.figure(1)
96 plt.clf()
97 plt.imshow(
98     Z,
99     interpolation="nearest",
100     extent=(xx.min(), xx.max(), yy.min(), yy.max()),
101     cmap=plt.cm.Paired,
102     aspect="auto",
103     origin="lower",
104 )
105
106 plt.plot(reduced_data[:, 0], reduced_data[:, 1], "k.", markersize=2)
107 # Plot the centroids as a white X
108 centroids = kmeans.cluster_centers_
109 plt.scatter(
110     centroids[:, 0],
111     centroids[:, 1],
112     marker="x",

```

```

113         s=169,
114         linewidths=3,
115         color="w",
116         zorder=10,
117     )
118     plt.title(
119         "K-means clustering on the excal dataset (PCA-reduced data)"
120     )
121     plt.xlim(x_min, x_max)
122     plt.ylim(y_min, y_max)
123     plt.xticks(())
124     plt.yticks(())
125     plt.show()

```

如图所示，为结果图

init	time	inertia	homo	compl	v-meas	ARI	AMI	silhouette
k-means++	0.048s	29011	0.248	0.179	0.208	0.179	0.207	0.163
random	0.036s	29005	0.248	0.178	0.207	0.205	0.206	0.162
PCA-based	0.020s	29012	0.240	0.171	0.200	0.170	0.199	0.157

Figure 2: KMeans 各指标结果图

K-means clustering on the excal dataset (PCA-reduced data)

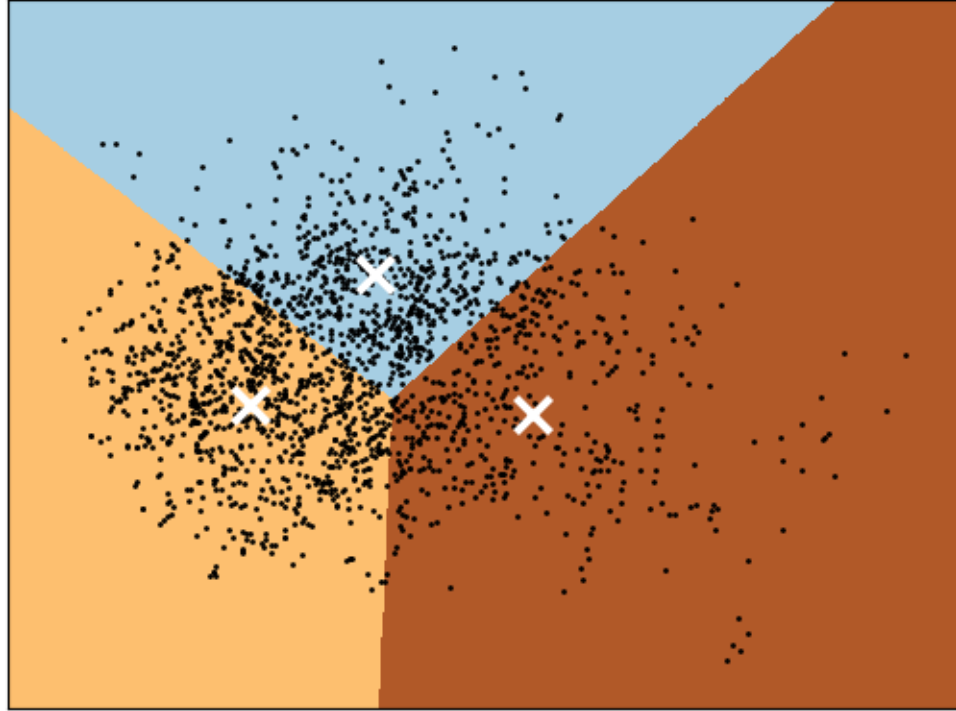


Figure 3: KMeans 分类结果图