

Решающие деревья

In [1]:

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

started 17:38:44 2020-03-14, finished in 808ms

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:

$$\frac{|X_l|}{|X|} H(X_l) + \frac{|X_r|}{|X|} H(X_r).$$

Реализация критериев информативности.

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

Ответ. Чем меньше разнообразие целевой переменной, тем меньше должно быть значение критерия информативности - и, соответственно, мы будем пытаться минимизировать его значение. Тогда функционал качества $Q(R_m, j, t)$ будет показывать выигрыш при выборе разбиения по предикату $\{x^j < t\}$. Его мы пытаемся максимизировать по всем таким предикатам.

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

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

In [2]:

```
1  # Код функций, реализующих критерии разбиения.
2
3  def mean_square_criterion(y):
4      ''' Критерий для квадратичной функции потерь. '''
5
6      return np.mean((y - np.mean(y))**2)
7
8
9  def mean_abs_criterion(y):
10     ''' Критерий для абсолютной функции потерь. '''
11
12     return np.mean(np.absolute(y - np.median(y)))
13
14
15  def get_probs_by_y(y):
16     ''' Возвращает вектор частот для каждого класса выборки. '''
17
18     _, counts = np.unique(y, return_counts=True)
19     return counts / np.sum(counts)
20
21
22  def gini_criterion(y):
23     ''' Критерий Джини. '''
24
25     probs = get_probs_by_y(y)
26     return np.sum(probs * (1 - probs))
27
28
29  def entropy_criterion(y):
30     ''' Энтропийный критерий. '''
31
32     probs = get_probs_by_y(y)
33     return -np.sum(probs * np.log(probs))
```

started 17:38:45 2020-03-14, finished in 9ms

Протестируйте реализованные функции.

Тесты для распределения вероятностей на классах.

In [3]:

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

started 17:38:45 2020-03-14, finished in 12ms

Тесты для критериев разбиения.

In [4]:

```
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)
```

started 17:38:45 2020-03-14, finished in 6ms

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 1. `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 2. `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 3. `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 4. `DecisionTreeClassifier` - наследник класса `BaseDecisionTree`, в котором реализованы функции для предсказаний при решении задачи классификации.

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

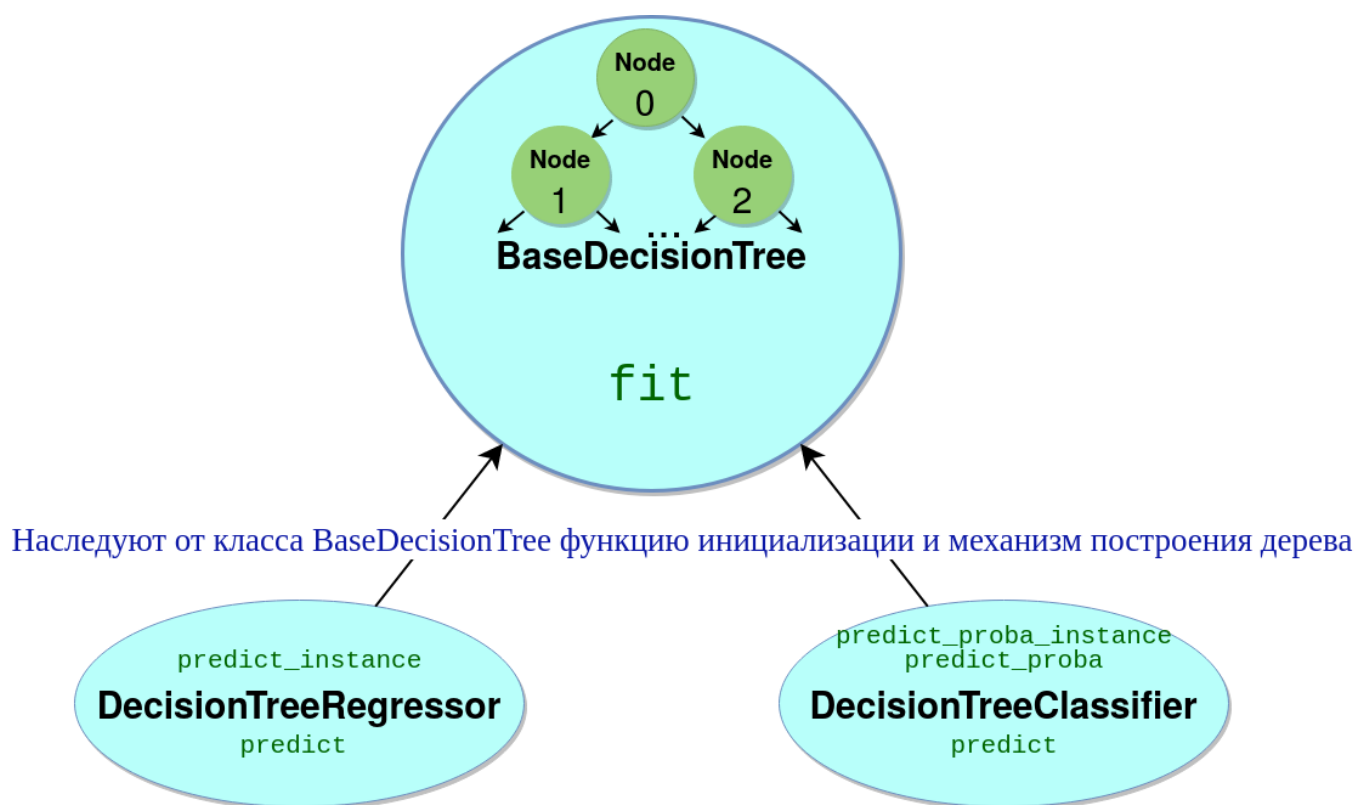
2) `predict_proba` - предсказание распределение вероятностей по классам для всех элементов выборки,

3) `predict` - 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: в качестве пороговых значений каждого из признаков нужно перебирать не все возможные действительные числа (тем более, это было бы невозможно, так как тогда пришлось бы перебирать бесконечное число вариантов), а только те значения признаков, которые содержатся в обучающей выборке. Другие значения не имеет смысла рассматривать в качестве пороговых, так как каждое из них будет находится между двумя соседними (в порядке возрастания) значениями признака и это никак не поменяет разбиение выборки на 2 части.

Структура решающего дерева



In [5]:

```
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 != np.nan]
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,
43                 criterion,
44                 max_depth=np.inf,
45                 min_samples_split=2):
46         '''
47         Функция инициализации решающего дерева.
48
49         Параметры.
50         1) criterion - критерий информативности,
51         2) max_depth - максимальная глубина дерева,
52         3) min_samples_split - минимальное количество элементов
53         обучающей выборки, которое должно попасть в вершину,
54         чтобы потом происходило разбиение этой вершины.
55         '''
56
57         self.criterion = criterion
58         self.max_depth = max_depth
59         self.min_samples_split = min_samples_split
```

```

60     # Список всех вершин дерева. В самом начале
61     # работы алгоритма есть только одна
62     # вершина - корень.
63     self.nodes = [Node()]
64     # Количество классов. Актуально только
65     # при решении задачи классификации.
66     self.class_count = 1
67     # Сюда нужно будет записать все значения
68     # для каждого из признаков датасета.
69     self.feature_values = None
70
71     def build_(self, v, X, y, depth=0):
72         """
73         Рекурсивная функция построения решающего дерева.
74
75         Параметры.
76         1) v - номер рассматриваемой вершины
77         2) X, y - обучающая выборка, попавшая в текущую вершину
78         3) depth - глубина вершины с номером v
79         """
80
81         if depth == self.max_depth or len(y) < self.min_samples_split:
82             # Если строим дерево для классификации, то
83             # сохраняем метки классов всех элементов выборки,
84             # попавших в вершину.
85             if callable(getattr(self, "set_class_count", None)):
86                 self.nodes[v].y_values = y.copy()
87                 # Для регрессии сразу вычислим среднее всех
88                 # элементов вершины.
89             else:
90                 self.nodes[v].y_values = np.mean(y)
91             return
92
93         best_criterion_value = np.inf
94         best_feature, best_threshold = 0, 0
95         sample_size, feature_count = X.shape
96
97         # переберём все возможные признаки и значения порогов,
98         # найдём оптимальный признак и значение порога
99         # и запишем их в best_feature, best_threshold
100        for feature_id in range(feature_count):
101            for threshold in self.feature_values[feature_id]:
102                # делим вершину по рассматриваемому признаку
103                # и пороговому значению
104                y_l = y[np.nan_to_num(X[:, feature_id], threshold) < threshold]
105                y_r = y[np.nan_to_num(X[:, feature_id], threshold - 1) >= threshold]
106                if len(y_l) == 0 or len(y_r) == 0:
107                    continue
108
109                left_fraction = len(y_l) / float(sample_size)
110                criterion_left = left_fraction * self.criterion(y_l)
111                criterion_right = (1 - left_fraction) * self.criterion(y_r)
112                # если для рассматриваемого признака и
113                # порога значение критерия лучше, чем для
114                # всех предыдущих пар (признак, порог),
115                # то обновляем оптимальный признак и
116                # порог разбиения
117                if criterion_left + criterion_right < best_criterion_value:
118                    best_criterion_value = criterion_left + criterion_right
119                    best_feature = feature_id
120                    best_threshold = threshold

```

```

121
122     # сохраним найденные параметры в класс текущей вершины
123     self.nodes[v].split_feature = best_feature
124     self.nodes[v].split_threshold = best_threshold
125     # разделим выборку на 2 части по порогу
126     left_indices = np.nan_to_num(X[:, best_feature],
127                                   best_threshold) < best_threshold
128     right_indices = np.nan_to_num(X[:, best_feature],
129                                   best_threshold + 1) >= best_threshold
130     left_indices_or_nans = np.logical_or(left_indices,
131                                           X[:, best_feature] == None)
132     right_indices_or_nans = np.logical_or(right_indices,
133                                           X[:, best_feature] == None)
134
135     X_l, y_l = X[left_indices_or_nans, :], y[left_indices_or_nans]
136     X_r, y_r = X[right_indices_or_nans, :], y[right_indices_or_nans]
137     self.nodes[v].left_prob = np.sum(left_indices) / len(y)
138     self.nodes[v].right_prob = np.sum(right_indices) / len(y)
139
140     # создаём левую и правую дочерние вершины,
141     # и кладём их в массив self.nodes
142     self.nodes.append(Node())
143     self.nodes.append(Node())
144     # сохраняем индексы созданных вершин в
145     # качестве левого и правого сына вершины v
146     self.nodes[v].left_son, self.nodes[v].right_son \
147         = len(self.nodes)-2, len(self.nodes)-1
148     # рекурсивно вызываем алгоритм построения
149     # дерева для дочерних вершин
150     self.build_(self.nodes[v].left_son, X_l, y_l, depth+1)
151     self.build_(self.nodes[v].right_son, X_r, y_r, depth+1)
152
153     def fit(self, X, y):
154         ...
155         Функция, из которой запускается построение
156         решающего дерева по обучающей выборке.
157
158         Параметры.
159         X, y - обучающая выборка
160         ...
161
162         # сохраним заранее все пороги для каждого
163         # из признаков обучающей выборки
164         X, y = np.array(X), np.array(y)
165         self.feature_values = []
166         for feature_id in range(X.shape[1]):
167             self.feature_values.append(
168                 np.unique(get_not_nans(X[:, feature_id]))
169             )
170
171         set_class_count = getattr(self, "set_class_count", None)
172         # если строится дерево для классификации,
173         # то нужно посчитать количество классов
174         if callable(set_class_count):
175             set_class_count(y)
176         self.build_(0, X, y)

```

started 17:38:45 2020-03-14, finished in 26ms

Теперь, когда общий код решающего дерева написан, нужно сделать обёртки над `BaseDecisionTree` - `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 [6]:

```
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         if self.nodes[v].left_son == None:
14             return np.mean(self.nodes[v].y_values)
15
16         # если у объекта x значение признака по
17         # которому происходит разделение, меньше
18         # порогового, то спускаемся в левое поддерево,
19         # иначе - в правое
20         if x[self.nodes[v].split_feature] < self.nodes[v].split_threshold:
21             return self.predict_instance(x, self.nodes[v].left_son)
22         elif x[self.nodes[v].split_feature] >= self.nodes[v].split_threshold:
23             return self.predict_instance(x, self.nodes[v].right_son)
24         # а если у элемента отсутствует значение
25         # разделяющего признака, то будем спускаться
26         # в оба поддерева
27         else:
28             left_predict = self.predict_instance(x, self.nodes[v].left_son)
29             right_predict = self.predict_instance(x, self.nodes[v].right_son)
30
31             return self.nodes[v].left_prob * left_predict + \
32                    self.nodes[v].right_prob * right_predict
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]
```

started 17:38:45 2020-03-14, finished in 9ms

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 [7]:

```
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(self.nodes[v].y_values,
29                                         return_counts=True)
30             result[classes.astype(int)] = counts
31             return result / np.sum(result)
32
33         # если у x значение признака по которому происходит разделение,
34         # меньше порогового, то спускаемся в левое поддерево, иначе - в правое
35         if x[self.nodes[v].split_feature] < self.nodes[v].split_threshold:
36             return self.predict_proba_instance(x, self.nodes[v].left_son)
37         elif x[self.nodes[v].split_feature] >= self.nodes[v].split_threshold:
38             return self.predict_proba_instance(x, self.nodes[v].right_son)
39         # а если у элемента отсутствует значение разделяющего признака,
40         # то будем спускаться в оба дерева
41         else:
42             left_predict \
43                 = self.predict_proba_instance(x, self.nodes[v].left_son)
44             right_predict = \
45                 self.predict_proba_instance(x, self.nodes[v].right_son)
46
47             return self.nodes[v].left_prob * left_predict + \
48                 self.nodes[v].right_prob * right_predict
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)

```

started 17:38:45 2020-03-14, finished in 14ms

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 делать 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 [8]:

```

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

```

started 17:38:45 2020-03-14, finished in 25ms

Для критерия Джини и энтропийного критерия найдите оптимальные параметры обучения дерева - `max_depth` и `min_samples_split`. Используйте для этого кросс-валидацию.

In [9]:

```

1  classification_criteria = [gini_criterion, entropy_criterion]
2  criterion_names = ['gini', 'entropy']

```

started 17:38:45 2020-03-14, finished in 3ms

Используя кросс-валидацию, найдите оптимальные параметры на обучающей выборке и выведите их вместе с метрикой качества.

С начала надо разбить выборку на `train` и `test`.

In [10]:

```

1  X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=777)

```

started 17:38:45 2020-03-14, finished in 17ms

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

In [11]:

```
1 ▼ for criterion, criterion_name in zip(classification_criteria,
2                                       criterion_names):
3 ▼     dt = GridSearchCV(
4         estimator=DecisionTreeClassifier(criterion),
5         param_grid={
6             'max_depth': np.arange(2, 11),
7             'min_samples_split': np.arange(3, 7)
8         },
9         scoring='accuracy',
10        cv=5, # разбиение выборки на 5 фолдов
11        verbose=3, # насколько часто печатать сообщения
12        n_jobs=2, # кол-во параллельных процессов
13    )
14
15    dt.fit(X_train, y_train)
16    predictions = dt.predict(X_test)
17
18 ▼    print('Оптимальные параметры для {}'.format(criterion_name),
19          dt.best_params_)
20    accuracy = accuracy_score(y_test, predictions)
21    print('Точность на тесте:', accuracy)
22    assert accuracy >= 0.85, "Something is wrong with your classifier"
```

started 17:38:45 2020-03-14, finished in 7m 5s

Fitting 5 folds for each of 36 candidates, totalling 180 fits

[Parallel(n_jobs=2)]: Using backend LokyBackend with 2 concurrent workers.

[Parallel(n_jobs=2)]: Done 28 tasks | elapsed: 7.6s

[Parallel(n_jobs=2)]: Done 124 tasks | elapsed: 1.4min

[Parallel(n_jobs=2)]: Done 180 out of 180 | elapsed: 3.0min finished

Оптимальные параметры для gini: {'max_depth': 3, 'min_samples_split': 3}

Точность на тесте: 0.9333333333333333

Fitting 5 folds for each of 36 candidates, totalling 180 fits

[Parallel(n_jobs=2)]: Using backend LokyBackend with 2 concurrent workers.

[Parallel(n_jobs=2)]: Done 28 tasks | elapsed: 12.4s

[Parallel(n_jobs=2)]: Done 124 tasks | elapsed: 2.3min

[Parallel(n_jobs=2)]: Done 180 out of 180 | elapsed: 4.0min finished

Оптимальные параметры для entropy: {'max_depth': 4, 'min_samples_split': 3}

Точность на тесте: 0.9333333333333333

Построение графиков.

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

Поскольку при кросс-валидации наилучшее значение `min_samples_split` для обоих критериев равно 3, будем строить графики зависимости accuracy от `max_depth` при фиксации `min_samples_split = 3`.

In [12]:

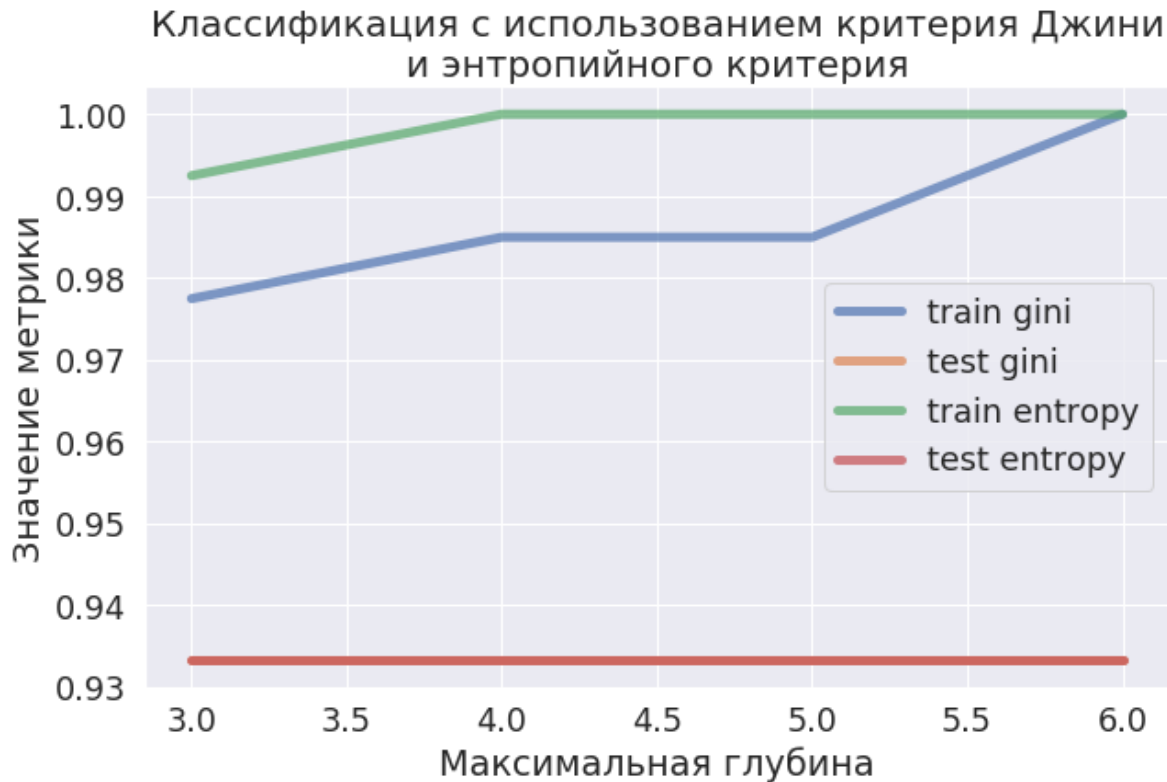
```
1  def plot_graphs(tree_estimator, criterions, scorer, title,
2                  max_depth_list=np.arange(3, 7), min_samples_split=3):
3      '''
4      Функция для рисования графиков.
5
6      Параметры.
7      1) tree_estimator - экземпляр класса решающего дерева,
8      2) criterions - критерии разбиения дерева,
9      3) scorer - метрика качества,
10     4) title - заголовок графика,
11     5) max_depth_list - список рассматриваемых значений макс.
12     глубины дерева,
13     6) min_samples_split - значение параметра min_samples_split.
14     '''
15
16     plt.figure(figsize=(10, 6))
17
18     for criterion in criterions:
19         score_train = []
20         score_test = []
21
22         for max_depth in max_depth_list:
23             estimator = tree_estimator(criterion[0], max_depth=max_depth,
24                                       min_samples_split=min_samples_split)
25             estimator.fit(X_train, y_train)
26             score_train.append(scorer(estimator.predict(X_train), y_train))
27             score_test.append(scorer(estimator.predict(X_test), y_test))
28
29         plt.plot(max_depth_list, score_train,
30                 lw=5, alpha=0.7, label='train {}'.format(criterion[1]))
31         plt.plot(max_depth_list, score_test,
32                 lw=5, alpha=0.7, label='test {}'.format(criterion[1]))
33
34     plt.xlabel('Максимальная глубина')
35     plt.ylabel('Значение метрики')
36     plt.legend()
37     plt.title(title, fontsize=20)
38     plt.show()
```

started 17:45:50 2020-03-14, finished in 7ms

In [13]:

```
1 plot_graphs(DecisionTreeClassifier,  
2             [(gini_criterion, 'gini'), (entropy_criterion, 'entropy')],  
3             accuracy_score,  
4             'Классификация с использованием критерия Джини\n' \  
5             + 'и энтропийного критерия')
```

started 17:45:50 2020-03-14, finished in 16.6s



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

Вывод.

На обоих графиках видно, что с некоторого значения максимальной глубины точность классификации на тестовой выборке перестаёт расти. Это указывает на то, что, возможно, при больших значениях макс. глубины, чем это, происходит переобучение дерева.

Решающее дерево, использующее энтропийный критерий при разбиении, дало более высокий результат как на тестовой, так и на обучающей выборке.

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 [14]:

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

started 17:46:07 2020-03-14, finished in 14ms

In [15]:

```
1 regression_criteria = [mean_square_criterion, mean_abs_criterion]
2 criterion_names = ['mean_square', 'mean_abs']
```

started 17:46:07 2020-03-14, finished in 7ms

Разобьём выборку на обучение и тест.

In [16]:

```
1 X_train, X_test, y_train, y_test = train_test_split(boston_X, boston_y,
2                                                    random_state=777)
```

started 17:46:07 2020-03-14, finished in 7ms

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

In [17]:

```
1 ▼ for criterion, criterion_name in zip(regression_criteria, criterion_names):
2 ▼     dt_regressor = GridSearchCV(
3         estimator=DecisionTreeRegressor(criterion),
4 ▼         param_grid={
5             'max_depth': np.arange(2, 11)
6         },
7         scoring='r2',
8         cv=5, # разбиение выборки на 5 фолдов
9         verbose=3, # насколько часто печатать сообщения
10        n_jobs=2, # кол-во параллельных процессов
11    )
12
13    dt_regressor.fit(X_train, y_train)
14    predictions = dt_regressor.predict(X_test)
15
16 ▼    print('Оптимальные параметры для {}'.format(criterion_name),
17          dt.best_params_)
18    print('r2_score на тесте:', r2_score(y_test, predictions))
```

started 17:46:07 2020-03-14, finished in 9m 14s

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=2)]: Using backend LokyBackend with 2 concurrent workers.

[Parallel(n_jobs=2)]: Done 28 tasks | elapsed: 1.2min

[Parallel(n_jobs=2)]: Done 45 out of 45 | elapsed: 4.0min finished

Оптимальные параметры для mean_square: {'max_depth': 4, 'min_samples_split': 3}

r2_score на тесте: 0.6752210376459766

Fitting 5 folds for each of 9 candidates, totalling 45 fits

[Parallel(n_jobs=2)]: Using backend LokyBackend with 2 concurrent workers.

[Parallel(n_jobs=2)]: Done 28 tasks | elapsed: 1.6min

[Parallel(n_jobs=2)]: Done 45 out of 45 | elapsed: 5.1min finished

Оптимальные параметры для mean_abs: {'max_depth': 4, 'min_samples_split': 3}

r2_score на тесте: 0.6397486631555982

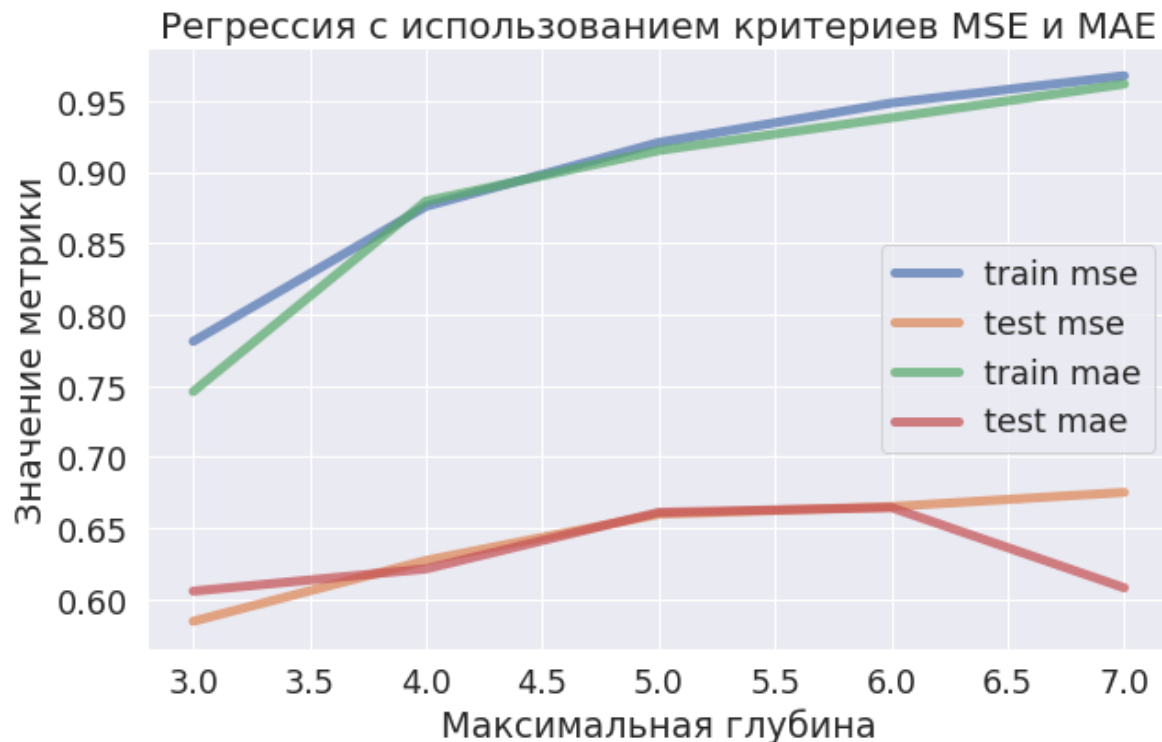
Вывод.

В данном случае оптимальные параметры для обоих критериев разбиения совпадают. Однако, видно, что критерий, основанный на mse, даёт гораздо более высокий результат.

In [18]:

```
1 plot_graphs(DecisionTreeRegressor,  
2             [(mean_square_criterion, 'mse'), (mean_abs_criterion, 'mae')],  
3             r2_score,  
4             'Регрессия с использованием критериев MSE и MAE',  
5             max_depth_list=range(3, 8))
```

started 17:55:21 2020-03-14, finished in 1m 4.18s



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.

Качество результатов, выдаваемых решающим деревом, построенным с использованием mae критерия, оказалось значительно выше, чем даёт дерево, использующее mse критерий. Но, тем не менее, на втором графике видно, что при $\text{max_depth} > 5$ решающее дерево начинает переобучаться, так как score на обучающей выборке растёт, а на тестовой - падает.

Можно ускорить построение решающего дерева, снизив асимптотику разбиения по признаку с $O(u^2)$ до $O(u)$, где u - количество различных значений признака.

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 [19]:

```
1 column_names = [  
2     'age', 'workclass', 'fnlwgt', 'education1',  
3     'education2', 'marital-status', 'occupation',  
4     'relationship', 'race', 'sex', 'capital-gain',  
5     'capital-loss', 'hours-per-week', 'native-country', 'target'  
6 ]
```

started 17:56:25 2020-03-14, finished in 3ms

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 2000 строк данных.

In [20]:

```
1 adult_df = pd.read_csv('adult.data', header=None)[:2000]  
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()
```

started 17:56:25 2020-03-14, finished in 175ms

Out[20]:

	age	workclass	fnlwgt	education1	education2	marital-status	occupation	relationship	race
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black

In [21]:

```
1 adult_df.shape
```

started 17:56:25 2020-03-14, finished in 3ms

Out[21]:

(2000, 14)

Предобработаем датасет, заменив категориальные признаки one-hot векторами.

In [22]:

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

started 17:56:25 2020-03-14, finished in 38ms

Out[22]:

	age	fnlwgt	education2	capital-gain	capital-loss	hours-per-week	workclass_?	workclass_Federal-gov	workclass_Local-gov	wo
0	39	77516	13	2174	0	40	0	0	0	
1	50	83311	13	0	0	13	0	0	0	
2	38	215646	9	0	0	40	0	0	0	
3	53	234721	7	0	0	40	0	0	0	
4	28	338409	13	0	0	40	0	0	0	

5 rows × 102 columns

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 [23]:

```
1 for feature in column_names:
2     if f'{feature}_?' in adult_df.columns:
3         none_indices = np.arange(adult_df.shape[0])[
4             adult_df[f'{feature}_?'] == 1
5         ]
6
7     for dummy_feature in adult_df.columns:
8         if dummy_feature.startswith(f'{feature}_') \
9             and dummy_feature != f'{feature}_?':
10             adult_df[dummy_feature][none_indices] = np.nan
11
12     adult_df = adult_df.drop(f'{feature}_?', axis=1)
```

started 17:56:25 2020-03-14, finished in 109ms

Посмотрим на распределение пропущенных значений по признакам.

In [24]:

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

started 17:56:26 2020-03-14, finished in 8ms

Out[24]:

```
age                                0
fnlwgt                             0
education2                         0
capital-gain                       0
capital-loss                       0
..
native-country_ Taiwan             39
native-country_ Thailand           39
native-country_ Trinidad&Tobago    39
native-country_ United-States      39
native-country_ Yugoslavia         39
Length: 99, dtype: int64
```

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

In [25]:

```
1 X_train, X_test, y_train, y_test = train_test_split(
2     adult_df.values, target.values, random_state=777
3 )
```

started 17:56:26 2020-03-14, finished in 10ms

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

In [26]:

```
1  ▼ for criterion, criterion_name in zip(classification_criteria,
2                                         criterion_names):
3  ▼     dt_classifier = GridSearchCV(
4         estimator=DecisionTreeClassifier(criterion),
5         param_grid={
6             'max_depth': np.arange(2, 10),
7             'min_samples_split': np.arange(3, 7)
8         },
9         scoring='accuracy',
10        cv=5, # разбиение выборки на 5 фолдов
11        verbose=3, # насколько часто печатать сообщения
12        n_jobs=2, # кол-во параллельных процессов
13    )
14
15    dt_classifier.fit(X_train, y_train)
16    predictions = dt_classifier.predict(X_test)
17
18  ▼    print('Оптимальные параметры для {}'.format(criterion_name),
19          dt_classifier.best_params_)
20    accuracy = accuracy_score(y_test, predictions)
21    print('Точность на тесте:', accuracy)
```

started 17:56:26 2020-03-14, finished in 41m 58s

Fitting 5 folds for each of 32 candidates, totalling 160 fits

[Parallel(n_jobs=2)]: Using backend LokyBackend with 2 concurrent workers.

[Parallel(n_jobs=2)]: Done 28 tasks | elapsed: 28.1s

[Parallel(n_jobs=2)]: Done 124 tasks | elapsed: 10.3min

[Parallel(n_jobs=2)]: Done 160 out of 160 | elapsed: 19.1min finished

Оптимальные параметры для mean_square: {'max_depth': 4, 'min_samples_split': 3}

Точность на тесте: 0.832

Fitting 5 folds for each of 32 candidates, totalling 160 fits

[Parallel(n_jobs=2)]: Using backend LokyBackend with 2 concurrent workers.

[Parallel(n_jobs=2)]: Done 28 tasks | elapsed: 25.8s

[Parallel(n_jobs=2)]: Done 124 tasks | elapsed: 9.9min

[Parallel(n_jobs=2)]: Done 160 out of 160 | elapsed: 22.6min finished

Оптимальные параметры для mean_abs: {'max_depth': 4, 'min_samples_split': 3}

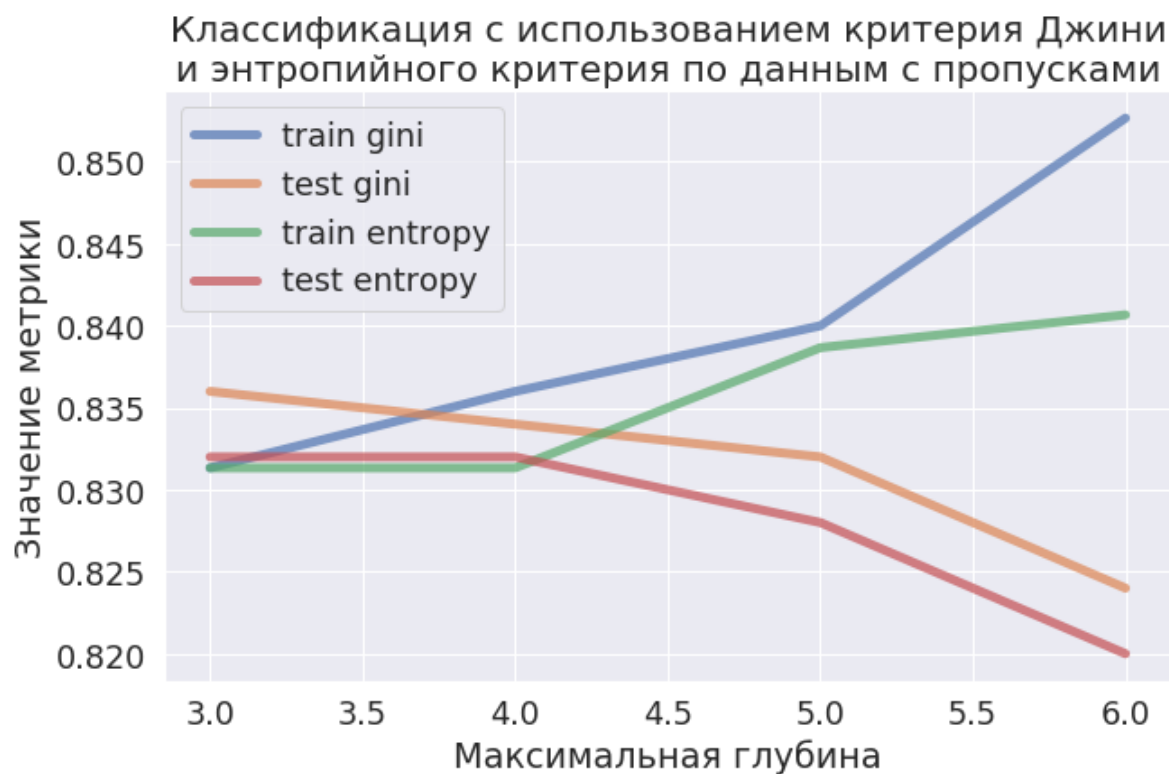
Точность на тесте: 0.832

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

In [27]:

```
1 plot_graphs(  
2     DecisionTreeClassifier,  
3     [(gini_criterion, 'gini'), (entropy_criterion, 'entropy')],  
4     accuracy_score,  
5     'Классификация с использованием критерия Джини\n'  
6     'и энтропийного критерия по данным с пропусками'  
7 )
```

started 18:38:23 2020-03-14, finished in 1m 17.5s



Вывод.

По графикам видно, что дерево начинается переобучаться для макс. глубины ≥ 5 . Также можно заметить, что использование критерия Джини даёт лучшие результаты как на обучающей, так и на тестовой выборке.