

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import warnings
4
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 from sklearn.model_selection import train_test_split, GridSearchCV
9 from sklearn.preprocessing import MinMaxScaler, StandardScaler, LabelEncoder
10 from sklearn.metrics import mean_squared_error, make_scorer
11 from sklearn.ensemble import GradientBoostingRegressor
12
13 from xgboost import XGBRegressor
14 from lightgbm import LGBMRegressor
15 from catboost import CatBoostRegressor
16
17 sns.set(font_scale=1.5, palette='Set2')
18 warnings.filterwarnings("ignore")
```

started 01:43:57 2020-03-27, finished in 2.11s

Задача 6.2

Реализация MAPE

In [2]:

```
1 def mape(y_true, y_pred):
2     return 100 * np.abs((y_true - y_pred) / y_true).mean()
```

started 01:43:59 2020-03-27, finished in 4ms

Считываем данные

In [3]:

```
1 train_df = pd.read_csv('houses_train.csv').drop(columns=['id', 'date'])
2 train_df.head()
```

started 01:43:59 2020-03-27, finished in 105ms

Out[3]:

	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	condition	gr
0	221900.0	3	1.00	1180	5650	1.0	0	0	3	
1	538000.0	3	2.25	2570	7242	2.0	0	0	3	
2	604000.0	4	3.00	1960	5000	1.0	0	0	5	
3	510000.0	3	2.00	1680	8080	1.0	0	0	3	
4	257500.0	3	2.25	1715	6819	2.0	0	0	3	

Посмотрим на их описания

In [4]:

1	train_df.describe()
started 01:43:59 2020-03-27, finished in 137ms	

Out[4]:

	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	
count	1.562600e+04	15626.000000	15626.000000	15626.000000	1.562600e+04	15626.000000	1!
mean	5.355339e+05	3.371240	2.111737	2073.132919	1.506655e+04	1.494304	
std	3.595051e+05	0.909872	0.769037	911.406092	4.235533e+04	0.539333	
min	7.500000e+04	0.000000	0.000000	290.000000	6.000000e+02	1.000000	
25%	3.200000e+05	3.000000	1.500000	1430.000000	5.060000e+03	1.000000	
50%	4.500000e+05	3.000000	2.250000	1910.000000	7.598500e+03	1.500000	
75%	6.400000e+05	4.000000	2.500000	2540.000000	1.057975e+04	2.000000	
max	7.700000e+06	11.000000	8.000000	13540.000000	1.651359e+06	3.500000	

In [5]:

1	train_df.columns
started 01:44:00 2020-03-27, finished in 6ms	

Out[5]:

```
Index(['price', 'bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors',  
      'waterfront', 'view', 'condition', 'grade', 'sqft_above',  
      'sqft_basement', 'yr_built', 'yr_renovated', 'zipcode', 'lat',  
      'long',  
      'sqft_living15', 'sqft_lot15'],  
      dtype='object')
```

Внимательно изучив признаки понимаем, что единственным категориальным признаком можно считать почтовый индекс `zipcode`

In [6]:

1	cat_features = ['zipcode']
started 01:44:00 2020-03-27, finished in 8ms	

Уникальные значения этого признака

In [7]:

```
1 train_df['zipcode'].unique()
```

started 01:44:00 2020-03-27, finished in 8ms

Out[7]:

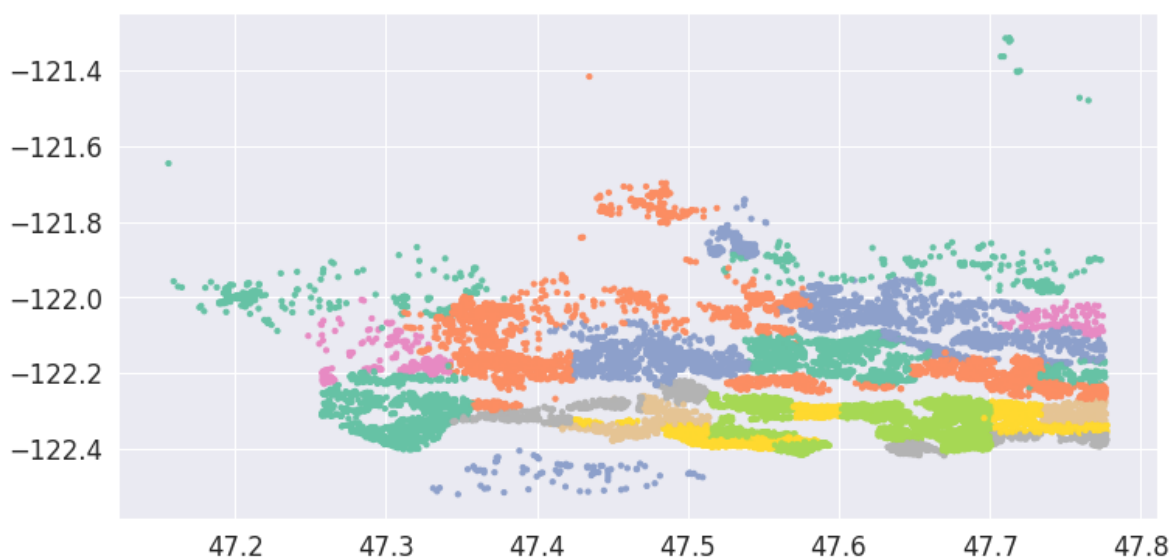
```
array([98178, 98125, 98136, 98074, 98003, 98198, 98146, 98038, 98007,
       98115, 98028, 98126, 98019, 98103, 98002, 98133, 98040, 98092,
       98030, 98119, 98112, 98052, 98027, 98117, 98058, 98001, 98056,
       98053, 98148, 98105, 98042, 98008, 98059, 98166, 98122, 98004,
       98005, 98116, 98023, 98199, 98032, 98045, 98102, 98077, 98168,
       98065, 98107, 98006, 98109, 98022, 98034, 98075, 98033, 98144,
       98177, 98155, 98108, 98118, 98011, 98024, 98010, 98106, 98070,
       98031, 98072, 98014, 98055, 98029, 98188, 98039])
```

Отобразим почтовый индекс на карте города, используя значения широты и долготы. Видим, что объекты разбиваются по районам города

In [8]:

```
1 plt.figure(figsize=(12, 6))
2 ▼ plt.scatter(train_df['lat'], train_df['long'],
3               c=train_df['zipcode'], cmap='Set2', s=10);
```

started 01:44:00 2020-03-27, finished in 1.57s



Разобьем данные на обучающую и валидационную выборки. На первой будем обучать все модели, на второй -- проверять качество моделей

In [9]:

```
1 X, labels = train_df.drop(columns=['price']), train_df['price']
2 ▼ x_train, x_valid, y_train, y_valid = train_test_split(X, labels,
3                                                         test_size=0.3)
```

started 01:44:01 2020-03-27, finished in 29ms

Для кодировки категориальных признаков будем использовать `MeanEncoder`, который заменяет значение категории на среднее значение таргета в этой категории. Если категория не встречалась в трейне, то на глобальное среднее таргета.

In [10]:

```
1 ▾ # https://github.com/AndreyKoceruba/mean-encoding/blob/master/mean\_encoder.py
2
3 from sklearn.base import BaseEstimator
4 from sklearn.base import TransformerMixin
5
6 ▾ class MeanEncoder(BaseEstimator, TransformerMixin):
7
8 ▾     def __init__(self, target_type='binary',
9                   encoding='likelihood', func=None):
10 ▾         if target_type == 'continuous' and encoding in ['woe', 'diff']:
11 ▾             raise ValueError(
12                 '{} target_type can\'t be used with {} encoding'.format(target_type, encoding))
13
14         self.target_type = target_type
15         self.encoding = encoding
16         self.func = func
17
18 ▾     def goods(self, x):
19         return np.sum(x == 1)
20
21 ▾     def bads(self, x):
22         return np.sum(x == 0)
23
24 ▾     def encode(self, X, y, agg_func):
25         self.means = dict()
26         self.global_mean = np.nan
27         X['target'] = y
28 ▾         for col in X.columns:
29 ▾             if col != 'target':
30                 col_means = X.groupby(col)['target'].agg(agg_func)
31                 self.means[col] = col_means
32         X.drop(['target'], axis=1, inplace = True)
33
34 ▾     def fit(self, X, y):
35 ▾         if self.encoding == 'woe':
36             self.encode(X, y, lambda x: np.log(self.goods(x) / self.bads(x)))
37             self.global_mean = np.log(self.goods(y) / self.bads(y)) * 100
38 ▾         elif self.encoding == 'diff':
39             self.encode(X, y, lambda x: self.goods(x) - self.bads(x))
40             self.global_mean = self.goods(y) - self.bads(y)
41 ▾         elif self.encoding == 'likelihood':
42             self.encode(X, y, np.mean)
43             self.global_mean = np.mean(y)
44 ▾         elif self.encoding == 'count':
45             self.encode(X, y, np.sum)
46             self.global_mean = np.sum(y)
47 ▾         elif self.encoding == 'function':
48             self.encode(X, y, lambda x: self.func(x))
49             self.global_mean = self.func(y)
50         return self
51
52 ▾     def transform(self, X):
53         X_new = pd.DataFrame()
54 ▾         for col in X.columns:
55             X_new[col] = X[col].map(self.means[col]).fillna(self.global_mean)
56         return X_new
57
58 ▾     def fit_transform(self, X, y):
59         self.fit(X, y)
```

```
60         return self.transform(X)
```

started 01:44:01 2020-03-27, finished in 22ms

Сохраним оригинальные данные

In [11]:

```
1 x_train_origin = x_train.copy()
2 x_valid_origin = x_valid.copy()
```

started 01:44:01 2020-03-27, finished in 9ms

Закодируем категориальный признак

In [12]:

```
1 encoder = MeanEncoder()
2 x_train[cat_features] = encoder.fit_transform(x_train[cat_features], y_train)
3 x_valid[cat_features] = encoder.transform(x_valid[cat_features])
```

started 01:44:01 2020-03-27, finished in 26ms

Некоторая вспомогательная функция для отрисовки графиков

In [14]:

```
1 colors = ['#FF3300', '#0099CC', '#00CC66', 'orange']
2
3 def plot_dependence_test(param_grid, mape_train, mape_valid, descr,
4                           param_label,
5                           title, ylim=(11, 16)):
6     '''
7     Функция для построения графиков зависимости целевой метрики
8     от некоторого параметра моделей на обучающей и валидационной выборках.
9
10    Параметры.
11    1) param_grid - значения исследуемого параметра,
12    2) mape_train - значения метрик на обучающей выборке,
13    3) mape_valid - значения метрик на валидационной выборке,
14    4) descr - описания моделей для легенды
15    5) param_label - названия параметра,
16    6) title - заголовок для графика.
17    '''
18
19    plt.figure(figsize=(16, 8))
20
21    for i in range(len(mape_train)):
22        plt.plot(param_grid, mape_train[i], color=colors[i],
23                linewidth=2, linestyle='--', alpha=0.7,
24                label=descr[i]+' train')
25        plt.plot(param_grid, mape_valid[i], color=colors[i],
26                linewidth=4, alpha=0.7, label=descr[i]+' valid')
27
28    plt.xlabel(param_label)
29    plt.ylabel('MAPE, %')
30    plt.legend(ncol=2)
31    plt.title(title)
32    plt.ylim(ylim)
33    plt.show()
```

started 01:53:26 2020-03-27, finished in 6ms

1. Зависимость от количества деревьев

Для каждой модели используем функции получения предсказания с ограничением количества деревьев при предсказаниях. Это позволяет не обучать каждый раз одни и те же модели.

In [13]:

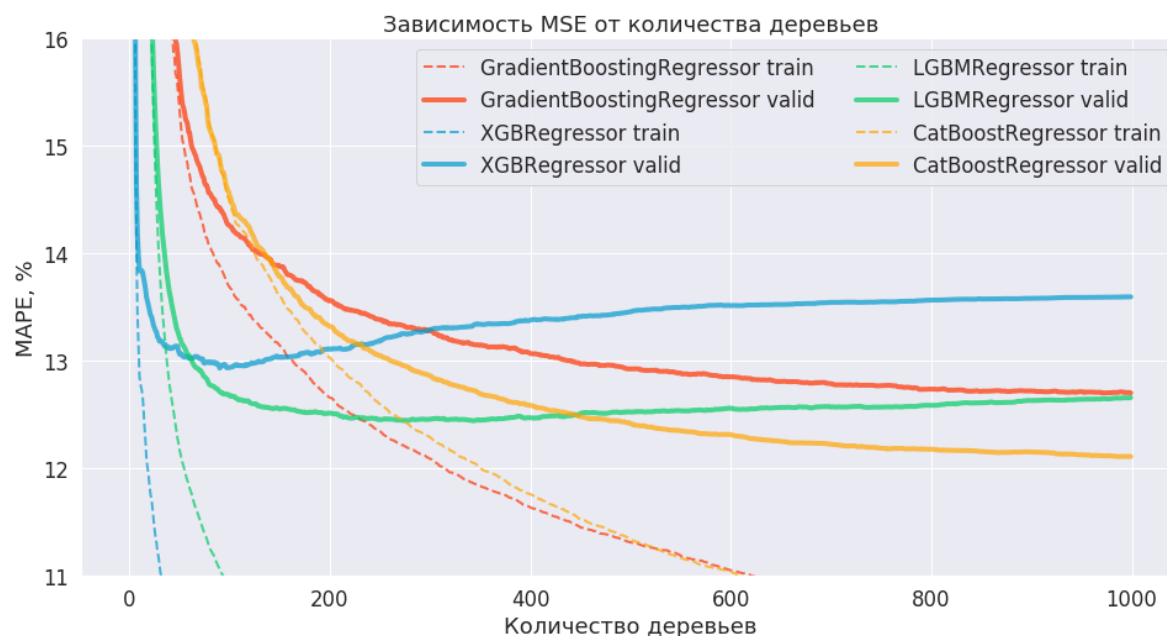
```
1  n_estimators_grid = np.arange(1, 1000)
2
3  regressor = GradientBoostingRegressor(n_estimators=max(n_estimators_grid))
4  regressor.fit(x_train, y_train)
5  mape_train = [[mape(y_train, prediction)
6                 for prediction in regressor.staged_predict(x_train)]]
7  mape_valid = [[mape(y_valid, prediction) for
8                 prediction in regressor.staged_predict(x_valid)]]
9
10 regressor = XGBRegressor(n_estimators=max(n_estimators_grid))
11 regressor.fit(x_train, y_train)
12 mape_train += [[mape(y_train, regressor.predict(x_train, ntree_limit=i+1))
13                for i in range(max(n_estimators_grid))]]
14 mape_valid += [[mape(y_valid, regressor.predict(x_valid, ntree_limit=i+1))
15                for i in range(max(n_estimators_grid))]]
16
17 regressor = LGBMRegressor(n_estimators=max(n_estimators_grid))
18 regressor.fit(x_train_origin, y_train, categorical_feature=cat_features)
19 mape_train += [[mape(y_train, regressor.predict(x_train_origin,
20                                                  num_iteration=i+1))
21                for i in range(max(n_estimators_grid))]]
22 mape_valid += [[mape(y_valid, regressor.predict(x_valid_origin,
23                                                  num_iteration=i+1))
24                for i in range(max(n_estimators_grid))]]
25
26 regressor = CatBoostRegressor(n_estimators=max(n_estimators_grid),
27                               cat_features=cat_features, verbose=0)
28 regressor.fit(x_train_origin, y_train)
29 mape_train += [[mape(y_train, regressor.predict(x_train_origin, ntree_end=i+1)
30                                                  for i in range(max(n_estimators_grid))]]
31 mape_valid += [[mape(y_valid, regressor.predict(x_valid_origin, ntree_end=i+1)
32                                                  for i in range(max(n_estimators_grid))]]
```

started 01:44:01 2020-03-27, finished in 9m 25s

In [15]:

```
1 plot_dependence_test(n_estimators_grid, mape_train, mape_valid,
2                       ['GradientBoostingRegressor', 'XGBRegressor',
3                       'LGBMRegressor', 'CatBoostRegressor'],
4                       'Количество деревьев',
5                       'Зависимость MSE от количества деревьев')
```

started 01:53:26 2020-03-27, finished in 574ms



Вывод: XGBoost и LightGBM позволяют быстро получить довольно приемлимое качество на небольшом количестве деревьев. Однако при увеличении количества деревьев они начинают переобучаться. CatBoost'у требуется большое количество деревьев для получения хорошего качества, причем при увеличении количества деревьев не наблюдается переобучения. Подобные свойства наблюдаются у реализации из sklearn. Наилучшее качество на валидационной выборке позволяет получить CatBoost.

2. Зависимость от максимальной глубины дерева

In [16]:

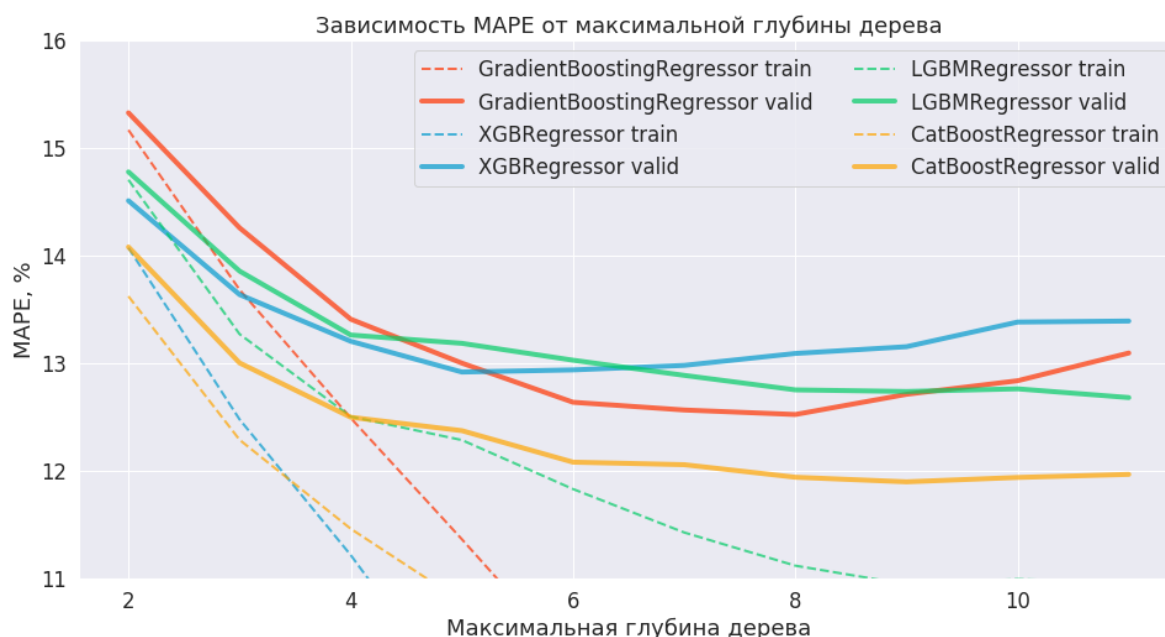
```
1 max_depth_grid = np.arange(2, 12)
2
3 mape_train, mape_valid = np.zeros((2, 4, len(max_depth_grid)))
4
5 for i, max_depth in enumerate(max_depth_grid):
6     regressor = GradientBoostingRegressor(max_depth=max_depth)
7     regressor.fit(x_train, y_train)
8     mape_train[0, i] = mape(y_train, regressor.predict(x_train))
9     mape_valid[0, i] = mape(y_valid, regressor.predict(x_valid))
10
11     regressor = XGBRegressor(max_depth=max_depth)
12     regressor.fit(x_train, y_train)
13     mape_train[1, i] = mape(y_train, regressor.predict(x_train))
14     mape_valid[1, i] = mape(y_valid, regressor.predict(x_valid))
15
16     regressor = LGBMRegressor(max_depth=max_depth)
17     regressor.fit(x_train_origin, y_train, categorical_feature=cat_features)
18     mape_train[2, i] = mape(y_train, regressor.predict(x_train_origin))
19     mape_valid[2, i] = mape(y_valid, regressor.predict(x_valid_origin))
20
21     regressor = CatBoostRegressor(max_depth=max_depth,
22                                   cat_features=cat_features, verbose=0)
23     regressor.fit(x_train_origin, y_train)
24     mape_train[3, i] = mape(y_train, regressor.predict(x_train_origin))
25     mape_valid[3, i] = mape(y_valid, regressor.predict(x_valid_origin))
```

started 01:53:27 2020-03-27, finished in 2m 29s

In [17]:

```
1 plot_dependence_test(max_depth_grid, mape_train, mape_valid,
2                       ['GradientBoostingRegressor', 'XGBRegressor',
3                       'LGBMRegressor', 'CatBoostRegressor'],
4                       'Максимальная глубина дерева',
5                       'Зависимость MAPE от максимальной глубины дерева')
```

started 01:55:56 2020-03-27, finished in 508ms



Вывод: Оптимальная глубина дерева для разных моделей равна от 5 до 8. Модели XGBoost и sklearn-реализация начинают переобучаться с увеличением глубины дерева. Модели LightGBM и CatBoost ведут

себя более устойчиво.

3. Зависимость от learning_rate

In [18]:

```
1 learning_rate_grid = np.linspace(0.0125, 0.6, 48)
2 n_estimators = 100
3
4 mape_train, mape_valid = np.zeros((2, 4, len(learning_rate_grid)))
5
6 for i, learning_rate in enumerate(learning_rate_grid):
7     regressor = GradientBoostingRegressor(learning_rate=learning_rate)
8     regressor.fit(x_train, y_train)
9     mape_train[0, i] = mape(y_train, regressor.predict(x_train))
10    mape_valid[0, i] = mape(y_valid, regressor.predict(x_valid))
11
12    regressor = XGBRegressor(learning_rate=learning_rate)
13    regressor.fit(x_train, y_train)
14    mape_train[1, i] = mape(y_train, regressor.predict(x_train))
15    mape_valid[1, i] = mape(y_valid, regressor.predict(x_valid))
16
17    regressor = LGBMRegressor(learning_rate=learning_rate)
18    regressor.fit(x_train_origin, y_train, categorical_feature=cat_features)
19    mape_train[2, i] = mape(y_train, regressor.predict(x_train_origin))
20    mape_valid[2, i] = mape(y_valid, regressor.predict(x_valid_origin))
21
22    regressor = CatBoostRegressor(learning_rate=learning_rate,
23                                  cat_features=cat_features, verbose=0)
24    regressor.fit(x_train_origin, y_train)
25    mape_train[3, i] = mape(y_train, regressor.predict(x_train_origin))
26    mape_valid[3, i] = mape(y_valid, regressor.predict(x_valid_origin))
```

started 01:55:56 2020-03-27, finished in 5m 36s

In [19]:

```
1 plot_dependence_test(learning_rate_grid, mape_train, mape_valid,
2                       ['GradientBoostingRegressor', 'XGBRegressor',
3                       'LGBMRegressor', 'CatBoostRegressor'],
4                       'learning_rate',
5                       'Зависимость MAPE от learning_rate')
```

started 02:01:32 2020-03-27, finished in 586ms



Вывод: Для каждой модели наблюдается некоторое оптимальное значение `learning_rate`. При меньших значениях модели недообучаются, при больших значениях наблюдается неустойчивость моделей при обучении --- графики начинают колбасить. Модель CatBoost показывает наилучшее качество. CatBoost лучший :)