

# Машинное обучение, DS-поток

## Домашнее задание 3А

### Правила:

- Дедлайн **09 марта 23:59**. После дедлайна работы не принимаются кроме случаев наличия уважительной причины.
- Выполненную работу нужно отправить на почту `mipt.stats@yandex.ru`, указав тему письма "[ml] Фамилия Имя - задание 3А". Квадратные скобки обязательны. Если письмо дошло, придет ответ от автоответчика.
- Прислать нужно ноутбук и его pdf-версию (без архивов). Названия файлов должны быть такими: 3А.N.ipynb и 3А.N.pdf, где N - ваш номер из таблицы с оценками.
- Решения, размещенные на каких-либо интернет-ресурсах не принимаются. Кроме того, публикация решения в открытом доступе может быть приравнена к предоставлению возможности списать.
- Для выполнения задания используйте этот ноутбук в качестве основы, ничего не удаляя из него.
- Никакой код из данного задания при проверке запускаться не будет.

Задание стоит 15 баллов.

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.base import BaseEstimator
4 from sklearn.metrics import accuracy_score, r2_score
5 from sklearn.model_selection import GridSearchCV, train_test_split
6 import matplotlib.pyplot as plt
7 import warnings
8 import seaborn as sns
9
10 warnings.filterwarnings('ignore')
11 sns.set(font_scale=1.6)
12 plt.rcParams['axes.facecolor'] = 'lightgrey'
```

Вспомним, как именно происходит построение решающего дерева. Для построения дерева в каждой нелистовой вершине происходит разбиение подвыборки на две части по некоторому признаку  $x_j$ . Этот признак и порог  $t$ , по которому будет происходить разбиение, мы хотим брать не произвольно, а основываясь на соображениях оптимальности. Для этого нам необходимо знать некоторый функционал качества, который будем оптимизировать при построении разбиения.

Обозначим через  $X_m$  -- множество объектов, попавших в вершину  $m$ , разбиваемую на данном шаге, а через  $X_l$  и  $X_r$  -- объекты, попадающие в левое и правое поддерево соответственно при заданном правиле  $I\{x_j < t\}$ . Пусть также  $H$  -- используемый критерий информативности (impurity criterion).

Выпишите функционал, который необходимо минимизировать при разбиении вершины:

Ответ: <...>

Реализация критериев информативности.

Вспомните еще раз, на какой общей идее основаны критерии информативности и какую характеристику выборки они стремятся оптимизировать?

Ответ: <...>

Перед тем, как непосредственно работать с решающими деревьями, реализуйте функции подсчёта значения критериев разбиения вершин решающих деревьев. Использовать готовые реализации критериев или классов для решающих деревьев из `sklearn` и из других библиотек **запрещено**. Также при реализации критериев по причине неэффективности **запрещается использовать циклы**. Воспользуйтесь библиотекой `numpy`.

Каждая функция принимает на вход одномерный `numpy`-массив размерности  $(n,)$  из значений отклика.

In [ ]:

```
1 # Код функций, реализующих критерии разбиения.
2
3 def mean_square_criterion(y):
4     ''' Критерий для квадратичной функции потерь. '''
5
6     return <...>
7
8
9 def mean_abs_criterion(y):
10     ''' Критерий для абсолютной функции потерь. '''
11
12     return <...>
13
14
15 def get_probs_by_y(y):
16     ''' Возвращает вектор частот для каждого класса выборки. '''
17
18     return <...>
19
20
21 def gini_criterion(y):
22     ''' Критерий Джини. '''
23
24     return <...>
25
26
27 def entropy_criterion(y):
28     ''' Энтропийный критерий. '''
29
30     return <...>
```

Протестируйте реализованные функции.

Тесты для распределения вероятностей на классах.

In [ ]:

```
1 assert np.allclose(get_probs_by_y([1, 1, 2, 2, 7]), np.array([0.4, 0.4, 0.2]))
2 assert np.allclose(get_probs_by_y([1]), np.array([1]))
```

Тесты для критериев разбиения.

In [ ]:

```
1 assert np.allclose(entropy_criterion([25]), 0)
2 assert np.allclose(gini_criterion([25]), 0)
3 assert np.allclose(mean_square_criterion([10, 10, 10]), 0)
4 assert np.allclose(mean_abs_criterion([10, 10, 10]), 0)
```

### Реализация класса решающего дерева.

Для того, чтобы лучше понять, как устроены решающие деревья и как именно устроен процесс их построения, вам предлагается реализовать класс `BaseDecisionTree`, реализующий базовые функции решающего дерева. Большая часть кода уже написана.

Используются следующие классы:

**Класс** `BaseDecisionTree` - класс для решающего дерева, в котором реализовано построение дерева. Все вершины дерева хранятся в списке `self.nodes`, при этом вершина с номером 0 - корень.

1) `__init__` - инициализация дерева. Здесь сохраняются гиперпараметры дерева: `criterion`, `max_depth`, `min_samples_split` и инициализируется список вершин, состоящий только из одной вершины - корневой,

2) `build_` - рекурсивная функция построения дерева. В ней при посещении каждой вершины дерева проверяются условия, стоит ли продолжать разбивать эту вершину. Если да, то перебираются все возможные признаки и пороговые значения и выбирается та пара (признак, значение), которой соответствует наименьшее значение критерия информативности,

3) `fit` - функция обучения дерева, принимающая на вход обучающую выборку. В этой функции происходит предподсчёт всех возможных пороговых значений для каждого из признаков, а затем вызывается функция `build_`.

**Класс** `Node` - класс вершины дерева. Внутри вершины, помимо разделяющего признака и порога хранятся `self.left_son`, `self.right_son` - номера дочерних вершин, а также `self.left_prob` и `self.right_prob` - вероятности попадания элемента в каждую из них. При этом в листовых вершинах хранятся также `self.y_values` - значения соответствующих элементов выборки, попавших в вершину.

1) `__init__` - инициализация вершины. Принимает в качестве аргументов разделяющий признак и пороговое значение и сохраняет их.

**Класс** `DecisionTreeRegressor` - наследник класса `BaseDecisionTree`, в котором реализованы функции для предсказаний при решении задачи регрессии.

1) `predict_instance` - получение предсказания для одного элемента выборки. Выполняется посредством спуска по решающему дереву до листовой вершины,

2) `predict` - получение предсказаний для всех элементов выборки.

**Класс** `DecisionTreeClassifier` - наследник класса `BaseDecisionTree`, в котором реализованы функции для предсказаний при решении задачи классификации.

1) `predict_proba_instance` - предсказание распределения вероятностей по классам для одного элемента выборки,

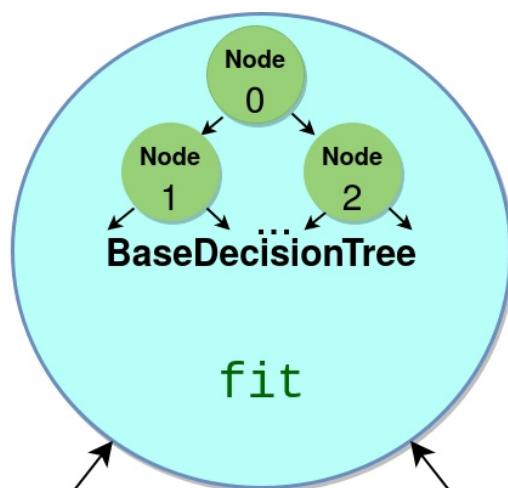
2) `predict_proba` - предсказание распределения вероятностей по классам для всех элементов выборки,

3) `predict_proba` - предсказание меток классов для всех элементов выборки.

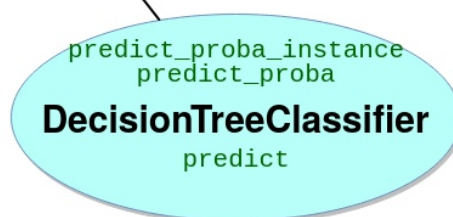
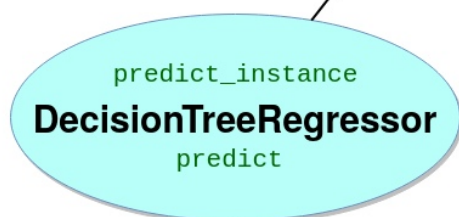
Перед написанием кода разбиения дерева, ответьте на вопрос, какие пороговые значения для каждого из признаков вы будете перебирать. Почему рассматривать другие значения в качестве пороговых не имеет смысла?

Ответ: <...>

Структура решающего дерева



Наследуют от класса `BaseDecisionTree` функцию инициализации и механизм построения дерева



In [ ]:

```
1 def get_not_nans(arr):
2     '''
3     Функция, которая создаёт и возвращает новый массив
4     из всех элементов переданного массива, не являющихся None.
5     '''
6
7     return arr.copy()[arr != None]
8
9
10 class Node(object):
11     def __init__(self, split_feature=None, split_threshold=None):
12         '''
13         Функция инициализации вершины решающего дерева.
14
15         Параметры.
16         1) split_feature - номер разделяющего признака
17         2) split_threshold - пороговое значение
18         '''
19
20         self.split_feature = split_feature
21         self.split_threshold = split_threshold
22         # По умолчанию считаем, что у вершины нет дочерних вершин.
23         self.left_son, self.right_son = None, None
24         # Вероятности попадания в каждую из дочерних вершин нужно поддерживать
25         # для корректной обработки данных с пропусками
26         self.left_prob, self.right_prob = 0, 0
27         # Массив значений y. Определён только для листовых вершин дерева
28         self.y_values = None
29
30
31 class BaseDecisionTree(BaseEstimator):
32     '''
33     Здесь содержится реализация всех основных функций для работы
34     с решающим деревом.
35
36     Наследование от класса BaseEstimator нужно для того, чтобы
37     в дальнейшем данный класс можно было использовать в
38     различных функциях библиотеки sklearn, например, в функциях
39     для кросс-валидации.
40     '''
41
42     def __init__(self, criterion, max_depth=np.inf, min_samples_split=2):
43         '''
44         Функция инициализации решающего дерева.
45
46         Параметры.
47         1) criterion - критерий информативности,
48         2) max_depth - максимальная глубина дерева,
49         3) min_samples_split - минимальное количество элементов
50         обучающей выборки, которое должно попасть в вершину,
51         чтобы потом происходило разбиение этой вершины.
52         '''
53
54         self.criterion = criterion
55         self.max_depth = max_depth
56         self.min_samples_split = min_samples_split
57         # Список всех вершин дерева. В самом начале
58         # работы алгоритма есть только одна
59         # вершина - корень.
```

```

60     self.nodes = [Node()]
61     # Количество классов. Актуально только
62     # при решении задачи классификации.
63     self.class_count = 1
64     # Сюда нужно будет записать все значения
65     # для каждого из признаков датасета.
66     self.feature_values = None
67
68 def build_(self, v, X, y, depth=0):
69     '''
70     Рекурсивная функция построения решающего дерева.
71
72     Параметры.
73     1) v - номер рассматриваемой вершины
74     2) X, y - обучающая выборка, попавшая в текущую вершину
75     3) depth - глубина вершины с номером v
76     '''
77
78     if depth == self.max_depth or len(y) < self.min_samples_split:
79         # Если строим дерево для классификации, то
80         # сохраняем метки классов всех элементов выборки,
81         # попавших в вершину.
82         if callable(getattr(self, "set_class_count", None)):
83             self.nodes[v].y_values = y.copy()
84             # Для регрессии сразу вычислим среднее всех
85             # элементов вершины.
86         else:
87             self.nodes[v].y_values = np.mean(y)
88         return
89
90     best_criterion_value = np.inf
91     best_feature, best_threshold = 0, 0
92     sample_size, feature_count = X.shape
93
94     # переберём все возможные признаки и значения порогов,
95     # найдём оптимальный признак и значение порога
96     # и запишем их в best_feature, best_threshold
97     for feature_id in range(feature_count):
98         for threshold in self.feature_values[feature_id]:
99             <...>
100
101     # сохраним найденные параметры в класс текущей вершины
102     self.nodes[v].split_feature = <...>
103     self.nodes[v].split_threshold = <...>
104     # разделим выборку на 2 части по порогу
105     <...>
106
107     # создаём левую и правую дочерние вершины,
108     # и кладём их в массив self.nodes
109     self.nodes.append(Node())
110     self.nodes.append(Node())
111     # сохраняем индексы созданных вершин в качестве
112     # левого и правого сына вершины v
113     self.nodes[v].left_son, self.nodes[v].right_son = \
114         len(self.nodes)-2, len(self.nodes)-1
115     # рекурсивно строим дерево для дочерних вершин
116     self.build_(self.nodes[v].left_son, X_l, y_l, depth+1)
117     self.build_(self.nodes[v].right_son, X_r, y_r, depth+1)
118
119 def fit(self, X, y):
120     '''

```

```

121     Функция, из которой запускается построение
122     решающего дерева по обучающей выборке.
123
124     Параметры.
125     X, y - обучающая выборка
126     ...
127
128     # сохраним заранее все пороги для каждого из
129     # признаков обучающей выборки
130     X, y = np.array(X), np.array(y)
131     self.feature_values = []
132     for feature_id in range(X.shape[1]):
133         self.feature_values.append(
134             np.unique(get_not_nans(X[:, feature_id]))
135         )
136
137     set_class_count = getattr(self, "set_class_count", None)
138     # если строится дерево для классификации,
139     # то нужно посчитать количество классов
140     if callable(set_class_count):
141         set_class_count(y)
142     self.build_(0, X, y)

```

Теперь, когда общий код решающего дерева написан, нужно сделать обёртки над `BaseDecisionTree` - классы `DecisionTreeRegressor` и `DecisionTreeClassifier` для использования решающего дерева в задачах регрессии и классификации соответственно.

Допишите функции `predict_instance` и `predict_proba_instance` в классах для регрессии и классификации соответственно. В этих функциях нужно для одного элемента  $x$  выборки промоделировать спуск в решающем дереве, а затем по листовой вершине, в которой окажется объект, посчитать для классификации - распределение вероятностей, а для регрессии - число  $y$ .

In [ ]:

```
1 class DecisionTreeRegressor(BaseDecisionTree):
2     def predict_instance(self, x, v):
3         '''
4         Рекурсивная функция, предсказывающая значение
5         у для одного элемента x из выборки.
6
7         Параметры.
8         1) x - элемент выборки, для которого
9         требуется предсказать значение y
10        2) v - рассматриваемая вершина дерева
11        '''
12
13        # если вершина - лист, возвращаем в качестве предсказания
14        # среднее всех элементов, содержащихся в ней
15        if self.nodes[v].left_son == None:
16            <...>
17
18        # если у объекта x значение признака по
19        # которому происходит разделение, меньше
20        # порогового, то спускаемся в левое поддерево,
21        # иначе - в правое
22        if x[self.nodes[v].split_feature] < self.nodes[v].split_threshold:
23            return <...>
24        elif x[self.nodes[v].split_feature] >= self.nodes[v].split_threshold:
25            return <...>
26        # а если у элемента отсутствует значение
27        # разделяющего признака, то будем спускаться
28        # в оба поддерева
29        else:
30            left_predict = self.predict_instance(x, self.nodes[v].left_son)
31            right_predict = self.predict_instance(x, self.nodes[v].right_son)
32            return <...>
33
34    def predict(self, X):
35        '''
36        Функция, предсказывающая значение
37        у для всех элементов выборки X.
38
39        Параметры.
40        X - выборка, для которой требуется
41        получить вектор предсказаний y
42        '''
43
44        return [self.predict_instance(x, 0) for x in X]
```

Для удобства реализации функции `predict_proba_instance` класса `DecisionTreeClassifier` будем считать, что все классы имеют целочисленные метки от 0 до  $k - 1$ , где  $k$  - количество классов. Если бы это условие не было выполнено, то нужно было бы сначала сделать предобработку меток классов в датасете.



In [ ]:

```
1 class DecisionTreeClassifier(BaseDecisionTree):
2     def set_class_count(self, y):
3         '''
4         Функция, вычисляющая количество классов
5         в обучающей выборке.
6
7         Параметры.
8         y - значения класса в обучающей выборке
9         '''
10
11         self.class_count = np.max(y) + 1
12
13     def predict_proba_instance(self, x, v):
14         '''
15         Рекурсивная функция, предсказывающая вектор
16         вероятностей принадлежности объекта x
17         к классам
18
19         Параметры.
20         1) x - элемент выборки, для которого
21         требуется предсказать значение y
22         2) v - вершина дерева, в которой
23         находится алгоритм
24         '''
25
26         if self.nodes[v].left_son == None:
27             result = np.zeros(self.class_count)
28             classes, counts = np.unique(
29                 self.nodes[v].y_values, return_counts=True)
30             result[classes.astype(int)] = counts
31             return result / np.sum(result)
32
33         # если у объекта x значение признака по которому
34         # происходит разделение, меньше порогового,
35         # то спускаемся в левое поддерево, иначе - в правое
36         if x[self.nodes[v].split_feature] < self.nodes[v].split_threshold:
37             return <...>
38         elif x[self.nodes[v].split_feature] >= self.nodes[v].split_threshold:
39             return <...>
40         # а если у объекта отсутствует значение
41         # разделяющего признака, то будем спускаться
42         # в оба поддерева
43         else:
44             left_predict = self.predict_proba_instance(
45                 x, self.nodes[v].left_son)
46             right_predict = self.predict_proba_instance(
47                 x, self.nodes[v].right_son)
48             return <...>
49
50     def predict_proba(self, X):
51         '''
52         Функция, предсказывающая вектор вероятностей
53         принадлежности объекта x к классам для
54         каждого x из X
55
56         Параметры.
57         X - выборка, для которой требуется получить вектор предсказаний y
58         '''
59
```

```

60         return [self.predict_proba_instance(x, 0) for x in X]
61
62     def predict(self, X):
63         '''
64         Функция, предсказывающая метку класса для
65         всех элементов выборки X.
66
67         Параметры.
68         X - выборка, для которой требуется получить
69         вектор предсказаний y
70         '''
71
72         return np.argmax(self.predict_proba(X), axis=1)

```

## Подбор параметров.

В этой части задания вам предлагается поработать с написанным решающим деревом, применив его к задаче классификации и регрессии, и в обеих задачах подобрать оптимальные параметры для построения.

Не забывайте писать выводы.

### 1. Задача классификации.

Теперь - самое время протестировать работу написанного нами решающего дерева. Делать мы это будем на датасете для классификации вина из `sklearn`.

In [ ]:

```

1 from sklearn.datasets import load_wine
2
3 X, y = load_wine(return_X_y=True)

```

Далее для критерия Джини и энтропийного критерия найдем оптимальные параметры обучения дерева - `max_depth` и `min_samples_split`.

In [ ]:

```

1 classification_criteria = <...>
2 criterion_names = <...>

```

С начала надо разбить выборку на train и test.

In [ ]:

```

1 X_train, X_test, y_train, y_test = <...>

```

Теперь проведите кросс-валидацию для каждого из критериев разбиения вершин.

In [ ]:

```

1 for criterion, criterion_name in zip(classification_criteria, criterion_names):
2     accuracy = <...>
3     print('accuracy:', accuracy)
4     assert(accuracy >= 0.85, "Something is wrong with your classifier")

```

## Построение графиков.

Постройте графики зависимости ассигасы от максимальной глубины дерева на обучающей и тестовой выборке для каждого критерия на train и на test. В качестве максимальной глубины используйте значения от 3 до 7. Значение `min_samples_split` фиксируйте.

In [ ]:

```
1 # код построения графиков
2 <...>
```

Сделайте выводы. Почему графики получились такими? Как соотносятся оптимальные значения параметров на обучающей и на тестовой выборках?

## Вывод.

<...>

## 2. Задача регрессии.

Прodelайте аналогичные шаги для задачи регрессии. В качестве датасете возьмите `boston` из `sklearn`, а в качестве критерия качества возьмите `r2_score`. Рассмотрим более широкий диапазон значений для `max_depth`: от 3 до 14.

In [ ]:

```
1 from sklearn.datasets import load_boston
2
3 boston_X, boston_y = load_boston(return_X_y=True)
```

In [ ]:

```
1 regression_criteria = <...>
2 criterion_names = <...>
```

In [ ]:

```
1 <...>
```

Разобьём выборку на обучение и тест.

In [ ]:

```
1 X_train, X_test, y_train, y_test = <...>
```

Проведите эксперименты, аналогичны тем, что были сделаны для задачи классификации.

In [ ]:

```
1 <...>
```

Сделайте вывод, в котором объясните, почему графики получились такими.

Скорее всего, вы заметили, что дерево в этих экспериментах строится довольно медленно. Как можно ускорить его построение? Можно ли ускорить нахождение оптимального разбиения по некоторому вещественному признаку?

**Вывод.**

<...>

## Обработка пропусков с использованием решающих деревьев.

А теперь рассмотрим датасет, в котором часть данных пропущена. В качестве примера возьмём датасет <https://archive.ics.uci.edu/ml/datasets/Adult> (<https://archive.ics.uci.edu/ml/datasets/Adult>) для определения категории дохода работников, по таким признакам, как возраст, образование, специальность, класс работы, пол, кол-во отработываемых часов в неделю и некоторым другим.

In [2]:

```
1 column_names = [  
2     'age', 'workclass', 'fnlwgt', 'education1', 'education2', 'marital-status',  
3     'occupation', 'relationship', 'race', 'sex', 'capital-gain',  
4     'capital-loss', 'hours-per-week', 'native-country', 'target'  
5 ]
```

Поскольку предсказание в дереве на данных с пропусками часто занимает сильно больше времени, чем в случае отсутствия пропусков (так как часто приходится спускаться разу в 2 поддерева), то для экономии времени сократим датасет, взяв из него только первые 10000 строк данных.

In [ ]:

```
1 adult_df = pd.read_csv('adult.data', header=None)[:10000]  
2 adult_df.columns = column_names  
3 target = adult_df['target'] == '>50K'  
4 adult_df = adult_df.drop(['target'], axis=1)  
5 adult_df.head()
```

Предобработаем датасет, заменив категориальные признаки one-hot векторами.

In [ ]:

```
1 adult_df = pd.get_dummies(adult_df)  
2 adult_df.head()
```

Поскольку все пропущенные значения относились к категориальным признакам и помечались в датасете знаком `?`, то для каждого категориального признака `feature` исходного датасета надо выполнить следующую процедуру: рассмотреть признак `feature_?` нового датасета и для всех строк, для которых выполнено `feature_?=1`, значениях всех признаков с префиксом `feature` установить в `None`.

In [ ]:

```
1 all_indices = np.arange(adult_df.shape[0])
2
3 for feature in column_names:
4     if f'{feature}_?' in adult_df.columns:
5         none_indices = all_indices[adult_df[f'{feature}_?'] == 1]
6
7         for dummy_feature in adult_df.columns:
8             if not dummy_feature.startswith(f'{feature}_ '):
9                 continue
10            if dummy_feature != f'{feature}_?':
11                adult_df[dummy_feature][none_indices] = None
12        adult_df = adult_df.drop(f'{feature}_?', axis=1)
```

Посмотрим на распределение пропущенных значений по признакам.

In [ ]:

```
1 np.sum(adult_df.isnull(), axis=0)
```

Разобьём данные на обучающую и тестовую выборки в отношении 3:1.

In [ ]:

```
1 X_adult_train, X_adult_test, y_adult_train, y_adult_test = train_test_split(
2     adult_df, target, random_state=777
3 )
```

При помощи кросс-валидации найдём оптимальные гиперпараметры для каждого из критериев разбиения деревьев для классификации.

In [ ]:

```
1 <...>
```

Проведите эксперименты с построением графиков, аналогичные тем, что были сделаны в предыдущем пункте для задач классификации и регрессии.

In [ ]:

```
1 <...>
```

**Вывод.**

<...>