

**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: В. В. Бирюков
Преподаватели: Д. В. Сошников
С. Х. Ахмед
Группа: М8О-307Б-19
Дата:
Оценка:
Подпись:

Москва, 2022

Лабораторная работа №2

Задача: В лабораторной работе требуется:

1. Используя модели которые вы реализовали в предыдущей лабораторной работе, реализовать два подхода для построения ансамблей: жесткое и мягкое голосование, однако учтите, некоторые модели не предусматривают оценку вероятностей, например SVM и потому вам необходимо будет оценивать вероятности
2. Реализовать дерево решений
3. Реализовать случайный лес
4. Воспользоваться готовой коробочной реализацией градиентного бустинга для решения вашей задачи
5. Для всех моделей провести fine-tuning.

1 Ход работы

Решающее дерево

Решающее дерево хранится как множество связанных между собой узлов. Каждый внутренний узел содержит номер признака, по которому проводится разбиение, и его значение. Листья — вероятности классов. Построение дерева осуществляется рекурсивно методом `process_node`. Асимптотика построения — $O(D \cdot q \cdot N \log N)$ по времени и $O(DN)$ по памяти, где D — глубина дерева, q — количество признаков, N — количество данных.

```
1 class DecisionTree(BaseEstimator, ClassifierMixin):
2     class Node:
3         def __init__(self):
4             self.feature = -1
5             self.value = None
6             self.left = None
7             self.right = None
8             self.size = 0
9
10    def __init__(self, min_leaf_size=5, max_depth=None, criterion=entropy, features=
    None):
11        self.min_leaf_size = min_leaf_size
12        self.max_depth = max_depth
13        self.criterion = criterion
14        self.features = features
15
16    def fit(self, data, labels):
17        self.root = self.Node()
18        self.classes = len(np.unique(labels))
19        self.process_node(data, labels, self.root, np.arange(len(labels)), 0)
20        return self
21
22    def process_node(self, data, labels, node, ids, depth):
23        X = data[ids]
24        Y = labels[ids]
25        n = len(X)
26        values, c = np.unique(Y, return_counts=True)
27        counts = np.zeros((self.classes, ))
28        counts[values] += c
29        node.size = n
30
31        if (self.max_depth is not None) and depth == self.max_depth or \
32            (self.min_leaf_size is not None) and n <= self.min_leaf_size or \
33            len(values) == 1:
34            node.value = counts / n
35            return
36
```

```

37         h = self.criterion(Y)
38         max_value = None
39         max_f = None
40         max_gain = -1
41         best_left_ids = None
42         best_right_ids = None
43         for f in (self.features if self.features is not None else range(data.shape[1])
):
44             sort_ids = X[:, f].argsort()
45             left = 1
46             left_counts = np.zeros(2)
47             left_counts[Y[sort_ids[0]]] = 1
48
49             while left < n:
50                 while left < n and X[sort_ids[left-1]][f] == X[sort_ids[left-2]][f]:
51                     left += 1
52                     left_counts[Y[sort_ids[left-1]]] += 1
53                 if left == n:
54                     break
55
56                 p = left_counts / left
57                 left_h = self.criterion(p, from_proba=True)
58                 p = (counts - left_counts) / (n - left)
59                 right_h = self.criterion(p, from_proba=True)
60
61                 gain = h - (left * left_h + (n - left) * right_h) / n
62                 if gain > max_gain:
63                     max_gain = gain
64                     max_value = X[sort_ids[left-1]][f]
65                     max_f = f
66                     best_left_ids = sort_ids[:left]
67                     best_right_ids = sort_ids[left:]
68
69                 left += 1
70                 left_counts[Y[sort_ids[left-1]]] += 1
71
72             if max_value is None:
73                 node.value = counts / n
74                 return
75
76             node.feature = max_f
77             node.value = max_value
78             node.left = self.Node()
79             node.right = self.Node()
80
81             self.process_node(X, Y, node.left, best_left_ids, depth+1)
82             self.process_node(X, Y, node.right, best_right_ids, depth+1)
83
84     def predict_proba(self, data):

```

```

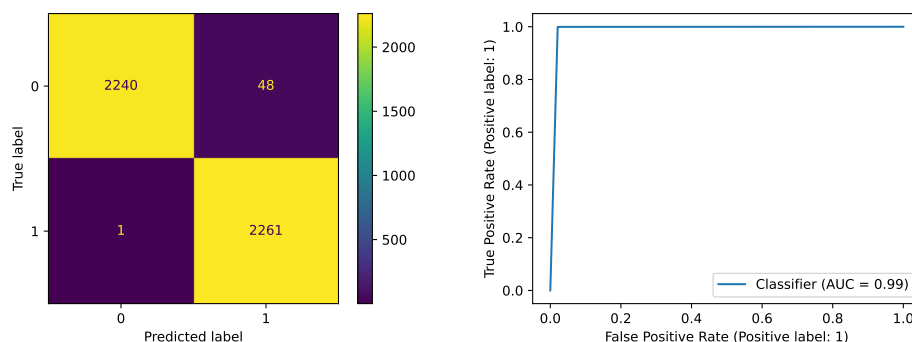
85     res = np.ndarray((data.shape[0], self.classes))
86     for i, obj in enumerate(data):
87         node = self.root
88         while node.feature != -1:
89             if obj[node.feature] > node.value:
90                 node = node.right
91             else:
92                 node = node.left
93         res[i] = node.value
94     return res
95
96 def predict(self, data):
97     return np.argmax(self.predict_proba(data), axis=1)

```

Оптимальным критерием информативности оказался критерий Джини, глубина дерева — 20, минимальный размер листа — 5.

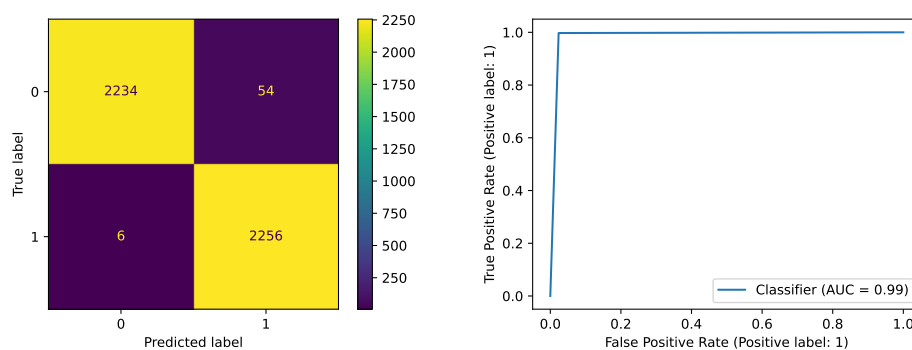
Результаты модели:

Accuracy: 0.9892307692307692
Precision: 0.9792117799913382
Recall: 0.9995579133510168

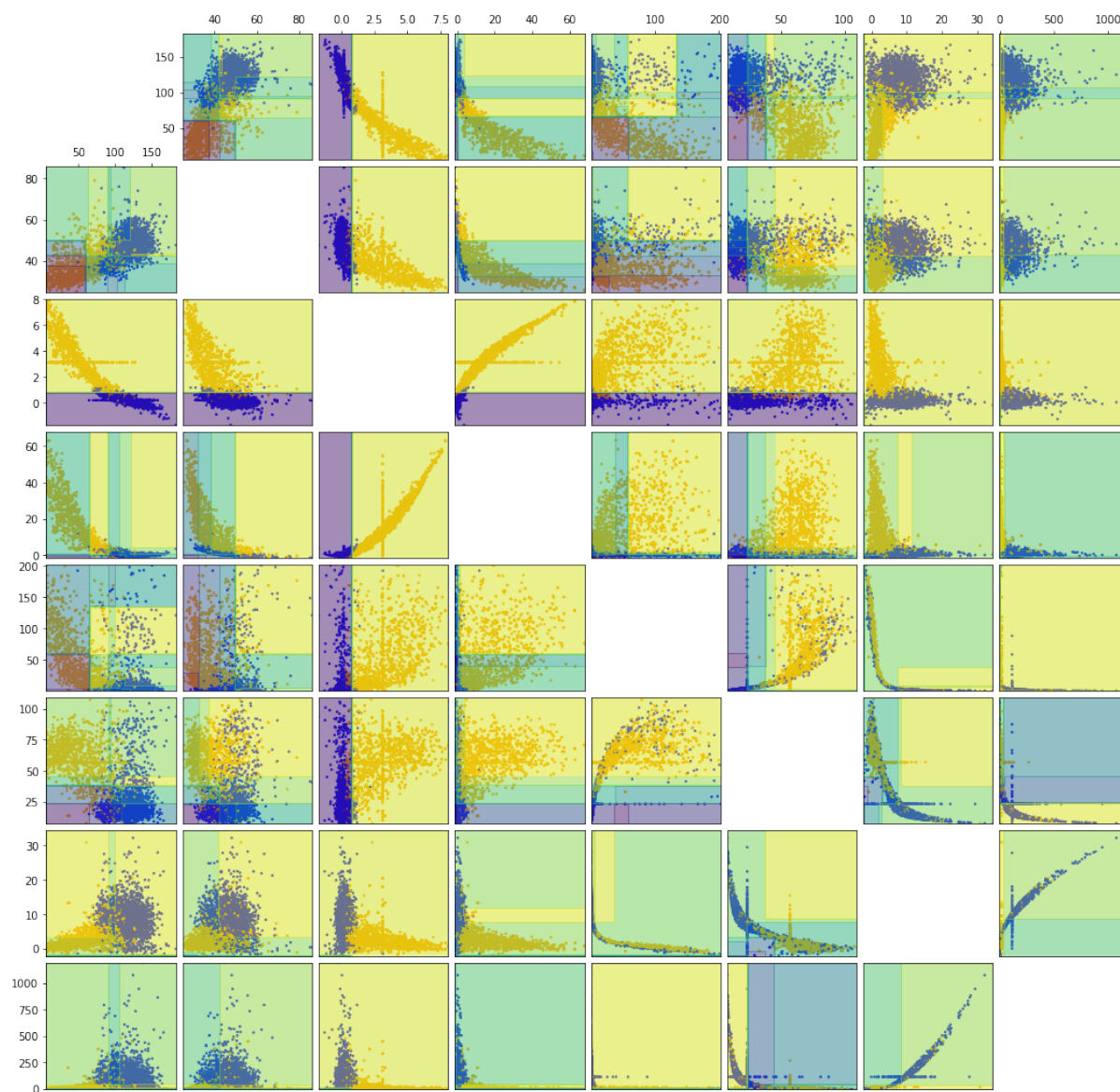


Готовый классификатор `DecisionTreeClassifier` показывает схожие результаты.

Accuracy: 0.9868131868131869
Precision: 0.9766233766233766
Recall: 0.9973474801061007



Визуализируем разделяющую поверхность решающего дерева.



Так как предсказания делаются всего по двум признакам, в узлах дерева с другими признаками ответ из обоих потомков усредняется. Поэтому это не отражает совсем полной картины разделения. Тем не менее, видно, что большая часть разделения происходит из-за признака, по которому классы очень хорошо разделяются одной точкой.

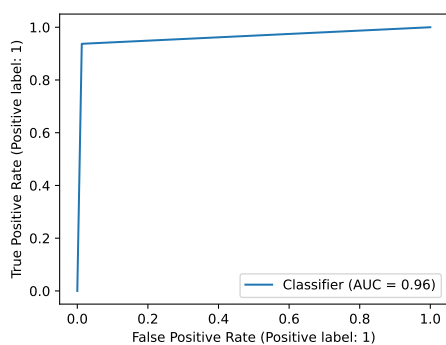
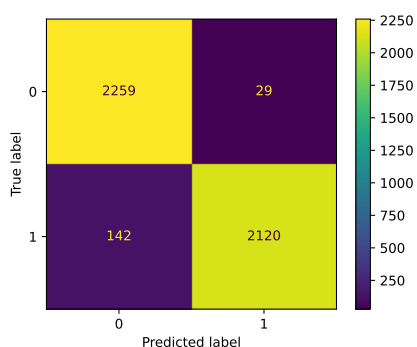
Голосование

Жесткое и мягкое голосование реализовано в рамках одного класса. Мягкое голосование усредняет предсказанные вероятности. Жесткое интерпретирует как вероятности доли моделей, предсказавших этот класс.

```
1 class Voting(BaseEstimator, ClassifierMixin):
2     def __init__(self, estimators, mode='soft', pretrained=False):
3         self.estimators = estimators
4         self.mode = mode
5         self.pretrained = pretrained
6
7     def fit(self, data, labels):
8         self.classes = len(np.unique(labels))
9         if not self.pretrained:
10             for estimator in self.estimators:
11                 estimator.fit(data, labels)
12             return self
13
14     def predict(self, data):
15         return self.predict_proba(data).argmax(axis=1)
16
17     def predict_proba(self, data):
18         if self.mode == 'hard':
19             pred = np.stack([est.predict(data) for est in self.estimators], axis=-1)
20             res = np.zeros((len(data), self.classes))
21             for i, p in enumerate(pred):
22                 v, c = np.unique(p, return_counts=True)
23                 res[i][v] += c
24                 res[i] /= len(self.estimators)
25             return res
26         else:
27             pred = np.stack([est.predict_proba(data) for est in self.estimators], axis
28                             =1)
29             return pred.mean(axis=1)
```

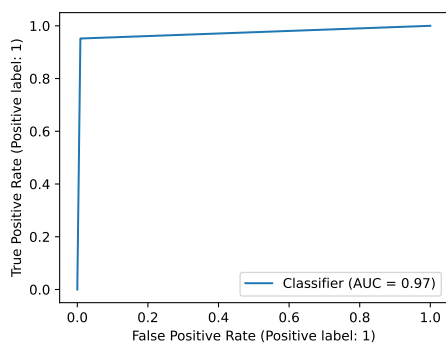
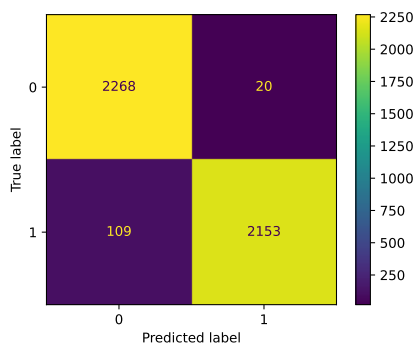

Объединение KNN, байесовского классификатора, логистической регрессии, SVM и решающего дерева. Жесткое голосование.

Accuracy: 0.9624175824175825
Precision: 0.9865053513261982
Recall: 0.9372236958443855



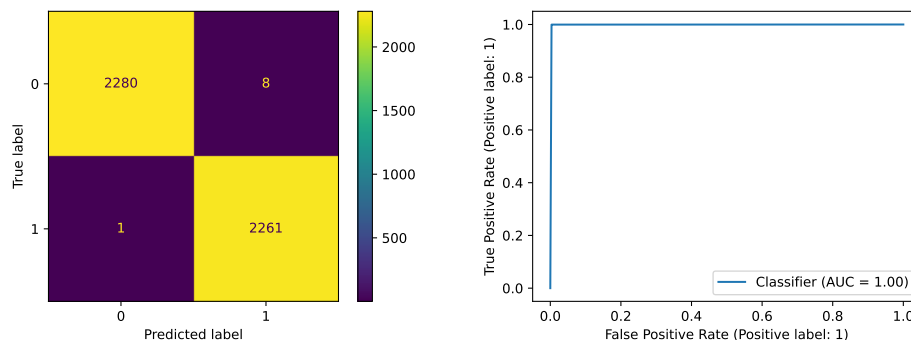
Мягкое голосование.

Accuracy: 0.9716483516483516
Precision: 0.9907961343764381
Recall: 0.9518125552608311



Наилучшие результаты показывает мягкое объединение двух самых лучших классификаторов — KNN и решающего дерева.

Accuracy: 0.998021978021978
Precision: 0.996474217717056
Recall: 0.9995579133510168



Случайный лес

Случайный лес использует множество решающих деревьев, каждое из которых обучается на случайной (с повторениями) подвыборке обучающих данных и подмножестве признаков. Предсказания деревьев объединяются по принципу мягкого голосования.

```

1 class RandomForest(BaseEstimator, ClassifierMixin):
2     def __init__(self, n_estimators=100, min_leaf_size=5, max_depth=None, criterion=
    entropy, max_features='sqrt', max_samples=0.8):
3         self.n_estimators = n_estimators
4         self.min_leaf_size = min_leaf_size
5         self.max_depth = max_depth
6         self.criterion = criterion
7         self.max_features = max_features
8         self.max_samples = max_samples
9
10    def fit(self, data, labels):
11        features = np.arange(data.shape[1])
12        indexes = np.arange(len(data))
13        samples = math.floor(self.max_samples * len(data))
14        if self.max_features == 'sqrt':
15            max_features = math.floor(np.sqrt(len(features)))
16        else:
17            max_features = math.floor(len(features) * self.max_features)
18        self.estimators = []
19        for _ in range(self.n_estimators):
20            np.random.shuffle(features)
21            self.estimators.append(DecisionTree(min_leaf_size=self.min_leaf_size,
    max_depth=self.max_depth, criterion=self.criterion, features=features[:
    max_features]))
22            idx = np.random.choice(indexes, (samples, ))
23            self.estimators[-1].fit(data[idx], labels[idx])
24
25    def predict_proba(self, data):

```

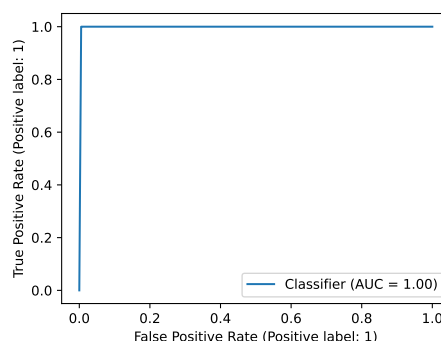
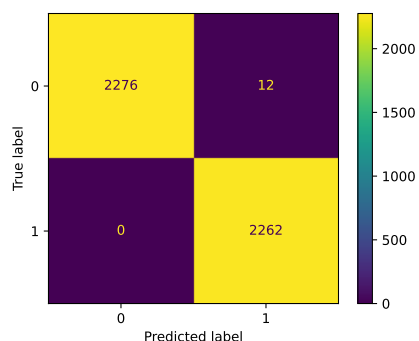
```

26 |         pred = np.stack([est.predict_proba(data) for est in self.estimators], axis=1)
27 |         return pred.mean(axis=1)
28 |
29 |     def predict(self, data):
30 |         return self.predict_proba(data).argmax(axis=1)

```

Результаты модели с 30 деревьями:

Accuracy: 0.9973626373626374
Precision: 0.9947229551451188
Recall: 1.0

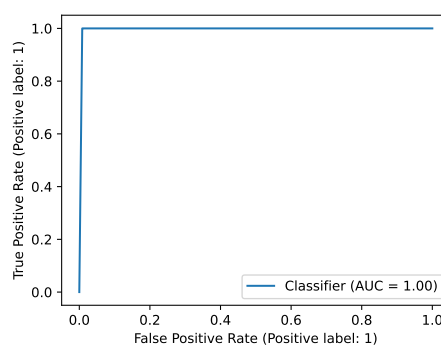
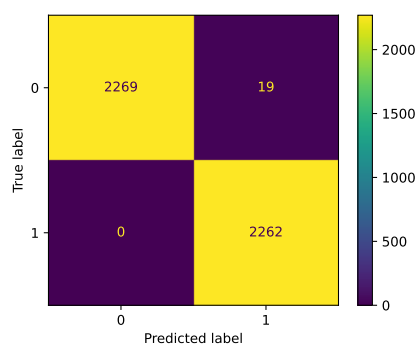


Градиентный бустинг

В качестве градиентного бустинга испытаны готовые реализации из библиотек `sklearn` и `catboost`.

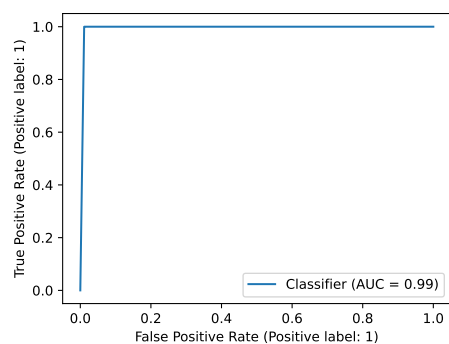
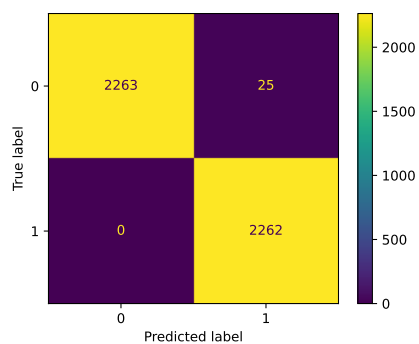
`GradientBoostingClassifier`:

Accuracy: 0.9958241758241758
Precision: 0.9916703200350724
Recall: 1.0



CatBoostClassifier:

Accuracy: 0.9945054945054945
Precision: 0.9890686488850022
Recall: 1.0



2 Выводы

В ходе выполнения лабораторной работы я познакомился с ансамблями и решающими деревьями. Решающее дерево оказалось довольно хорошей моделью для моей задачи. Объединение же самых лучших моделей в одну привело к самому лучшему результату за обе лабораторных работы.