

**Московский авиационный институт
(национальный исследовательский университет)**

**Факультет информационных технологий и прикладной
математики**

Кафедра вычислительной математики и программирования

Лабораторная работа №2 по курсу «Искусственный интеллект»

Студент: К. М. Воронов
Преподаватель: Самир Ахмед
Группа: М8О-307Б-19
Дата:
Оценка:
Подпись:

Москва, 2022

Задача

Вы построили базовые (слабые) модели машинного обучения под вашу задачу. Некоторые задачи показали себя не очень, некоторые показали себя хорошо. Как выяснилось, вашим инвесторам показалось этого мало и они хотят, чтобы вы построили модели посерьезней и поточнее. Вы вспомнили, что когда то вы проходили курс машинного обучения и слышали что есть способ улучшить результаты вашей задачи: ансамбли: беггинг, пастинг, бустинг и стекинг, а также классификация путем жесткого и мягкого голосования и вы решили это опробовать. Требования к написанным классам вы оставляете теми же, что и в предыдущей работе. Будьте аккуратны в оптимизации целевой метрики и учитывайте несбалансированность классов. Требования к отчету сохраняются такими же.

1 Описание

Дерево принятия решений выглядит следующим образом:

```
1 class DecisionTree(BaseEstimator, ClassifierMixin):
2
3     def entropy(self, x):
4         return np.sum(-x * np.log2(x + (np.abs(x) < 1e-12)))
5
6     def __init__(self, leaf_size = 1, depth = 4, features = None):
7         self.leaf_size = leaf_size
8         self.depth = depth
9         self.features = features
10
11     class Node:
12         def __init__(self):
13             self.left = None
14             self.right = None
15             self.value = None
16             self.number = None
17
18     def fit(self, data, labels):
19         self.root = self.Node()
20         self.grow_tree(self.root, data, labels, np.arange(len(labels)), 0)
21         return self
22
23     def ans(self, y):
24         k = 0
25
26         for i in range(len(y)):
27             if y[i] == 1:
28                 k += 1
29         return k / len(y)
30
31
32     def grow_tree(self, node, data, labels, idx, depth):
33
34         x = data[idx]
35         y = labels[idx]
36
37         if (not self.depth is None) and (self.depth < depth):
38             node.value = self.ans(y)
39             return
40
41         if self.leaf_size >= len(x):
42             node.value = self.ans(y)
43             return
44
45         if (len(np.unique(y)) == 1):
46             node.value = self.ans(y)
```

```

47         return
48
49     maximum = -1
50     split_id = 0
51     leftidx = np.array([])
52     rightidx = np.array([])
53     index = 0
54
55     for i in (self.features if self.features is not None else range(data.shape[1]))
56         :
57         left = 1
58         sortid = x[:, i].argsort()
59         left_classes = np.zeros(2)
60         left_classes[y[sortid[0]]] = 1
61
62         classes = np.zeros(2)
63         for j in range(len(sortid)):
64             classes[y[sortid[j]]] += 1
65
66         xm = self.entropy(classes / len(x))
67
68         lenxm = len(sortid)
69
70         right_classes = np.zeros(2)
71
72         for j in range(left, len(sortid)):
73             right_classes[y[sortid[j]]] += 1
74
75         while left < len(x):
76
77             while left < len(x) and abs(x[sortid[left-1]][i] - x[sortid[left-2]][i])
78                 < 1e-6:
79                 left += 1
80                 left_classes[y[sortid[left - 1]]] += 1
81                 right_classes[y[sortid[left - 1]]] -= 1
82
83             if left == len(x):
84                 break
85
86             p = left_classes / left
87             x1 = self.entropy(p)
88             p = right_classes / (len(x) - left)
89             xr = self.entropy(p)
90
91             lenx1 = left
92             lenxr = len(x) - left
93
94             if maximum < xm * lenxm - lenx1 * x1 - lenxr * xr:

```

```

94         maximum = xm * lenxm - lenxl * xl - lenxr * xr
95         split_id = left - 1
96         leftidx = sortid[0:left]
97         rightidx = sortid[left:]
98         index = i
99
100         left += 1
101
102         left_classes[y[sortid[left - 1]]] += 1
103         right_classes[y[sortid[left - 1]]] -= 1
104
105
106     if len(leftidx) == 0 or len(rightidx) == 0:
107         node.value = self.ans(y)
108         return
109
110     node.number = index
111     node.value = x[sortid[split_id]][index]
112     node.left = self.Node()
113     node.right = self.Node()
114
115     self.grow_tree(node.left, x, y, leftidx, depth + 1)
116     self.grow_tree(node.right, x, y, rightidx, depth + 1)
117
118
119     def print_tree(node, d = 0):
120         if node is None:
121             return
122         print(' '*4*d, node.number, node.value)
123         print_tree(node.left, d+1)
124         print_tree(node.right, d+1)
125
126     def predict_proba(self, data):
127         res = np.ndarray((data.shape[0], 2))
128         for i, st in enumerate(data):
129             node = self.root
130             while not node.number is None:
131                 if st[node.number] > node.value:
132                     node = node.right
133                 else:
134                     node = node.left
135             res[i][1] = node.value
136             res[i][0] = 1 - node.value
137         return res
138
139     def predict(self, data):
140         return np.argmax(self.predict_proba(data), axis=1)

```

Для хранения вершин я использовал класс Node, который включает в себя левую

правую вершину, номер признака, по которому надо смотреть и значение для сравнения. В случае, когда вершина является листом, номер признака принимает значение None, а значение value ответом. Метод ans считает вероятность единички в выборке. Метод def grow_tree непосредственно строит дерево. Метод predict_proba считает вероятности (нужно для случайного леса), а predict непосредственно ответ. В качестве параметров выступают минимальное число элементов в листе и глубина дерева.

Случайный лес тоже построен в качестве класса. Выглядит он следующим образом:

```

1 class RandomForest(BaseEstimator, ClassifierMixin):
2
3     def __init__(self, features = 1, leaf_size = 1, depth = None, n_estimators = 10):
4         self.leaf_size = leaf_size
5         self.depth = depth
6         self.features = features
7         self.n_estimators = n_estimators
8
9     def fit(self, data, labels):
10        features = np.arange(data.shape[1])
11        self.estimators = []
12        indexes = np.arange(len(data))
13        for i in range(self.n_estimators):
14            np.random.shuffle(features)
15            self.estimators.append(DecisionTree(leaf_size = self.leaf_size, depth =
16                self.depth, features = features[:self.features]))
17            idx = np.random.choice(indexes, (len(data),))
18            self.estimators[-1].fit(data[idx], labels[idx])
19
20        def predict_proba(self, data):
21            pred = np.stack([est.predict_proba(data) for est in self.estimators], axis=1)
22            return pred.mean(axis=1)
23
24        def predict(self, data):
25            return self.predict_proba(data).argmax(axis=1)

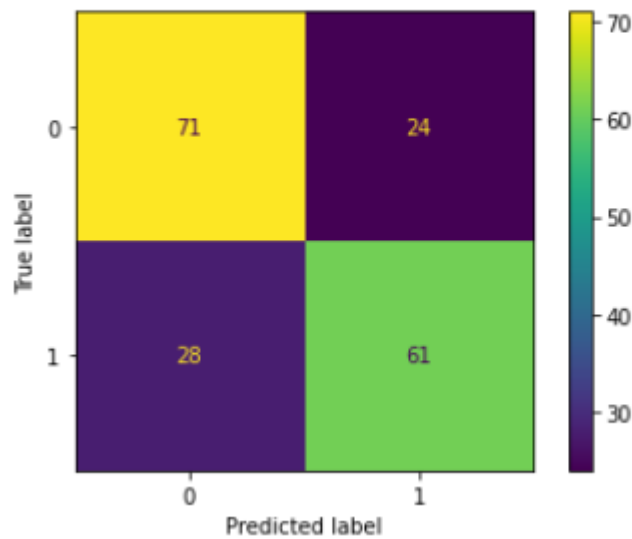
```

Здесь используется мягкое голосование. В качестве параметров выступают количество деревьев, количество фич, а также минимальное число элементов в листе и глубина деревьев.

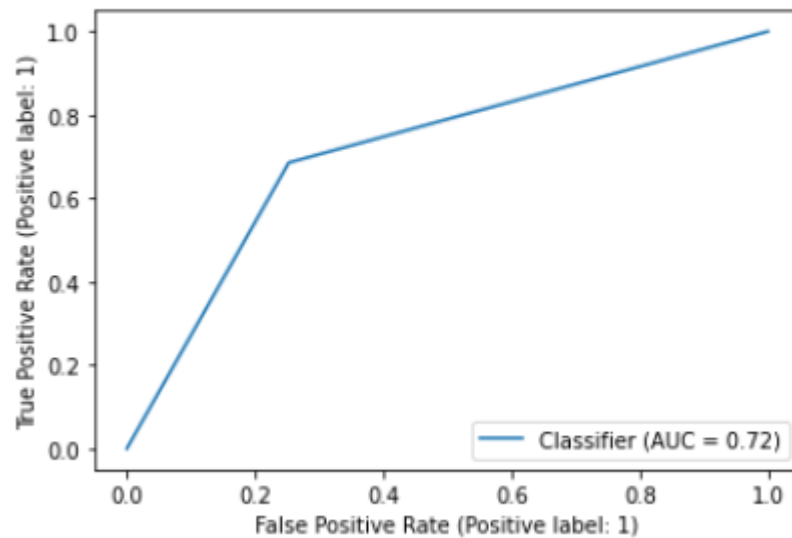
2 Результаты

Мои реализации отработали почему-то хуже, чем реализации sklearn. Для всех моделей я использовал кросс-валидацию. Мое дерево решений

```
{'depth': 10, 'leaf_size': 8}  
Accuracy train: 0.6104650079209766  
Accuracy tests: 0.717391304347826
```



```
Precision tests: 0.7176470588235294  
Recall tests: 0.6853932584269663
```

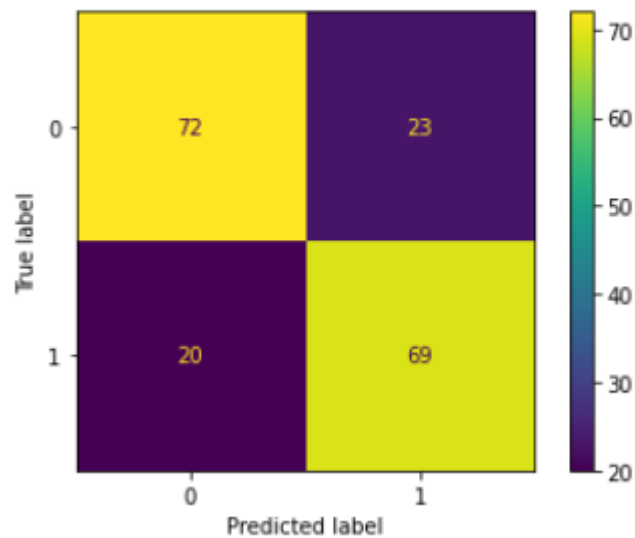


Дерево решений sklearn

```
{'max_depth': None, 'min_samples_leaf': 9}
```

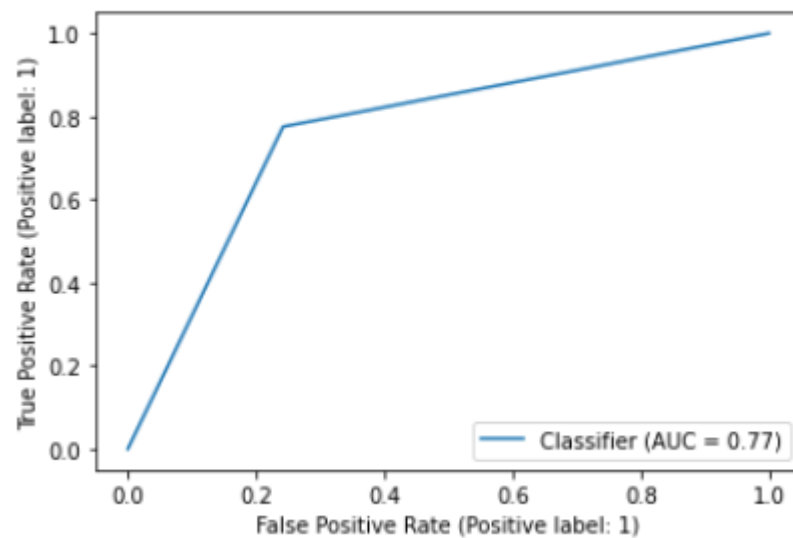
Accuracy train: 0.8432858074736744

Accuracy tests: 0.7663043478260869



Precision tests: 0.75

Recall tests: 0.7752808988764045

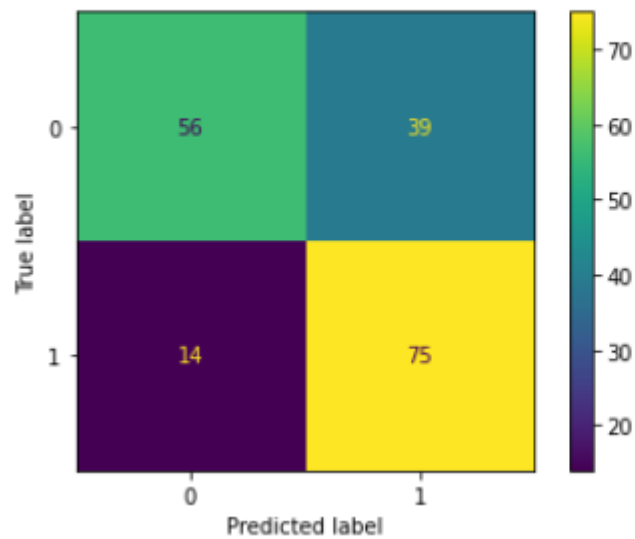


Мой случайный лес

```
{'depth': None, 'features': 3, 'leaf_size': 9, 'n_estimators': 30}
```

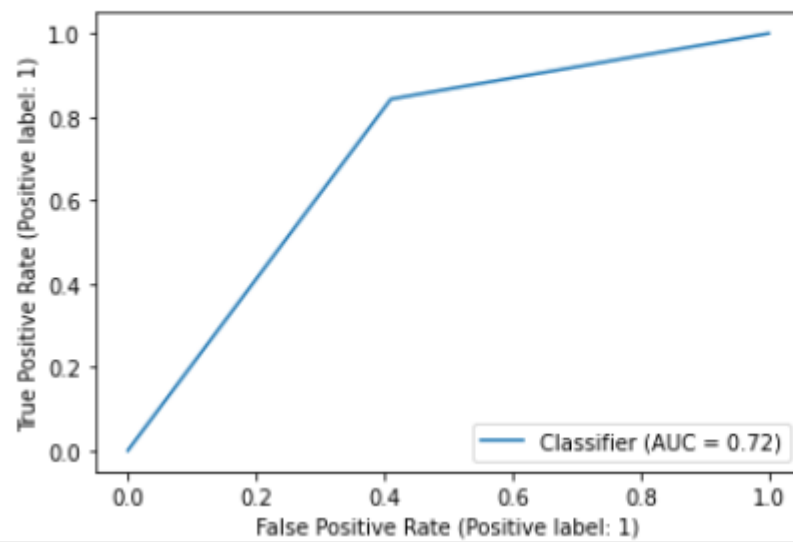
Accuracy train: 0.7737675892274718

Accuracy tests: 0.7119565217391305



Precision tests: 0.6578947368421053

Recall tests: 0.8426966292134831

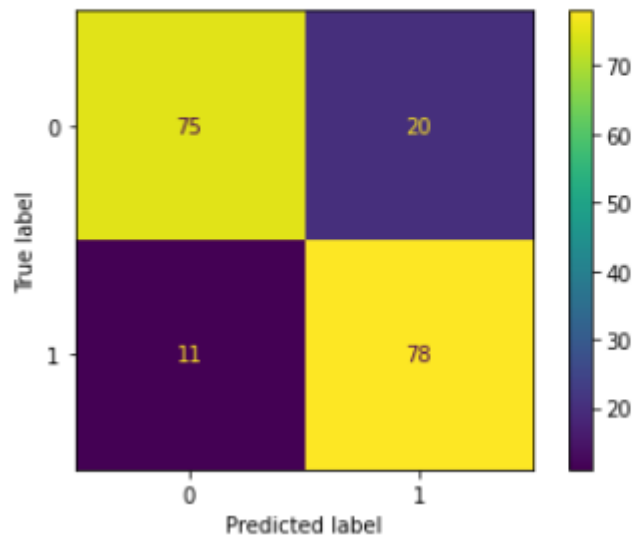


Случайный лес sklearn

```
{'max_features': None, 'min_samples_leaf': 3, 'n_estimators': 30}
```

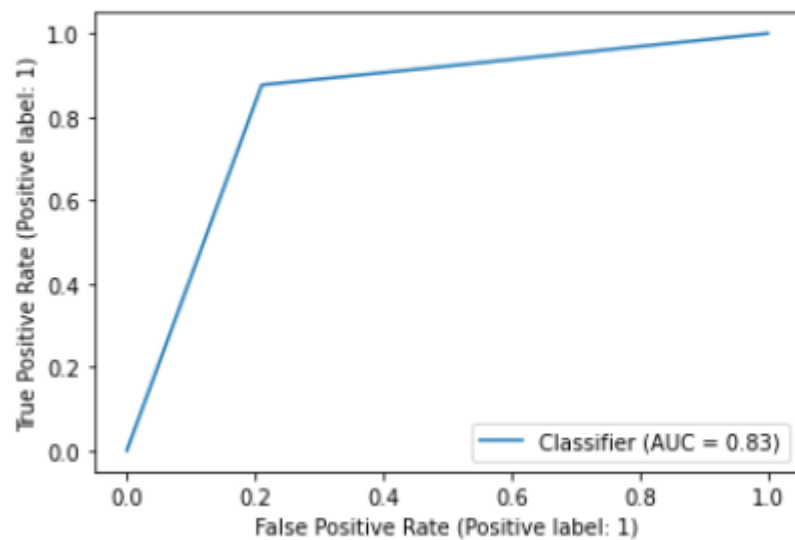
Accuracy train: 0.8678501537601342

Accuracy tests: 0.8315217391304348



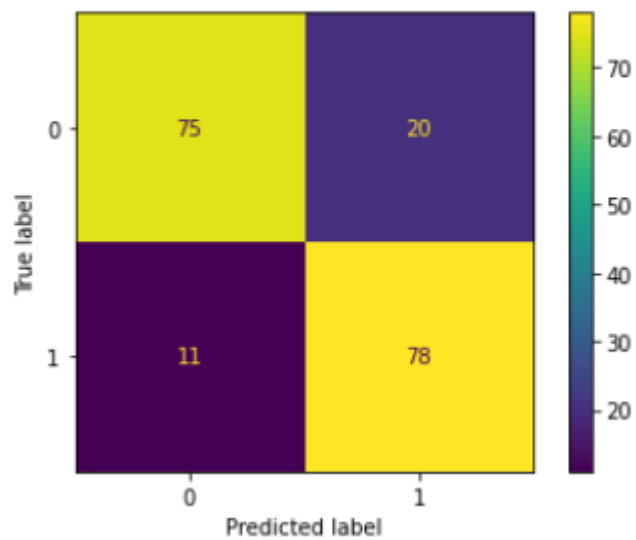
Precision tests: 0.7959183673469388

Recall tests: 0.8764044943820225

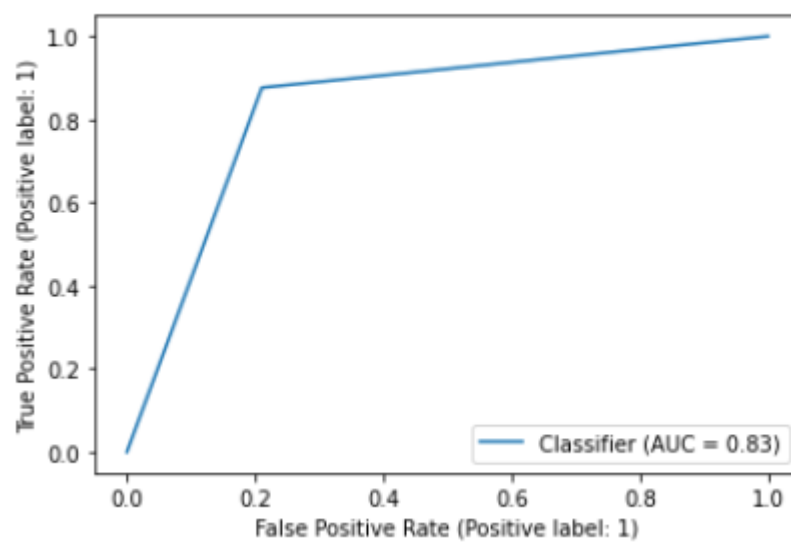


Также я воспользовался бустингом от sklearn. Он дал очень хорошие результаты, близкие к методу Байеса из прошлой лабораторной.

```
{'learning_rate': 0.1, 'n_estimators': 50}  
Accuracy train: 0.882816140154692  
Accuracy tests: 0.8315217391304348
```



```
Precision tests: 0.7959183673469388  
Recall tests: 0.8764044943820225
```



3 Выводы

Данная лабораторная работа была непростой. Дерево решений вызывало у меня трудности в написании. Я сталкивался с разными проблемами, например, деление категориальных признаков (они находятся в одной точке). Однако результаты вышли не совсем плохие. Точность в районе 70 % выглядит не так плохо, с учётом, что у меня почти всегда больше Recall, что для моей области это важнее, так как лучше определить здорового человека больным, чем наоборот. Бустинг, как уже говорил, дал хорошие результаты в 83 процента правильных ответов на тестовой выборке. В целом, если бы времени было бы больше, я бы попробовал большее количество моделей.