

# Рубежный контроль №2

Мирсонов Вячеслав Александрович

Группа РТ5-61Б

Вариант 9

## Постановка задачи

### Условие:

Для заданного набора данных (по Вашему варианту) постройте модели классификации или регрессии (в зависимости от конкретной задачи, рассматриваемой в наборе данных). Для построения моделей используйте методы 1 и 2 (по варианту для Вашей группы). Оцените качество моделей на основе подходящих метрик качества (не менее двух метрик). Какие метрики качества Вы использовали и почему? Какие выводы Вы можете сделать о качестве построенных моделей? Для построения моделей необходимо выполнить требуемую предобработку данных: заполнение пропусков, кодирование категориальных признаков, и т.д.

### Методы:

Метод №1: "Дерево решений". Метод №2: "Градиентный бустинг".

Набор данных: Houses to rent in Brazil.

<https://www.kaggle.com/rubenssjr/brasilian-houses-to-rent>

## Импорт библиотек:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

In [36]:

## Загрузка и первичная подготовка данных:

```
# Загрузка данных
data = pd.read_csv("houses_to_rent_v2.csv")
data.shape
```

In [38]:

```
(499, 13)
```

Out[38]:

```
data
```

In [39]:

Out[39]:

	city	area	rooms	bathroom	parking spaces	floor	animal	furniture	hoa (R\$)	rent amount (R\$)	property tax (R\$)	fire insurance (R\$)	total (R\$)
0	São Paulo	70	2	1	1	7	accept	furnished	2065	3300	211	42	5618
1	São Paulo	320	4	4	0	20	accept	not furnished	1200	4960	1750	63	7973
2	Porto Alegre	80	1	1	1	6	accept	not furnished	1000	2800	0	41	3841
3	Porto Alegre	51	2	1	0	2	accept	not furnished	270	1112	22	17	1421
4	São Paulo	25	1	1	0	1	not accept	not furnished	0	800	25	11	836
...	...	...	...	...	...	...	...	...	...	...	...	...	...
494	São Paulo	120	3	3	2	9	accept	furnished	1200	3500	109	45	4854
495	Rio de Janeiro	90	2	2	0	5	accept	not furnished	830	2500	84	33	3447
496	Porto Alegre	125	4	2	0	7	accept	not furnished	1006	1632	150	14	2802
497	São Paulo	200	3	4	3	6	accept	furnished	1800	5058	1200	65	8123
498	São Paulo	190	3	3	2	5	accept	not furnished	1800	3260	292	42	5394

499 rows × 13 columns

```
# уникальные значения столбца 'city'  
data['city'].unique()
```

In [41]:

```
array(['São Paulo', 'Porto Alegre', 'Rio de Janeiro', 'Campinas',  
      'Belo Horizonte'], dtype=object)
```

Out[41]:

```
# список колонок с типами данных  
data.dtypes
```

In [42]:

```
city                object  
area                int64  
rooms              int64  
bathroom           int64  
parking spaces     int64  
floor              int64  
animal             object  
furniture          object  
hoa (R$)           int64  
rent amount (R$)   int64  
property tax (R$)  int64  
fire insurance (R$) int64  
total (R$)         int64  
dtype: object
```

Out[42]:

Обработка пропусков в данных:

```
# проверим, есть ли пропущенные значения  
data.isnull().sum()
```

In [43]:

```
city          0
area          0
rooms         0
bathroom      0
parking spaces 0
floor         0
animal        0
furniture     0
hoa (R$)      0
rent amount (R$) 0
property tax (R$) 0
fire insurance (R$) 0
total (R$)    0
dtype: int64
```

Пропущенных данных нет, можем приступать к масштабированию данных

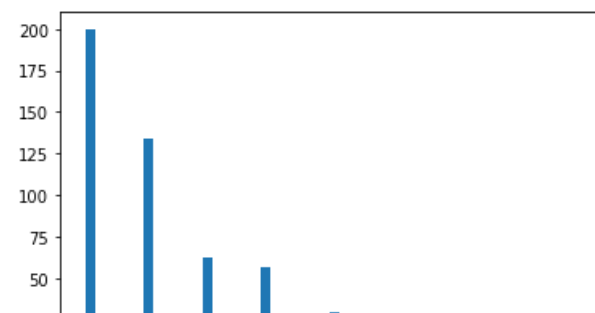
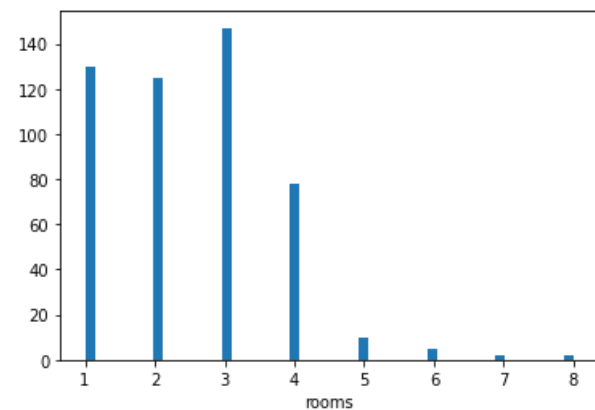
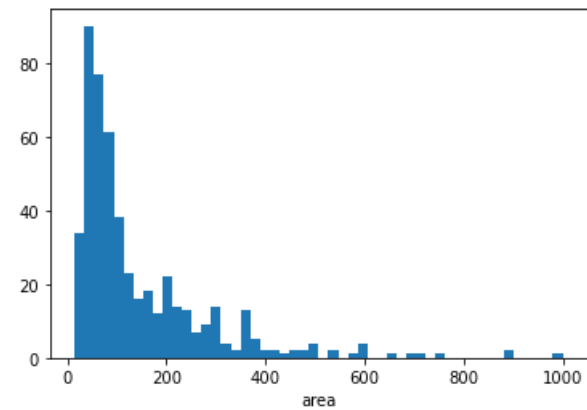
## Масштабирование данных:

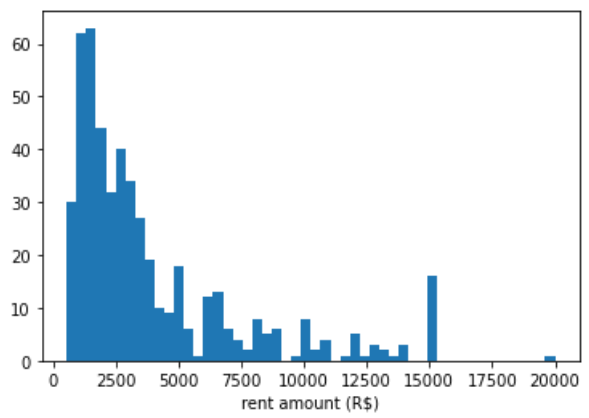
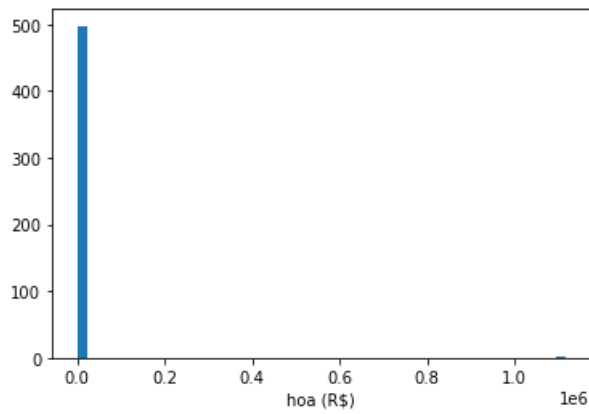
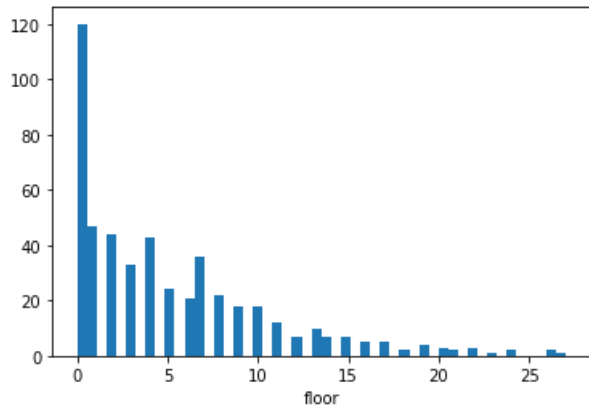
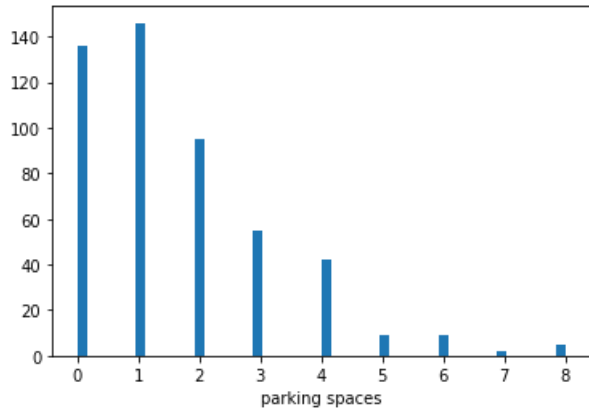
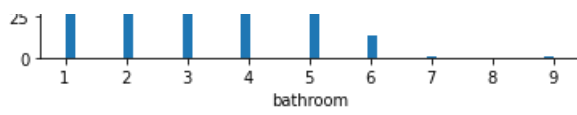
### MinMax масштабирование:

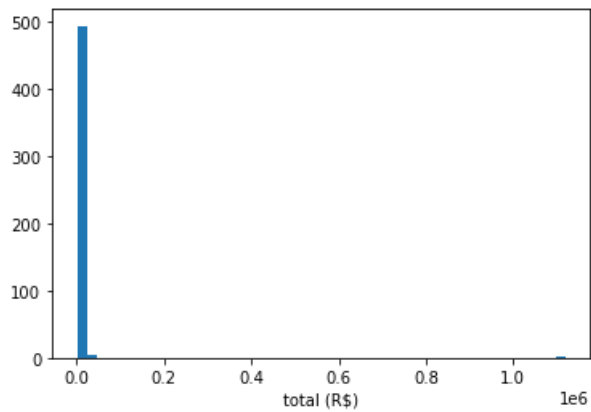
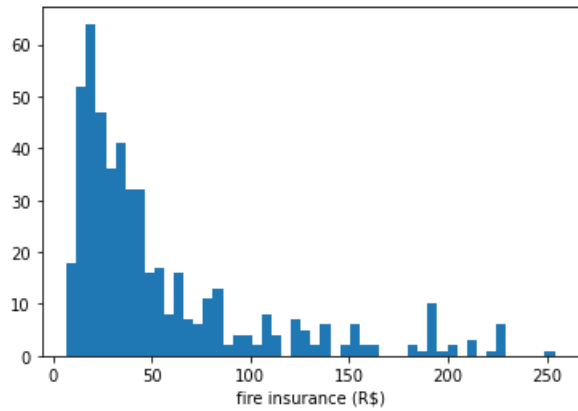
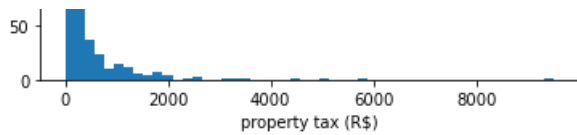
$$x_{\text{новый}} = \frac{x_{\text{старый}} - \min(X)}{\max(X) - \min(X)}$$

```
from sklearn.preprocessing import MinMaxScaler
```

```
# СПИСОК КОЛОНОК С ЧИСЛОВЫМИ ДАННЫМИ
num_cols = ['area', 'rooms', 'bathroom', 'parking spaces', 'floor', 'hoa (R$)', 'rent amount (R$)', 'property tax (R$)', 'fire insurance (R$)', 'total (R$)']
# Гистограмма по признакам
for col in data[num_cols]:
    plt.hist(data[col], 50)
    plt.xlabel(col)
    plt.show()
```

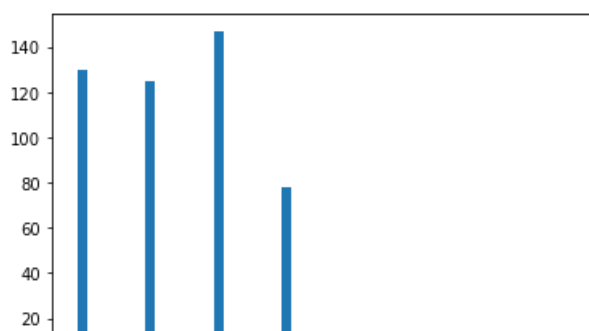
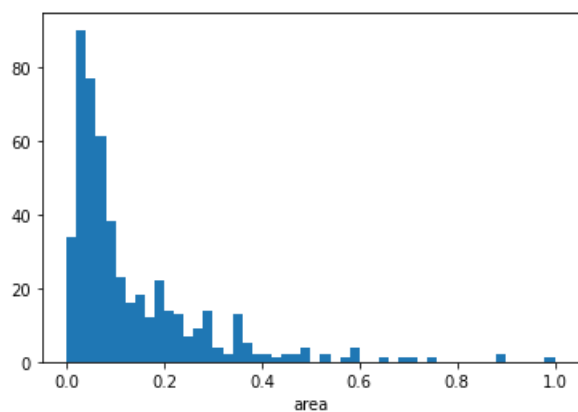






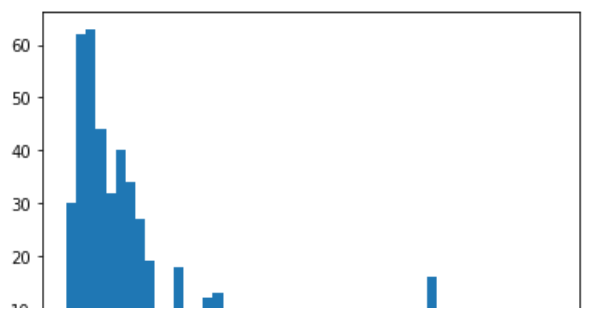
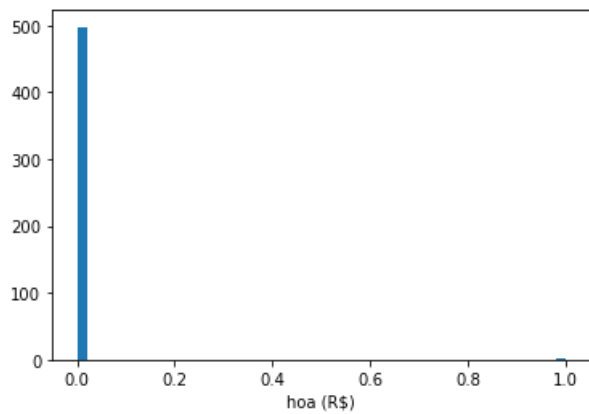
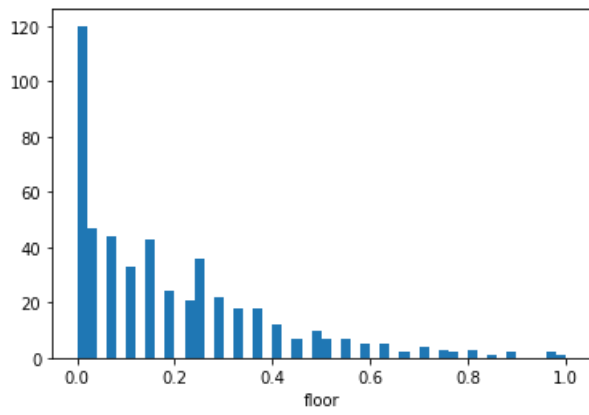
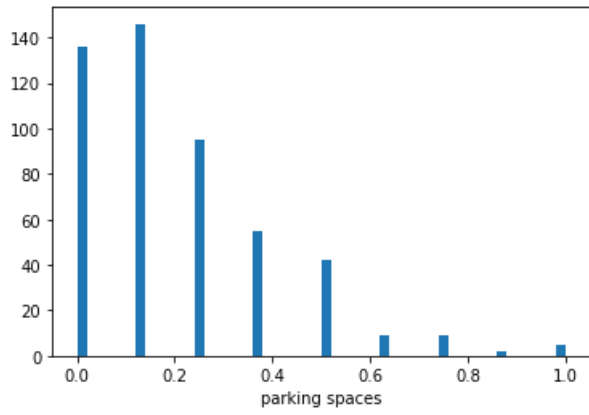
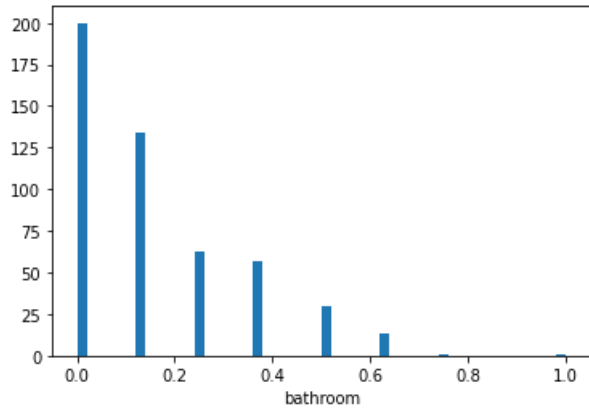
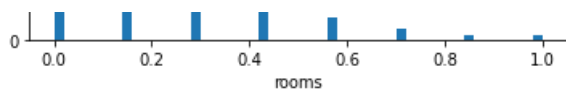
```
# MinMax масштабирование
sc1 = MinMaxScaler()
for item in num_cols:
    data.loc[:, item] = sc1.fit_transform(data[[item]])
```

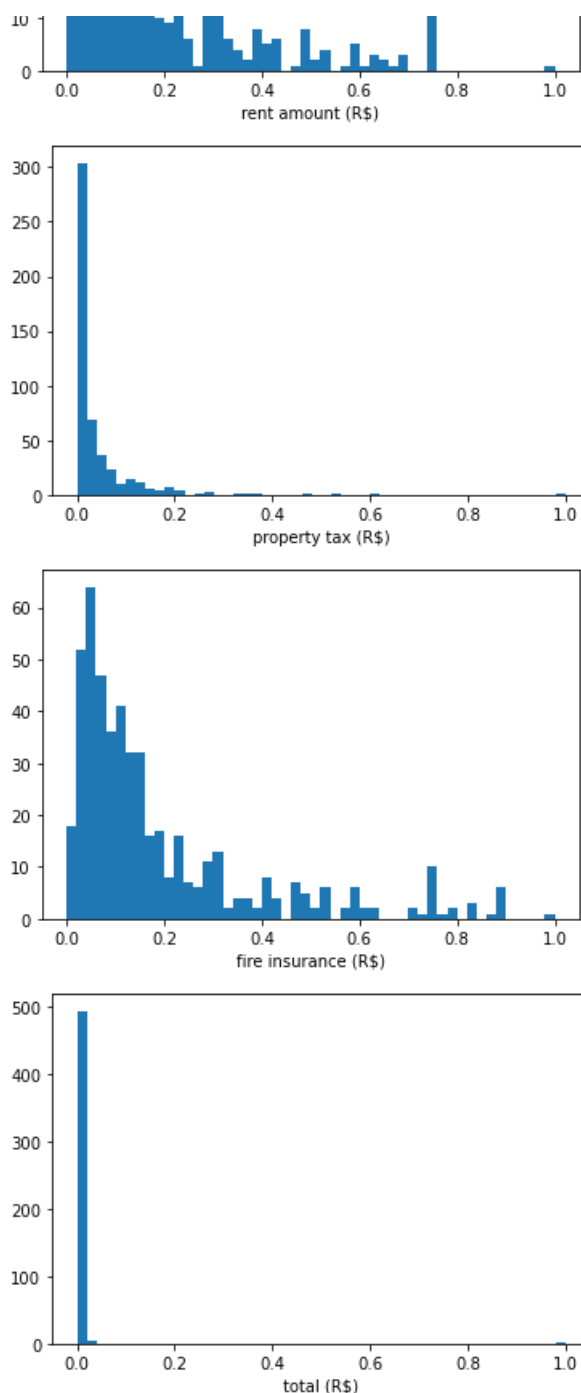
```
# Гистограмма по отмасштабированным признакам
for col in data[num_cols]:
    plt.hist(data[col], 50)
    plt.xlabel(col)
    plt.show()
```



In [49]:

In [50]:





### Кодирование категориальных признаков:

In [51]:

```
# Выберем категориальные колонки
# Цикл по колонкам датасета
for col in data.columns:
    dt = str(data[col].dtype)
    if dt=='object':
        temp_un = data[col].nunique()
        print('Колонка {}. Тип данных {}. Количество уникальных значений {}'.format(col, dt, temp_un))
```

Колонка city. Тип данных object. Количество уникальных значений 5.  
 Колонка animal. Тип данных object. Количество уникальных значений 2.  
 Колонка furniture. Тип данных object. Количество уникальных значений 2.

### Кодирование категорий наборами бинарных значений - one-hot encoding

In [52]:

```
cat_cols = ['city', 'animal', 'furniture']
one_hot = pd.get_dummies(data[cat_cols].astype(str))
one_hot.head()
```

Out[52]:

	city_Belo Horizonte	city_Campinas	city_Porto Alegre	city_Rio de Janeiro	city_São Paulo	animal_acept	animal_not_acept	furniture_furnished	furniture_not_furnished
0	0	0	0	0	1	1	0	1	0
1	0	0	0	0	1	1	0	0	1
2	0	0	1	0	0	1	0	0	1
3	0	0	1	0	0	1	0	0	1
4	0	0	0	0	1	0	1	0	1

In [53]:

```
# Замена исходных категориальных колонок наборами бинарных значений
data = data.join(one_hot)
data.drop(columns=cat_cols, inplace=True)
```

In [54]:

```
# первые 5 строк получившегося набора данных
data.head()
```

Out[54]:

	area	rooms	bathroom	parking spaces	floor	hoa (R\$)	rent amount (R\$)	property tax (R\$)	fire insurance (R\$)	total (R\$)	city_Belo Horizonte	city_Campinas	city_Porto Alegre	city_Rio de Janeiro
0	0.055894	0.142857	0.000	0.125	0.259259	0.001849	0.141388	0.022211	0.141700	0.004473	0	0	0	C
1	0.309959	0.428571	0.375	0.000	0.740741	0.001074	0.226735	0.184211	0.226721	0.006577	0	0	0	C
2	0.066057	0.000000	0.000	0.125	0.222222	0.000895	0.115681	0.000000	0.137652	0.002886	0	0	1	C
3	0.036585	0.142857	0.000	0.000	0.074074	0.000242	0.028895	0.002316	0.040486	0.000724	0	0	1	C
4	0.010163	0.000000	0.000	0.000	0.037037	0.000000	0.012853	0.002632	0.016194	0.000201	0	0	0	C

## Построение моделей:

### Разделение выборки на обучающую и тестовую

In [55]:

```
from sklearn.model_selection import train_test_split
data_train, data_test, data_y_train, data_y_test = train_test_split(data[data.columns.drop('total (R$)')], data_y,
```

### Модель "Дерево решений"

In [56]:

```
from sklearn.tree import DecisionTreeRegressor
dtt = DecisionTreeRegressor(random_state=1).fit(data_train, data_y_train)
data_test_predicted_dtt = dtt.predict(data_test)
```

### Модель "Градиентный бустинг"

In [57]:

```
from sklearn.ensemble import GradientBoostingRegressor
gbdt = GradientBoostingRegressor(random_state=10).fit(data_train, data_y_train)
data_test_predicted_gbdt = gbdt.predict(data_test)
```

### Оценка качества моделей:

В качестве метрик для оценки качества моделей будем использовать Mean squared error (средняя квадратичная ошибка), как наиболее часто используемую метрику для оценки качества регрессии, и метрику  $R^2$  (коэффициент детерминации), потому что эта метрика является нормированной.

In [58]:

```
from sklearn.metrics import mean_squared_error, r2_score
# Mean squared error - средняя квадратичная ошибка
print('Метрика MSE:\nДерево решений: {} \nГрадиентный бустинг: {}'.format(mean_squared_error(data_y_test, data_train_predicted_dtt), mean_squared_error(data_y_test, data_train_predicted_gbdt)))
```

Метрика MSE:

Дерево решений: 4.009125408728235e-07

Градиентный бустинг: 6.742776208219548e-07

In [60]:

```
#Метрика R2 или коэффициент детерминации
```



```
print('Метрика R\u00B2:\nДерево решений: {}\nГрадиентный бустинг: {}'.format(r2_score(data_y_test, data_test_1,
```

Метрика R<sup>2</sup>:

Дерево решений: 0.974242110132498

Градиентный бустинг: 0.9566789089225258

## Выводы о качестве построенных моделей:

Исходя из оценки качества построенных моделей, можно увидеть, что модель "Дерево решений" лучше справляется с задачей по сравнению с моделью "Градиентный бустинг", что может свидетельствовать о переобучении модели "Градиентный бустинг".