

## AdaBoost 分类器代码：

jupyter AdaBoost.py ✓ 6 minutes ago Logout

File Edit View Language Python

```
1  #!/usr/bin/env python
2  # coding: utf-8
3
4  # In[1]:
5
6
7  # Import required packages
8  import pandas as pd
9  import numpy as np
10 import math
11 from sklearn.tree import DecisionTreeClassifier
12 from sklearn.base import BaseEstimator, ClassifierMixin
13
14
15 # In[2]:
16
17
18 """ Summary: The model will be a AdaBoost classification.
19
20 We will build AdaBoost class to make classification. There will be 9 hyper-parameters in
21 the model. The model will contain fit(), predict(), get_params(), set_params(),
22 score(), decision_function(), __init__().
23
24 """
25
26 class AdaBoostClassifier(BaseEstimator, ClassifierMixin):
27
28     """ A class used to make classification.
29
30     Attributes
31
32     - - - - -
33     n_estimators : int
34         The number of weak classifiers
35     base_estimator : class
36         The base estimator we choose to train the AdaBoost.
37     estimators_ : list
38         The list of base_estimators
39
40     """
41
42     def __init__(self, n_estimators=20, base_estimator=DecisionTreeClassifier(), estimators_=None):
43         """
44
45         Parameters
46         - - - - -
47         n_estimators : int
48             The number of weak classifiers
49         base_estimator : class
50             The base estimator we choose to train the AdaBoost.
51         estimators_ : list
52             The list of base_estimators
53
54         """
55
56         self.n_estimators = n_estimators
57         self.base_estimator = base_estimator
58         self.estimators_ = estimators_
59
60     def fit(self, X, y, n_trees=None,
61            sample_weight=None,
62            sample_size=None,
63            max_depth=None,
```

```

63         max_depth=None,
64         max_features=None,
65         min_samples_split=None,
66         min_samples_leaf=None,
67         min_weight_fraction_leaf=None,
68         max_leaf_nodes=None,
69     ):
70
71     """Train the model.
72
73
74     Parameters
75     - - - - -
76     n_trees: int
77         The number of bootstrapped trees
78     sample_weight: list
79         The weight for each samples
80     sample_size: int
81         The size of bootstrapped samples.
82     max_depth: int
83         The maximum depth of a single tree.
84     max_features: int
85         The maximum number of features used to train a single tree.
86     min_samples_split: int
87         The minimum number of samples required to continue splitting.
88     min_samples_leaf: int
89         The minimum number of samples in a leaf.
90     min_weight_fraction_leaf: int
91         The minimum weight fraction.
92     max_leaf_nodes: int
93         The maximum number of nodes.
94

```

```

94
95     Returns
96     - - -
97     self
98
99     """
100
101     # If the number of samples in a leaf is below the Min_samples_Leaf, the leaf would be pruned.
102     # the leaf nodes, and if it is less than this value, it will be pruned together with the sibling nodes.
103     y = y.reshape(-1, 1)
104     # Convert value of 0 to -1 in y.
105     index = (y==0)
106     y[index] = -1
107
108     # The processes below set the initial value of each parameter.
109     # if the 'if' function is true, the default value will be used,
110     # but if False, it means users has input their own values.
111     if self.estimators_ == [None]:
112         self.estimators_ = []
113     if n_trees == None:
114         n_trees = 20
115     if sample_weight == None:
116         sample_weight = np.full(np.shape(y), 1/len(y)).reshape(-1, 1)
117     if sample_size == None:
118         sample_size = X.shape[0]
119     if max_depth == None:
120         max_depth = 5
121     if max_features == None:
122         max_features = None
123     if min_samples_split == None:
124         min_samples_split = 2
125     if min_samples_leaf == None:

```

```

125     if min_samples_leaf == None:
126         min_samples_leaf = 1
127     if min_weight_fraction_leaf == None:
128         min_weight_fraction_leaf = 0
129     if max_leaf_nodes == None:
130         max_leaf_nodes = None
131
132     # Set attributes:
133     # trees_ represents the list of bootstrapped sample trees.
134     # tree_preds records the prediction of each bootstrapped sample trees.
135     # estimator_errors represents saves the errors of each base estimators.
136     # estimator_weights represents saves the updated voting value of each base estimator.
137     trees_ = []
138     tree_preds_ = []
139     self.estimator_errors_ = []
140     self.estimator_weights_ = []
141
142
143     # The processes below will show how we train the weak classifiers.
144     # diffs is the list of values of the difference of the real value and predicted
145     # value for each bootstrapped tree.
146     # For example, If diffs[0] is [1, 1, 0], means the first two prediction is different
147     # from the real value.
148     diffs = []
149     for n in range(n_trees):
150         # Generate m bootstrap samples from a dataset
151         rows = np.random.choice(X.shape[0], sample_size)
152         X_sample = X[rows]
153         y_sample = y[rows]
154
155         clf = DecisionTreeClassifier(max_depth=max_depth,
156                                     max_features=max_features,
157                                     min_samples_split=min_samples_split,
158                                     min_samples_leaf=min_samples_leaf,
159                                     min_weight_fraction_leaf=min_weight_fraction_leaf,
160                                     max_leaf_nodes=max_leaf_nodes
161                                     )
162         clf.fit(X_sample, y_sample)
163
164         tree_pred = clf.predict(X)
165         tree_pred = np.array(tree_pred).reshape(-1,1)
166
167         tree_preds_.append(tree_pred)
168         diff = ((tree_pred-y)/2)**2
169         diffs.append(diff)
170         trees_.append(clf)
171
172     # The iteration process.
173     for i in range(self.n_estimators):
174         # Calculate misclassification rate for each classifier
175         # The choose the minimum misclassification rate.
176         errors = []
177         for t in range(len(tree_preds_)):
178             error = sum(sample_weight*diffs[t])[0]
179             errors.append(error)
180         min_error = min(errors)
181
182         # Stop criteria 1: if the minimum error is Larger than 0.5,
183         # this means the trees can not get better classification than random guess so we'd better stop.
184         if min_error>=0.5:
185             return self
186         # We get the index of the selected base estimator using minerror_index.
187         # And pick the classifier with the lowest error.

```

```

187         # Then take out the diff with Lowest error.
188         minerror_index = errors.index(min_error)
189         base_estimator_ = trees_[minerror_index]
190         self.min_error_diff = diffs[minerror_index]
191
192         # Calculate voting power for the best classifier
193         min_error_2 = np.clip(min_error, 1e-15, 1-1e-15).copy()
194         a = 0.5 * (math.log((1-min_error_2)/min_error_2, 10))
195         # Update the base estimator.
196         # Update the voting power.
197         self.estimators_.append(base_estimator_)
198         self.estimator_weights_.append(a)
199         self.estimator_errors_.append(min_error)
200
201         #Stop point 2: if we got perfect prediction, then we stop.
202         if min_error==0:
203             return self
204         # Update the data sample weights using e of the most recent classifier.
205         # If we make a mistake in the prediction,
206         # then the weight of the sample will get larger then before,
207         # if we predict right, then the weight will get smaller.
208         sample_weight_new = []
209         for i in range(len(self.min_error_diff)):
210             if self.min_error_diff[i]==0:
211                 new = (1/(2*(1-min_error)))*sample_weight[i]
212                 sample_weight_new.append(new)
213             if self.min_error_diff[i]==1:
214                 new = (1/(2*min_error))*sample_weight[i]
215                 sample_weight_new.append(new)
216         sample_weight_new = np.array(sample_weight_new)
217         sample_weight = sample_weight_new
218
219
220         return self
221
222         # We are supposed to return the information of parameters set at the begining.
223         # get_params returns a dict.
224         def get_params(self, deep=True):
225             """ Print the parameter information of the class.
226
227             Returns
228             - - -
229             { 'n_estimators': self.n_estimators,
230               'base_estimator': self.base_estimator, 'estimators_':self.estimators_
231             }
232
233             """
234             return { 'n_estimators': self.n_estimators,
235                       'base_estimator': self.base_estimator, 'estimators_':self.estimators_
236                     }
237
238         # We make prediction in this process.
239         # We just input X.
240         def predict(self, X):
241             """ Prdict the y based on the input X
242
243
244             Returns
245             - - -
246             y
247
248             """

```

```

249
250     preds = []
251     for i in range(len(self.estimators_)):
252         u = self.estimators_[i].predict(X)
253         u = u.reshape(-1, 1)
254         pred = u * self.estimator_weights_[i]
255         preds.append(pred)
256
257     final_preds = np.array(preds).sum(axis=0)
258     y = self.decision_funtion(final_preds)
259
260     return y #predict returns an ndarray
261
262     # score function will give the accuracy of a prediction.
263     def score(self, X, y):
264         """ Calculate the accuracy of the prediction to X and y
265
266             Returns
267             - - - -
268             score
269
270         """
271
272         pred = self.predict(X)
273         pred = np.array(pred).reshape(-1,1)
274
275         y = y.reshape(-1, 1)
276
277         index = (y==0)
278         y[index]=-1
279
280         diff = ((pred-u)/2)**2
281
282         diff = ((pred-y)/2)**2
283         score = 1 - diff.sum()/len(diff)
284         return score #score fuction returns a float
285         # We manually set value of parameters in this step.
286         # set_params has single arg called params, which is a **kwargs.
287         # set_params returns self.
288         def set_params(self, **params):
289             """ Set the parameters for the class.
290
291             Returns
292             - - - -
293             self
294
295             """
296
297             for parameter, value in params.items():
298                 setattr(self, parameter, value)
299             return self
300             #set_params has single arg called params, which is a **kwargs
301             #set_params returns self)
302             # We refer to a sign() function in decision_function()
303             def decision_funtion(self, X):
304                 """ Assign 1 and -1 to the prediction
305
306                 Returns
307                 - - - -
308                 y
309
310                 """
311
312                 y = np.sign(X)
313                 return y
314
315     # In[ ]:
316
317
318
319
320
321

```