

一、利用 sklearn 库和 iris 数据集实现分类

```
18 import numpy as np
19 import pandas as pd
20 from sklearn.datasets import load_iris
21 from sklearn.model_selection import train_test_split
22 import matplotlib.pyplot as plt
23
24 # data
25 def create_data():
26     iris = load_iris()
27     df = pd.DataFrame(iris.data, columns=iris.feature_names)
28     df['label'] = iris.target
29     df.columns = ['sepal length', 'sepal width', \
30                  'petal length', 'petal width', 'label']
31     data = np.array(df.iloc[:100, [0, 1, -1]])
32     for i in range(len(data)):
33         if data[i, -1] == 0:
34             data[i, -1] = -1
35     # print(data)
36     return data[:, :2], data[:, -1]
37
38 X, y = create_data()
39 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
40
41 plt.scatter(X[:50, 0], X[:50, 1], label='0')
42 plt.scatter(X[50:, 0], X[50:, 1], label='1')
43 plt.legend()
44
45 from sklearn.ensemble import AdaBoostClassifier
46 clf = AdaBoostClassifier(n_estimators=100, learning_rate=0.5)
47 clf.fit(X_train, y_train)
48
49 print (clf.score(X_test, y_test))
```

二、自己实现 adaboost 算法

```
5 import numpy as np
6 import pandas as pd
7 from sklearn.datasets import load_iris
8 from sklearn.model_selection import train_test_split
9 import matplotlib.pyplot as plt
10
11 # data
12 def create_data():
13     iris = load_iris()
14     df = pd.DataFrame(iris.data, columns=iris.feature_names)
15     df['label'] = iris.target
16     df.columns = ['sepal length', 'sepal width', \
17                  'petal length', 'petal width', 'label']
18     data = np.array(df.iloc[:100, [0, 1, -1]])
19     for i in range(len(data)):
20         if data[i, -1] == 0:
21             data[i, -1] = -1
22     # print(data)
23     return data[:, :2], data[:, -1]
```

```

25 X, y = create_data()
26 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
27
28 plt.scatter(X[:50,0],X[:50,1], label='0')
29 plt.scatter(X[50:,0],X[50:,1], label='1')
30 plt.legend()
31
32 class AdaBoost:
33     def __init__(self, n_estimators=50, learning_rate=1.0):
34         self.clf_num = n_estimators
35         self.learning_rate = learning_rate
36
37     def init_args(self, datasets, labels):
38         self.X = datasets
39         self.Y = labels
40         self.M, self.N = datasets.shape
41
42         # 弱分类器数目和集合
43         self.clf_sets = []
44
45         # 初始化weights
46         self.weights = [1.0/self.M]*self.M
47
48         # G(x)系数 alpha
49         self.alpha = []
50
51     def _G(self, features, labels, weights):
52         m = len(features)
53         error = 100000.0 # 无穷大
54         best_v = 0.0
55         # 单维features
56         features_min = min(features)
57         features_max = max(features)
58         n_step = (features_max - features_min + self.learning_rate) // self.learning_rate
59         # print('n_step:{}'.format(n_step))
60         direct, compare_array = None, None
61         for i in range(1, int(n_step)):
62             v = features_min + self.learning_rate * i
63
64             if v not in features:
65                 # 误分类计算
66                 compare_array_positive = np.array([1 if features[k] > v else -1 for k in range(m)])
67                 weight_error_positive = sum([weights[k] for k in range(m) \
68                     if compare_array_positive[k] != labels[k]])
69
70                 compare_array_nagetive = np.array([-1 if features[k] > v else 1 for k in range(m)])
71                 weight_error_nagetive = sum([weights[k] for k in range(m) \
72                     if compare_array_nagetive[k] != labels[k]])
73
74                 if weight_error_positive < weight_error_nagetive:
75                     weight_error = weight_error_positive
76                     _compare_array = compare_array_positive
77                     direct = 'positive'
78                 else:
79                     weight_error = weight_error_nagetive
80                     _compare_array = compare_array_nagetive
81                     direct = 'nagetive'

```

```

83         # print('v:{} error:{}'.format(v, weight_error))
84         if weight_error < error:
85             error = weight_error
86             compare_array = _compare_array
87             best_v = v
88         return best_v, direct, error, compare_array
89
90     # 计算alpha
91     def _alpha(self, error):
92         return 0.5 * np.log((1-error)/error)
93
94     # 规范化因子
95     def _Z(self, weights, a, clf):
96         return sum([weights[i]*np.exp(-1*a*self.Y[i]*clf[i]) for i in range(self.M)])
97
98     # 权值更新
99     def _w(self, a, clf, Z):
100         for i in range(self.M):
101             self.weights[i] = self.weights[i]*np.exp(-1*a*self.Y[i]*clf[i])/ Z
102
103     # G(x) 的线性组合
104     def _f(self, alpha, clf_sets):
105         pass
106
107     def G(self, x, v, direct):
108         if direct == 'positive':
109             return 1 if x > v else -1
110         else:
111             return -1 if x > v else 1
112
113
114 def fit(self, X, y):
115     self.init_args(X, y)
116     for epoch in range(self.clf_num):
117         best_clf_error, best_v, clf_result = 100000, None, None
118         # 根据特征维度, 选择误差最小的
119         for j in range(self.N):
120             features = self.X[:, j]
121             # 分类阈值, 分类误差, 分类结果
122             v, direct, error, compare_array = self._G(features, self.Y, self.weights)
123
124             if error < best_clf_error:
125                 best_clf_error = error
126                 best_v = v
127                 final_direct = direct
128                 clf_result = compare_array
129                 axis = j
130
131         # print('epoch:{}/{} feature:{} error:{} v:{}'.format(epoch, self.clf_num, j, error, best_v))
132         if best_clf_error == 0:
133             break
134
135     # 计算G(x)系数a
136     a = self._alpha(best_clf_error)
137     self.alpha.append(a)
138     # 记录分类器
139     self.clf_sets.append((axis, best_v, final_direct))
140     # 规范化因子
141     Z = self._Z(self.weights, a, clf_result)
142     # 权值更新
143     self._w(a, clf_result, Z)

```

```

145 # print('classifier:{}/{} error:{:.3f} v:{ } direct:{ } a:{:.5f}'. \
146 #format(epoch+1, self.clf_num, error, best_v, final_direct, a))
147 # print('weight:{ }'.format(self.weights))
148 # print('\n')
149
150 def predict(self, feature):
151     result = 0.0
152     for i in range(len(self.clf_sets)):
153         axis, clf_v, direct = self.clf_sets[i]
154         f_input = feature[axis]
155         result += self.alpha[i] * self.G(f_input, clf_v, direct)
156     # sign
157     return 1 if result > 0 else -1
158
159 def score(self, X_test, y_test):
160     right_count = 0
161     for i in range(len(X_test)):
162         feature = X_test[i]
163         if self.predict(feature) == y_test[i]:
164             right_count += 1
165
166     return right_count / len(X_test)
167
168 # 100次结果
169 result = []
170 for i in range(1, 101):
171     X, y = create_data()
172     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33)
173     clf = AdaBoost(n_estimators=100, learning_rate=0.2)
174     clf.fit(X_train, y_train)
175     r = clf.score(X_test, y_test)
176     print('{} / 100 score: {}'.format(i, r))
177     result.append(r)
178
179 print('average score:{:.3f}%'.format(sum(result)))

```

三、实现提升树算法

```

6 from collections import defaultdict
7 import numpy as np
8
9
10 class BoostingTree:
11     def __init__(self, epsilon=1e-2):
12         self.epsilon = epsilon
13         self.cand_splits = [] # 候选切分点
14         self.split_index = defaultdict(tuple) # 由于要多次切分数据集，故预先存储，切分后数据点的索引
15         self.split_list = [] # 最终各个基本回归树的切分点
16         self.c1_list = [] # 切分点左区域取值
17         self.c2_list = [] # 切分点右区域取值
18         self.N = None
19         self.n_split = None
20
21     def init_param(self, X_data):
22         # 初始化参数
23         self.N = X_data.shape[0]
24         for i in range(1, self.N):
25             self.cand_splits.append((X_data[i][0] + X_data[i - 1][0]) / 2)
26         self.n_split = len(self.cand_splits)
27         for split in self.cand_splits:
28             left_index = np.where(X_data[:, 0] <= split)[0]
29             right_index = list(set(range(self.N)) - set(left_index))
30             self.split_index[split] = (left_index, right_index)
31         return

```

```
33 def _cal_err(self, split, y_res):
34     # 计算每个切分点的误差
35     inds = self.split_index[split]
36     left = y_res[inds[0]]
37     right = y_res[inds[1]]
38
39     c1 = np.sum(left) / len(left)
40     c2 = np.sum(right) / len(right)
41     y_res_left = left - c1
42     y_res_right = right - c2
43     res = np.hstack([y_res_left, y_res_right])
44     res_square = np.apply_along_axis(lambda x: x ** 2, 0, res).sum()
45     return res_square, c1, c2
46
47 def best_split(self, y_res):
48     # 获取最佳切分点，并返回对应的残差
49     best_split = self.cand_splits[0]
50     min_res_square, best_c1, best_c2 = self._cal_err(best_split, y_res)
51
52     for i in range(1, self.n_split):
53         res_square, c1, c2 = self._cal_err(self.cand_splits[i], y_res)
54         if res_square < min_res_square:
55             best_split = self.cand_splits[i]
56             min_res_square = res_square
57             best_c1 = c1
58             best_c2 = c2
59
60     self.split_list.append(best_split)
61     self.c1_list.append(best_c1)
62     self.c2_list.append(best_c2)
63     return
64
65 def _fx(self, X):
66     # 基于当前组合树，预测X的输出值
67     s = 0
68     for split, c1, c2 in zip(self.split_list, self.c1_list, self.c2_list):
69         if X < split:
70             s += c1
71         else:
72             s += c2
73     return s
74
75 def update_y(self, X_data, y_data):
76     # 每添加一颗回归树，就要更新y，即基于当前组合回归树的预测残差
77     y_res = []
78     for X, y in zip(X_data, y_data):
79         y_res.append(y - self._fx(X[0]))
80     y_res = np.array(y_res)
81     res_square = np.apply_along_axis(lambda x: x ** 2, 0, y_res).sum()
82     return y_res, res_square
83
84 def fit(self, X_data, y_data):
85     self.init_param(X_data)
86     y_res = y_data
87     while True:
88         self.best_split(y_res)
89         y_res, res_square = self.update_y(X_data, y_data)
90         if res_square < self.epsilon:
91             break
92     return
93
94 def predict(self, X):
95     return self._fx(X)
```

```
98 if __name__ == '__main__':
99     data = np.array(
100         [[1, 5.56], [2, 5.70], [3, 5.91], [4, 6.40], [5, 6.80], \
101          [6, 7.05], [7, 8.90], [8, 8.70], [9, 9.00], [10, 9.05]])
102     X_data = data[:, :-1]
103     y_data = data[:, -1]
104     BT = BoostingTree(epsilon=0.18)
105     BT.fit(X_data, y_data)
106     print(BT.split_list)
107     print(BT.c1_list)
108     print(BT.c2_list)
109     '''
110     X_data_raw = np.linspace(-5, 5, 100)
111     X_data = np.transpose([X_data_raw])
112     y_data = np.sin(X_data_raw)
113     BT = BoostingTree(epsilon=0.1)
114     BT.fit(X_data, y_data)
115     y_pred = [BT.predict(X) for X in X_data]
116
117     import matplotlib.pyplot as plt
118
119     p1 = plt.scatter(X_data_raw, y_data, color='r')
120     p2 = plt.scatter(X_data_raw, y_pred, color='b')
121     plt.legend([p1, p2], ['real', 'pred'])
122     plt.show()
123     '''
```