

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import statsmodels.api as sm
5 import seaborn as sns
6 from tqdm import tqdm_notebook
7 import warnings
8 from statsmodels.tsa.stattools import kpss
9 from statsmodels.stats.multitest import multipletests
10 from statsmodels.tsa.holtwinters import ExponentialSmoothing
11 from sklearn.preprocessing import StandardScaler
12
13 from xgboost import XGBClassifier
14 import lightgbm as lgb
15 from sklearn.model_selection import StratifiedKFold, KFold
16 from sklearn.metrics import mean_absolute_error
17 import time
18 from catboost import CatBoostRegressor, CatBoostClassifier
19
20 import lightgbm as lgb
21 from sklearn.model_selection import train_test_split
22 from sklearn.metrics import f1_score
23
24 import gc
25
26 from feature_engineering import reduce_mem_usage, add_rolling_features
27 from feature_engineering import exponential_smoothing, signal_shifts
28 from feature_engineering import batch_stats2, add_minus_signal
29 from feature_engineering import delete_objects_after_rolling
30 from feature_engineering import add_quantiles, add_target_encoding
31
32 from bayes_opt import BayesianOptimization
33
34 warnings.filterwarnings('ignore')
35 sns.set_style('whitegrid')
```

```

In [2]: 1 def pred_proc(pred):
2         pred = np.round(np.clip(pred, 0, 10))
3         return pred.astype(int)
4
5 def MacroF1Metric(preds, dtrain):
6     labels = dtrain.get_label()
7     preds = np.round(np.clip(preds, 0, 10)).astype(int)
8     score = f1_score(labels, preds, average = 'macro')
9     return ('MacroF1Metric', score, True)
10
11 def prepare_df(df, window_sizes, alphas, shifts, batch_sizes):
12     df = reduce_mem_usage(df)
13     df = add_rolling_features(df, window_sizes)
14     df = reduce_mem_usage(df)
15     df = exponential_smoothing(df, alphas)
16     df = reduce_mem_usage(df)
17     df = signal_shifts(df, shifts)
18     df = reduce_mem_usage(df)
19     df = batch_stats2(df, batch_sizes)
20     df = reduce_mem_usage(df)
21     df = add_minus_signal(df)
22     df = reduce_mem_usage(df)
23
24     if 'open_channels' in df.columns:
25         y = df['open_channels']
26         df = df.drop(columns=['time'])
27         return df, y
28     else:
29         df = df.drop(columns=['time'])
30         return df

```

## 0. Загрузка данных

```

In [2]: 1 train = pd.read_csv('data/train.csv')
2         test = pd.read_csv('data/test.csv')
3
4         print(train.shape)
5         print(test.shape)
6         print(train.head())

```

```

(5000000, 3)
(2000000, 2)
   time  signal  open_channels
0  0.0001 -2.7600             0
1  0.0002 -2.8557             0
2  0.0003 -2.4074             0
3  0.0004 -3.1404             0
4  0.0005 -3.1525             0

```

## 1. Подготовка датасета

```
In [3]: 1 from feature_engineering import reduce_mem_usage, add_rolling_features, expo
```

```
In [4]: 1 def prepare_df(df, window_sizes, alphas, shifts):
2     df = reduce_mem_usage(df)
3     df = add_rolling_features(df, window_sizes)
4     df = reduce_mem_usage(df)
5     df = exponential_smoothing(df, alphas)
6     df = reduce_mem_usage(df)
7     df = signal_shifts(df, shifts)
8     df = reduce_mem_usage(df)
9
10    if 'open_channels' in df.columns:
11        y = df['open_channels']
12        df = df.drop(columns=['time', 'open_channels'])
13        return df, y
14    else:
15        df = df.drop(columns=['time'])
16        return df
```

```
In [5]: 1 window_sizes = [5, 100, 5000]
2     alphas = [0.5, 0.1]
3     shifts = [1, 2, -1, -2]
4
5     X_train, y_train = prepare_df(train, window_sizes, alphas, shifts)
6     X_test = prepare_df(test, window_sizes, alphas, shifts)
```

Mem. usage decreased to 23.84 Mb (79.2% reduction)

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Mem. usage decreased to 286.10 Mb (74.6% reduction)

Mem. usage decreased to 305.18 Mb (15.8% reduction)

Mem. usage decreased to 343.32 Mb (16.3% reduction)

Mem. usage decreased to 7.63 Mb (75.0% reduction)

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Mem. usage decreased to 112.53 Mb (74.1% reduction)

Mem. usage decreased to 120.16 Mb (16.0% reduction)

Mem. usage decreased to 135.42 Mb (9.0% reduction)

## 2. Первая наивная попытка: LGBMRegressor

Попробуем обучить LGBMRegressor, разбивая данные на 5 фолдов, обучаясь как в кросс-валидации, при этом для каждого разбиения считать прогноз лучшей итерации на тесте.

```
In [6]: 1 X_train = np.array(X_train)
2 y_train = np.array(y_train)
3 X_test = np.array(X_test)
4
5 scaler = StandardScaler()
6 X_train = scaler.fit_transform(X_train)
7 X_test = scaler.transform(X_test)
```

```
In [4]: 1 params = {'num_leaves': 128,
2             'min_data_in_leaf': 64,
3             'objective': 'huber',
4             'max_depth': -1,
5             'learning_rate': 0.005,
6             "boosting": "gbdt",
7             "bagging_freq": 5,
8             "bagging_fraction": 0.8,
9             "bagging_seed": 11,
10            "metric": 'mae',
11            "verbosity": -1,
12            'reg_alpha': 0.1,
13            'reg_lambda': 0.3}
```

```
In [5]: 1 n_fold = 5
2 folds = KFold(n_splits=n_fold, shuffle=True, random_state=42)
```

```
In [10]: 1 oof = np.zeros(len(X_train))
2 prediction = np.zeros(len(X_test))
3 scores = []
4
5 for fold_n, (train_index, valid_index) in tqdm_notebook(enumerate(folds.split(X_train, y_train))):
6     print('Fold', fold_n, 'started at', time.ctime())
7     X_train_, X_valid = X_train[train_index], X_train[valid_index]
8     y_train_, y_valid = y_train[train_index], y_train[valid_index]
9
10    model = lgb.LGBMRegressor(**params, n_estimators = 5000, n_jobs = -1)
11    model.fit(X_train_, y_train_,
12              eval_set=[(X_train_, y_train_), (X_valid, y_valid)], eval_metric='mae',
13              verbose=500, early_stopping_rounds=200)
14
15    y_pred_valid = model.predict(X_valid)
16    y_pred = model.predict(X_test, num_iteration=model.best_iteration_)
17
18    oof[valid_index] = y_pred_valid.reshape(-1,)
19    scores.append(mean_absolute_error(y_valid, y_pred_valid))
20
21    prediction += y_pred
22
23 prediction /= n_fold
```

...

```
In [11]: 1 sample_df = pd.read_csv("data/sample_submission.csv", dtype={'time':str})
2
3 sample_df['open_channels'] = np.round(prediction).astype(np.int)
4 sample_df.to_csv("submission.csv", index=False, float_format='%.4f')
```

**Результат:** 0.71 на public LB

### 3. Возьмем другой датасет

```
In [3]: 1 def prepare_df(df, window_sizes, alphas, shifts, batch_sizes):
2     df = reduce_mem_usage(df)
3     df = add_rolling_features(df, window_sizes)
4     df = reduce_mem_usage(df)
5     df = exponential_smoothing(df, alphas)
6     df = reduce_mem_usage(df)
7     df = signal_shifts(df, shifts)
8     df = reduce_mem_usage(df)
9     df = batch_stats2(df, batch_sizes)
10    df = reduce_mem_usage(df)
11    df = add_minus_signal(df)
12    df = reduce_mem_usage(df)
13
14    if 'open_channels' in df.columns:
15        y = df['open_channels']
16        df = df.drop(columns=['time', 'open_channels'])
17        return df, y
18    else:
19        df = df.drop(columns=['time'])
20        return df
```

```
In [4]: 1 train = pd.read_csv('data/train.csv')
2 test = pd.read_csv('data/test.csv')
3
4 window_sizes = [5, 100, 1000]
5 alphas = [0.5, 0.1]
6 shifts = [1, -1]
7 batch_sizes = [25000, 2500]
8
9
10 X_train, y_train = prepare_df(train, window_sizes, alphas, shifts, batch_size)
11 X_test = prepare_df(test, window_sizes, alphas, shifts, batch_sizes)
```

Mem. usage decreased to 23.84 Mb (79.2% reduction)

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Mem. usage decreased to 286.10 Mb (74.6% reduction)

Mem. usage decreased to 305.18 Mb (15.8% reduction)

Mem. usage decreased to 324.25 Mb (17.1% reduction)

Mem. usage decreased to 562.67 Mb (51.6% reduction)

Mem. usage decreased to 1096.73 Mb (0.0% reduction)

Mem. usage decreased to 7.63 Mb (75.0% reduction)

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Mem. usage decreased to 112.53 Mb (74.1% reduction)

Mem. usage decreased to 120.16 Mb (16.0% reduction)

Mem. usage decreased to 127.79 Mb (9.5% reduction)

Mem. usage decreased to 211.72 Mb (54.3% reduction)

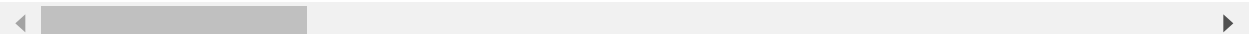
Mem. usage decreased to 413.89 Mb (0.0% reduction)

```
In [5]: 1 X_train.head()
```

Out[5]:

	signal	batch	rolling_mean_5	rolling_std_5	rolling_var_5	rolling_min_5	rolling_max_5	rolli
0	-2.759766	0	0.000000	0.000000	0.000000	0.000000	0.000000	
1	-2.855469	0	0.000000	0.000000	0.000000	0.000000	0.000000	
2	-2.408203	0	0.000000	0.000000	0.000000	0.000000	0.000000	
3	-3.140625	0	0.000000	0.000000	0.000000	0.000000	0.000000	
4	-3.152344	0	-2.863281	0.307617	0.094604	-3.152344	-2.408203	

5 rows × 108 columns



```
In [6]: 1 X_train = np.array(X_train)
        2 y_train = np.array(y_train)
        3 X_test = np.array(X_test)
```

## 4. Макро F1-метрика

Подсчет макро f1, которую будем выводить сразу же при обучении lgbm.

Подсмотрено в этом нутбуке: <https://www.kaggle.com/vbmokin/ion-switching-advanced-fe-lgb-xgb-confmatrix/notebook> (<https://www.kaggle.com/vbmokin/ion-switching-advanced-fe-lgb-xgb-confmatrix/notebook>)

```
In [3]: 1 def MacroF1Metric(preds, dtrain):
        2     labels = dtrain.get_label()
        3     preds = np.round(np.clip(preds, 0, 10)).astype(int)
        4     score = f1_score(labels, preds, average = 'macro')
        5     return ('MacroF1Metric', score, True)
```

## 5. Обучение lgbm

Параметры подбираем пока что чисто интуитивно.

```
In [10]: 1 num_iterations = 3000
2 X_train_, X_valid, y_train_, y_valid = train_test_split(X_train, y_train, te
3 params = {'learning_rate': 0.05,
4           'max_depth': -1,
5           'num_leaves': 200,
6           'metric': 'logloss',
7           'random_state': 17,
8           'n_jobs': -1,
9           'sample_fraction': 0.33}
10 model = lgb.train(params, lgb.Dataset(X_train_, y_train_),
11                   num_iterations,
12                   lgb.Dataset(X_valid, y_valid),
13                   verbose_eval=100,
14                   early_stopping_rounds=200,
15                   feval=MacroF1Metric)
16 model.save_model('model5.txt', num_iteration=model.best_iteration)
```

Training until validation scores don't improve for 200 rounds

```
[100] valid_0's MacroF1Metric: 0.931818
[200] valid_0's MacroF1Metric: 0.934456
[300] valid_0's MacroF1Metric: 0.935061
[400] valid_0's MacroF1Metric: 0.935205
[500] valid_0's MacroF1Metric: 0.935294
[600] valid_0's MacroF1Metric: 0.935436
[700] valid_0's MacroF1Metric: 0.935445
[800] valid_0's MacroF1Metric: 0.935532
[900] valid_0's MacroF1Metric: 0.935671
[1000] valid_0's MacroF1Metric: 0.935628
[1100] valid_0's MacroF1Metric: 0.935719
[1200] valid_0's MacroF1Metric: 0.935687
[1300] valid_0's MacroF1Metric: 0.935751
[1400] valid_0's MacroF1Metric: 0.935708
Early stopping, best iteration is:
[1287] valid_0's MacroF1Metric: 0.935771
```

Out[10]: 132

```
In [12]: 1 num_iterations = 2000
2 model = lgb.train(params, lgb.Dataset(X_train, y_train),
3                   num_iterations,
4                   verbose_eval=100,
5                   feval=MacroF1Metric)
```

```
In [14]: 1 %%time
2
3 y_lgb_pred = model.predict(X_test)
```

Wall time: 1min 6s

```
In [15]: 1 def pred_proc(pred):
2         pred = np.round(np.clip(pred, 0, 10))
3         return pred.astype(int)
```



```
In [16]: 1 y_lgb_pred = pred_proc(y_lgb_pred)
```

```
In [24]: 1 sample_df = pd.read_csv("data/sample_submission.csv", dtype={'time':str})
2
3 sample_df['open_channels'] = y_lgb_pred
4 sample_df.to_csv("model5.csv", index=False, float_format='%.4f')
```

Результат: 0.932 на public lb

## 6. LGBM с БОльшим кол-вом фичей

```
In [4]: 1 train = pd.read_csv('data/train.csv')
2 test = pd.read_csv('data/test.csv')
3
4 window_sizes = [5, 100, 1000, 5000]
5 alphas = [0.5, 0.2, 0.05]
6 shifts = [1, -1, 2, -2]
7 batch_sizes = [50000, 25000, 2500]
8
9
10 X_train, y_train = prepare_df(train, window_sizes, alphas, shifts, batch_size)
11 X_test = prepare_df(test, window_sizes, alphas, shifts, batch_sizes)
12
13 # X_train = np.array(X_train)
14 y_train = np.array(y_train)
15 # X_test = np.array(X_test)
```

Mem. usage decreased to 23.84 Mb (79.2% reduction)

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Mem. usage decreased to 371.93 Mb (74.7% reduction)

Mem. usage decreased to 400.54 Mb (17.6% reduction)

Mem. usage decreased to 438.69 Mb (13.2% reduction)

Mem. usage decreased to 782.01 Mb (53.9% reduction)

Mem. usage decreased to 1535.42 Mb (0.0% reduction)

Mem. usage decreased to 7.63 Mb (75.0% reduction)

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Mem. usage decreased to 146.87 Mb (74.3% reduction)

Mem. usage decreased to 158.31 Mb (17.8% reduction)

Mem. usage decreased to 173.57 Mb (7.1% reduction)

Mem. usage decreased to 299.45 Mb (55.8% reduction)

Mem. usage decreased to 589.37 Mb (0.0% reduction)

```

In [28]: 1 %%time
2
3 num_iterations = 3000
4 X_train_, X_valid, y_train_, y_valid = train_test_split(X_train, y_train, te
5 params = {'learning_rate': 0.05,
6           'max_depth': -1,
7           'num_leaves': 200,
8           'metric': 'logloss',
9           'random_state': 17,
10          'n_jobs': -1,
11          'sample_fraction': 0.33}
12 model = lgb.train(params, lgb.Dataset(X_train_, y_train_),
13                  num_iterations,
14                  lgb.Dataset(X_valid, y_valid),
15                  verbose_eval=100,
16                  early_stopping_rounds=200,
17                  feval=MacroF1Metric)
18 model.save_model('model6.txt', num_iteration=model.best_iteration)

```

Training until validation scores don't improve for 200 rounds

```

[100] valid_0's MacroF1Metric: 0.933017
[200] valid_0's MacroF1Metric: 0.935251
[300] valid_0's MacroF1Metric: 0.935518
[400] valid_0's MacroF1Metric: 0.935647
[500] valid_0's MacroF1Metric: 0.935838
[600] valid_0's MacroF1Metric: 0.93596
[700] valid_0's MacroF1Metric: 0.936037
[800] valid_0's MacroF1Metric: 0.936087
[900] valid_0's MacroF1Metric: 0.936211
[1000] valid_0's MacroF1Metric: 0.936281
[1100] valid_0's MacroF1Metric: 0.936396
[1200] valid_0's MacroF1Metric: 0.936453
[1300] valid_0's MacroF1Metric: 0.936567
[1400] valid_0's MacroF1Metric: 0.936519
[1500] valid_0's MacroF1Metric: 0.936491
Early stopping, best iteration is:
[1369] valid_0's MacroF1Metric: 0.936582
Wall time: 40min 54s

```

Out[28]: <lightgbm.basic.Booster at 0x27b90bcee48>

```

In [29]: 1 %%time
2
3 num_iterations = 2000
4 model = lgb.train(params, lgb.Dataset(X_train, y_train),
5                  num_iterations,
6                  verbose_eval=100,
7                  feval=MacroF1Metric)
8
9 y_lgb_pred = model.predict(X_test)

```

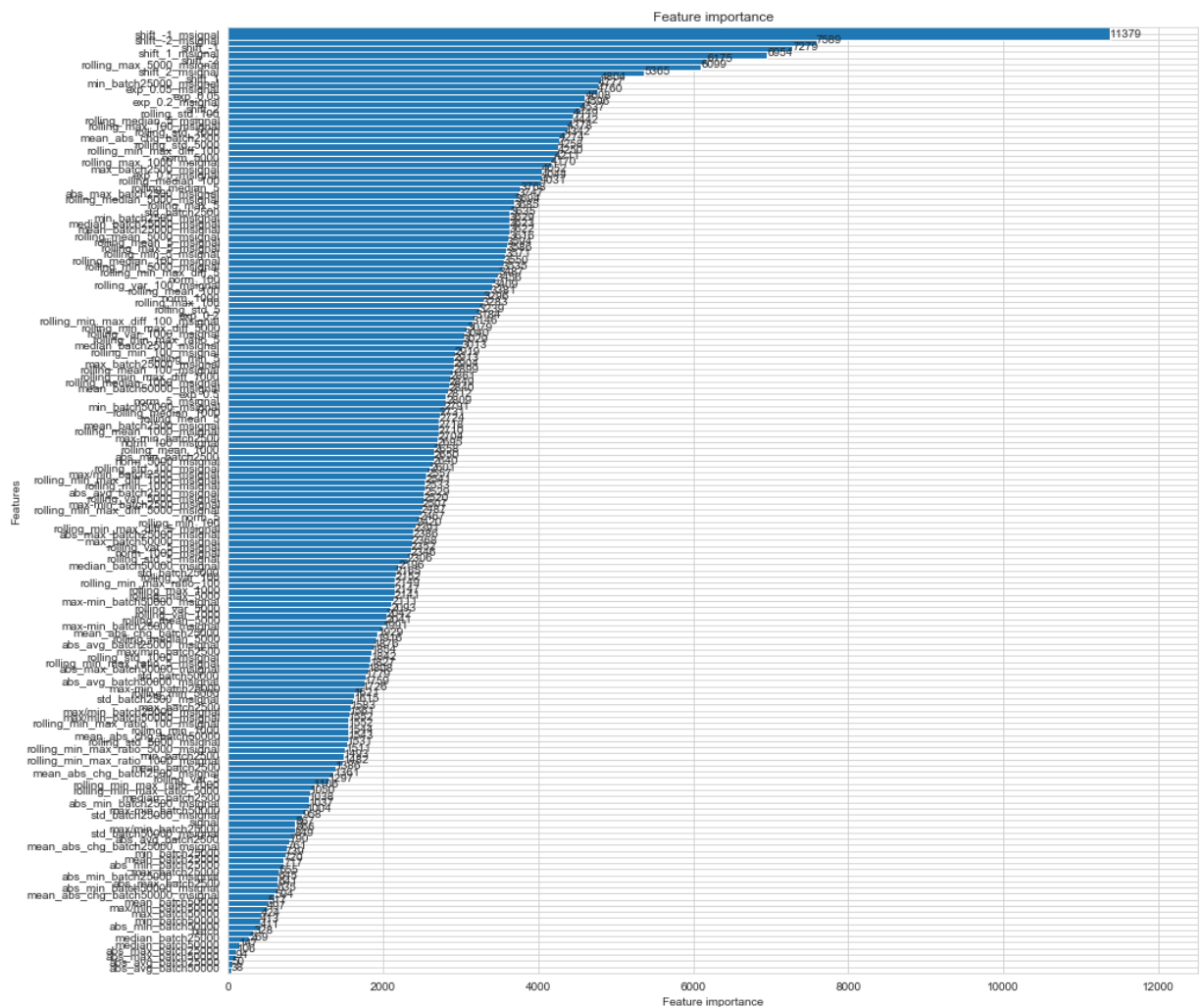
Wall time: 21min 15s

```
In [30]: 1 %%time
2 y_lgb_pred = pred_proc(y_lgb_pred)
3 sample_df = pd.read_csv("data/sample_submission.csv", dtype={'time':str})
4 sample_df['open_channels'] = y_lgb_pred
5 sample_df.to_csv("model6.csv", index=False, float_format='%0.4f')
```

Wall time: 4.39 s

**Резултат:** 0.936 на public LB

```
In [33]: 1 fig = plt.figure(figsize = (15,15))
          2 axes = fig.add_subplot(111)
          3 lgb.plot_importance(model,ax = axes,height = 2.)
          4 plt.show()
```



Как видим, здорово себя показывают эксп. сглаживания, сдвиги и роллинги. Сделаем упор на них в следующей модели.

## 7. Еще больше фичей

```

In [5]: 1 train = pd.read_csv('data/train.csv')
        2 test = pd.read_csv('data/test.csv')
        3
        4 window_sizes = [5, 100, 1000, 5000, 10000]
        5 alphas = [0.5, 0.2, 0.05, 0.7]
        6 shifts = [1, -1, 2, -2, 3, -3]
        7 batch_sizes = [50000, 25000, 2500]
        8
        9
       10 X_train, y_train = prepare_df(train, window_sizes, alphas, shifts, batch_size)
       11 X_test = prepare_df(test, window_sizes, alphas, shifts, batch_sizes)
       12
       13 # X_train = np.array(X_train)
       14 y_train = np.array(y_train)
       15 # X_test = np.array(X_test)

```

Mem. usage decreased to 23.84 Mb (79.2% reduction)

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Mem. usage decreased to 457.76 Mb (74.7% reduction)  
 Mem. usage decreased to 495.91 Mb (18.7% reduction)  
 Mem. usage decreased to 553.13 Mb (10.8% reduction)  
 Mem. usage decreased to 896.45 Mb (50.5% reduction)  
 Mem. usage decreased to 1764.30 Mb (0.0% reduction)  
 Mem. usage decreased to 7.63 Mb (75.0% reduction)

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Mem. usage decreased to 181.20 Mb (74.5% reduction)  
 Mem. usage decreased to 196.46 Mb (18.9% reduction)  
 Mem. usage decreased to 219.35 Mb (5.7% reduction)  
 Mem. usage decreased to 345.23 Mb (52.2% reduction)  
 Mem. usage decreased to 680.92 Mb (0.0% reduction)

```

In [15]: 1 X_train = delete_objects_after_rolling(X_train, 10000)
        2 add_quantiles(X_train, X_test, n_bins=7)

```

```

In [16]: 1 tmp = []
        2 for i, y in enumerate(y_train):
        3     if i % 500000 < 10000:
        4         continue
        5     else:
        6         tmp.append(y)

```

```

In [17]: 1 y_train = np.array(tmp)

```

```

In [18]: 1 %%time
2
3 num_iterations = 3000
4 X_train_, X_valid, y_train_, y_valid = train_test_split(X_train, y_train, te
5 params = {'learning_rate': 0.05,
6           'max_depth': -1,
7           'num_leaves': 200,
8           'metric': 'logloss',
9           'random_state': 17,
10          'n_jobs': -1,
11          'sample_fraction': 0.33}
12 model = lgb.train(params, lgb.Dataset(X_train_, y_train_),
13                  num_iterations,
14                  lgb.Dataset(X_valid, y_valid),
15                  verbose_eval=100,
16                  early_stopping_rounds=200,
17                  feval=MacroF1Metric)
18 model.save_model('model7.txt', num_iteration=model.best_iteration)

```

Training until validation scores don't improve for 200 rounds

```

[100] valid_0's MacroF1Metric: 0.93366
[200] valid_0's MacroF1Metric: 0.935862
[300] valid_0's MacroF1Metric: 0.936162
[400] valid_0's MacroF1Metric: 0.936317
[500] valid_0's MacroF1Metric: 0.936515
[600] valid_0's MacroF1Metric: 0.936509
[700] valid_0's MacroF1Metric: 0.936523
[800] valid_0's MacroF1Metric: 0.936641
[900] valid_0's MacroF1Metric: 0.936617
[1000] valid_0's MacroF1Metric: 0.936656
[1100] valid_0's MacroF1Metric: 0.936691
[1200] valid_0's MacroF1Metric: 0.936778
[1300] valid_0's MacroF1Metric: 0.936782
[1400] valid_0's MacroF1Metric: 0.936859
[1500] valid_0's MacroF1Metric: 0.936927
[1600] valid_0's MacroF1Metric: 0.936938
[1700] valid_0's MacroF1Metric: 0.936913
[1800] valid_0's MacroF1Metric: 0.936934
[1900] valid_0's MacroF1Metric: 0.936877
[2000] valid_0's MacroF1Metric: 0.936855
Early stopping, best iteration is:
[1841] valid_0's MacroF1Metric: 0.936951
Wall time: 2h 30min 13s

```

Out[18]: <lightgbm.basic.Booster at 0x2c90007a9c8>

```
In [19]: 1 %%time
2
3 num_iterations = 2000
4 model = lgb.train(params, lgb.Dataset(X_train, y_train),
5                       num_iterations,
6                       verbose_eval=100,
7                       feval=MacroF1Metric)
8
9 y_lgb_pred = model.predict(X_test)
```

Wall time: 24min 10s

```
In [23]: 1 %%time
2 y_lgb_pred = pred_proc(y_lgb_pred)
3 sample_df = pd.read_csv("data/sample_submission.csv", dtype={'time':str})
4 sample_df['open_channels'] = y_lgb_pred
5 sample_df.to_csv("model7.csv", index=False, float_format='%.4f')
```

Wall time: 4.2 s

**Результат:** 0.935 на public LB.

## 9. Обучим такой же lgbm, но на данных без дрифта

Датасет без дрифта взят отсюда: <https://www.kaggle.com/c/liverpool-ion-switching/discussion/135480> (<https://www.kaggle.com/c/liverpool-ion-switching/discussion/135480>)

```
In [3]: 1 train = pd.read_csv('data/train_clean.csv')
2 test = pd.read_csv('data/test_clean.csv')
3
4 window_sizes = [5, 100, 1000, 5000]
5 alphas = [0.5, 0.2, 0.05]
6 shifts = [1, -1, 2, -2]
7 batch_sizes = [50000, 25000, 2500]
8
9
10 X_train, y_train = prepare_df(train, window_sizes, alphas, shifts, batch_size=batch_sizes)
11 X_test = prepare_df(test, window_sizes, alphas, shifts, batch_size=batch_sizes)
12
13 # X_train = np.array(X_train)
14 y_train = np.array(y_train)
15 # X_test = np.array(X_test)
16
17 # X_train = delete_objects_after_rolling(X_train, 10000)
18 # add_quantiles(X_train, X_test, n_bins=7)
```

...

```

In [8]: 1 %%time
2
3 num_iterations = 3000
4 X_train_, X_valid, y_train_, y_valid = train_test_split(X_train, y_train, te
5 params = {'learning_rate': 0.05,
6           'max_depth': -1,
7           'num_leaves': 200,
8           'metric': 'logloss',
9           'random_state': 17,
10          'n_jobs': -1,
11          'sample_fraction': 0.33}
12 model = lgb.train(params, lgb.Dataset(X_train_, y_train_),
13                  num_iterations,
14                  lgb.Dataset(X_valid, y_valid),
15                  verbose_eval=100,
16                  early_stopping_rounds=200,
17                  feval=MacroF1Metric)
18 model.save_model('model9.txt', num_iteration=model.best_iteration)

```

Training until validation scores don't improve for 200 rounds

```

[100] valid_0's MacroF1Metric: 0.937193
[200] valid_0's MacroF1Metric: 0.937541
[300] valid_0's MacroF1Metric: 0.937565
[400] valid_0's MacroF1Metric: 0.937608
[500] valid_0's MacroF1Metric: 0.93763
[600] valid_0's MacroF1Metric: 0.937583
[700] valid_0's MacroF1Metric: 0.937556
Early stopping, best iteration is:
[514] valid_0's MacroF1Metric: 0.937659
Wall time: 23min 19s

```

Out[8]: <lightgbm.basic.Booster at 0x2538008a8c8>

```

In [10]: 1 %%time
2
3 num_iterations = 600
4 model = lgb.train(params, lgb.Dataset(X_train, y_train),
5                  num_iterations)
6
7 y_lgb_pred = model.predict(X_test)

```

Wall time: 8min 18s

```

In [11]: 1 y_lgb_pred = pred_proc(y_lgb_pred)
2 sample_df = pd.read_csv("data/sample_submission.csv", dtype={'time':str})
3 sample_df['open_channels'] = y_lgb_pred
4 sample_df.to_csv("model9.csv", index=False, float_format='%.4f')

```

Результат: 0.939 на public lb

## 10. Добавим фичи с таргет енкодингом

```
In [33]: 1 train = pd.read_csv('data/train_clean.csv')
2 test = pd.read_csv('data/test_clean.csv')
3
4 window_sizes = [5, 100, 1000, 5000]
5 alphas = [0.5, 0.2, 0.05]
6 shifts = [1, -1, 2, -2]
7 batch_sizes = [50000, 25000, 2500]
8
9
10 X_train, y_train = prepare_df(train, window_sizes, alphas, shifts, batch_size=batch_sizes)
11 X_test = prepare_df(test, window_sizes, alphas, shifts, batch_size=batch_sizes)
12
13 y_train = np.array(y_train)
14
15 add_quantiles(X_train, X_test, [3, 7, 15])
16 add_target_encoding(X_train, X_test, [3, 7, 15])
17
18 X_train = reduce_mem_usage(X_train)
19 X_test = reduce_mem_usage(X_test)
20
21 X_train = X_train.drop(columns=['open_channels'])
22
23 gc.collect()
```

Mem. usage decreased to 23.84 Mb (79.2% reduction)

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Mem. usage decreased to 381.47 Mb (74.0% reduction)  
 Mem. usage decreased to 410.08 Mb (17.3% reduction)  
 Mem. usage decreased to 448.23 Mb (13.0% reduction)  
 Mem. usage decreased to 762.94 Mb (55.3% reduction)  
 Mem. usage decreased to 1497.27 Mb (0.0% reduction)  
 Mem. usage decreased to 7.63 Mb (75.0% reduction)

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Mem. usage decreased to 154.50 Mb (73.0% reduction)  
 Mem. usage decreased to 165.94 Mb (17.1% reduction)  
 Mem. usage decreased to 181.20 Mb (6.9% reduction)  
 Mem. usage decreased to 307.08 Mb (55.2% reduction)  
 Mem. usage decreased to 604.63 Mb (0.0% reduction)

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.



the notebook. You may be able to create the widget by running the appropriate cells.

Mem. usage decreased to 1587.87 Mb (15.9% reduction)

Mem. usage decreased to 640.87 Mb (15.8% reduction)

Out[33]: 6

```
In [4]: 1 %%time
2 gc.collect()
3 num_iterations = 3000
4 X_train_, X_valid, y_train_, y_valid = train_test_split(X_train, y_train, te
5 params = {'learning_rate': 0.05,
6           'max_depth': -1,
7           'num_leaves': 200,
8           'metric': 'logloss',
9           'random_state': 17,
10          'n_jobs': -1,
11          'sample_fraction': 0.33}
12
13 del X_train, y_train
14
15 train_set = lgb.Dataset(X_train_, y_train_)
16 val_set = lgb.Dataset(X_valid, y_valid)
17
18 del X_train_, X_valid, y_train_, y_valid
19
20 gc.collect()
21
22 model = lgb.train(params, train_set,
23                  num_iterations,
24                  val_set,
25                  verbose_eval=100,
26                  early_stopping_rounds=200,
27                  feval=MacroF1Metric)
28 model.save_model('model10.txt', num_iteration=model.best_iteration)
```

Training until validation scores don't improve for 200 rounds

[100] valid\_0's MacroF1Metric: 0.937142

[200] valid\_0's MacroF1Metric: 0.937563

[300] valid\_0's MacroF1Metric: 0.937675

[400] valid\_0's MacroF1Metric: 0.937667

[500] valid\_0's MacroF1Metric: 0.93772

[600] valid\_0's MacroF1Metric: 0.937751

[700] valid\_0's MacroF1Metric: 0.937672

[800] valid\_0's MacroF1Metric: 0.937705

Early stopping, best iteration is:

[609] valid\_0's MacroF1Metric: 0.937757

Wall time: 20min 12s

Out[4]: <lightgbm.basic.Booster at 0x263b56c1a48>



```
In [10]: 1 %%time
2
3 train_set = lgb.Dataset(X_train, y_train)
4
5 del X_train, y_train
6
7 gc.collect()
8
9 num_iterations = 600
10 model = lgb.train(params, train_set,
11                   num_iterations)
12
13 y_lgb_pred = model.predict(X_test)
14
15 np.save('model10_test_preds.npy', y_lgb_pred)
```

Wall time: 8min 14s

```
In [11]: 1 y_lgb_pred = pred_proc(y_lgb_pred)
2 sample_df = pd.read_csv("data/sample_submission.csv", dtype={'time':str})
3 sample_df['open_channels'] = y_lgb_pred
4 sample_df.to_csv("model10.csv", index=False, float_format='%.4f')
```

**Результат:** 0.939 на public lb

## 11. cv-loop

```

In [36]: 1 def lgb_cv_loop(X_train, y_train, X_test):
2         n_fold = 5
3         folds = KFold(n_splits=n_fold, shuffle=True, random_state=17)
4
5         num_iterations = 900
6
7         oof = np.zeros(len(X_train))
8         prediction = np.zeros(len(X_test))
9         scores = []
10
11        for training_index, validation_index in tqdm_notebook(folds.split(X_train, y_train)):
12            gc.collect()
13
14            # разбиение на трэйн и валидацию
15            X_train_ = X_train.iloc[training_index]
16            y_train_ = y_train[training_index]
17            X_valid = X_train.iloc[validation_index]
18            y_valid = y_train[validation_index]
19
20            train_set = lgb.Dataset(X_train_, y_train_)
21
22            del X_train_, y_train_
23
24            gc.collect()
25
26            # обучение модели
27            model = lgb.train(params, train_set, num_iterations)
28
29            # скор на валидации
30            preds = model.predict(X_valid)
31            oof[validation_index] = preds.reshape(-1,)
32            preds = np.round(np.clip(preds, 0, 10)).astype(int)
33            score = f1_score(y_valid, preds, average = 'macro')
34            scores.append(score)
35
36            # предсказание на тесте
37            preds = model.predict(X_test)
38            prediction += preds
39
40            print(f'score: {score}')
41
42        prediction /= n_fold
43        prediction = np.round(np.clip(prediction, 0, 10)).astype(int)
44
45        return scores, oof, prediction

```

```

In [22]: 1 params = {'learning_rate': 0.05,
2               'max_depth': -1,
3               'num_leaves': 200,
4               'metric': 'logloss',
5               'random_state': 17,
6               'n_jobs': -1,
7               'sample_fraction': 0.33}

```

```
In [23]: 1 %%time
          2 scores, oof, prediction = lgb_cv_loop(X_train, y_train, X_test)
```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

```
score: 0.9377723991122747
score: 0.9377418056780993
score: 0.9374989672286218
score: 0.9393080850181423
score: 0.9385557886223648
```

```
Wall time: 34min 57s
```

```
In [24]: 1 np.mean(scores)
```

```
Out[24]: 0.9381754091319007
```

```
In [25]: 1 np.save('model10_oof_preds.npy', oof)
```

## 12. Bayesian optimization

```

In [16]: 1 def bayesian_opt_lgbm(X, y, init_iter=3, n_iters=7, random_state=11, seed=10
2         dtrain = lgb.Dataset(data=X, label=y)
3
4         def hyp_lgbm(num_leaves, feature_fraction, bagging_fraction,
5                     max_depth, min_split_gain, min_child_weight, num_iterations
6
7                     params = {'application': 'regression',
8                               'learning_rate': 0.05,
9                               'early_stopping_round': 100,
10                              'metric': 'MacroF1Metric',
11                              'n_jobs': -1}
12                              #параметры для гри (убрать при использовании сри)
13                              #'device': 'gpu',
14                              #'gpu_platform_id': 0,
15                              #'gpu_device_id': 0}
16
17                     params["num_iterations"] = int(round(num_iterations))
18                     params["num_leaves"] = int(round(num_leaves))
19                     params['feature_fraction'] = max(min(feature_fraction, 1), 0)
20                     params['bagging_fraction'] = max(min(bagging_fraction, 1), 0)
21                     params['max_depth'] = int(round(max_depth))
22                     params['min_split_gain'] = min_split_gain
23                     params['min_child_weight'] = min_child_weight
24                     cv_results = lgb.cv(params, dtrain, nfold=5, seed=seed,
25                                         categorical_feature=[], stratified=False,
26                                         verbose_eval=None, feval=MacroF1Metric)
27
28                     return np.max(cv_results['MacroF1Metric-mean'])
29
30         pds = {'num_leaves': (80, 220),
31               'feature_fraction': (0.1, 0.9),
32               'bagging_fraction': (0.8, 1),
33               'max_depth': (15, 30),
34               'min_split_gain': (0.001, 0.1),
35               'min_child_weight': (10, 25),
36               'num_iterations': (100, 1000)
37               }
38
39         optimizer = BayesianOptimization(hyp_lgbm, pds, random_state=random_state)
40
41         optimizer.maximize(init_points=init_iter, n_iter=n_iters)
42
43         return optimizer

```

```
In [17]: 1 def lgbm_train(dtrain, best_params, num_iterations=200):
2
3     params = {'application': 'regression', 'num_iterations': num_iterations,
4               'learning_rate': 0.05,
5               'metric': 'MacroF1Metric',
6               'n_jobs': -1}
7
8     params["num_leaves"] = int(round(best_params['num_leaves']))
9     params['feature_fraction'] = max(min(best_params['feature_fraction'], 1)
10    params['bagging_fraction'] = max(min(best_params['bagging_fraction'], 1)
11    params['max_depth'] = int(round(best_params['max_depth']))
12    params['min_split_gain'] = best_params['min_split_gain']
13    params['min_child_weight'] = best_params['min_child_weight']
14
15    model = lgb.train(params, dtrain, 200, verbose_eval=100,
16                      feval=MacroF1Metric)
17
18    return model
```

```
In [18]: 1 op = bayesian_opt_lgbm(X_train, y_train, init_iter=5, n_iters=20,
2                        random_state=17, seed = 17)
```

iter	target	baggin...	featur...	max_depth	min_ch...	min
_sp...	num_it...	num_le...				
-----						
1	0.9381	0.8589	0.5245	17.87	11.02	0.
07891	690.7	169.3				
2	0.9383	0.9151	0.1313	20.37	24.19	0.
006944	877.6	202.8				
3	0.9379	0.8102	0.6219	23.28	18.96	0.
04887	354.7	121.7				
4	0.9377	0.9123	0.4168	26.83	16.28	0.
01525	235.8	87.73				
5	0.938	0.9436	0.3339	17.98	22.47	0.
05723	174.1	156.3				
6	0.938	0.9638	0.13	29.77	18.75	0.
05716	993.1	82.91				
7	0.9383	0.8483	0.4234	15.07	21.99	0.
02238	873.9	201.8				
8	0.9377	0.8511	0.7458	18.2	15.4	0.
001702	100.0	220.0				

### 13. Модель на результатах bayesian opt

```

In [32]: 1 %%time
2
3 params = {'learning_rate': 0.02,
4           'bagging_fraction': 0.88,
5           'feature_fraction': 0.53,
6           'min_child_weight': 12.4,
7           'min_split_gain': 0.06,
8           'max_depth': 27,
9           'num_leaves': 218,
10          'metric': 'logloss',
11          'random_state': 17,
12          'n_jobs': -1,
13          #'sample_fraction': 0.33
14          }
15
16 gc.collect()
17 num_iterations = 3000
18 X_train_, X_valid, y_train_, y_valid = train_test_split(X_train, y_train, te
19
20 del X_train, y_train
21
22 train_set = lgb.Dataset(X_train_, y_train_)
23 val_set = lgb.Dataset(X_valid, y_valid)
24
25 del X_train_, X_valid, y_train_, y_valid
26
27 gc.collect()
28
29 model = lgb.train(params, train_set,
30                  num_iterations,
31                  val_set,
32                  verbose_eval=100,
33                  early_stopping_rounds=200,
34                  feval=MacroF1Metric)
35 model.save_model('lgb13.txt', num_iteration=model.best_iteration)

```

Training until validation scores don't improve for 200 rounds

```

[100] valid_0's MacroF1Metric: 0.614335
[200] valid_0's MacroF1Metric: 0.936323
[300] valid_0's MacroF1Metric: 0.937467
[400] valid_0's MacroF1Metric: 0.937577
[500] valid_0's MacroF1Metric: 0.937678
[600] valid_0's MacroF1Metric: 0.937667
[700] valid_0's MacroF1Metric: 0.937722
[800] valid_0's MacroF1Metric: 0.937742
[900] valid_0's MacroF1Metric: 0.937774
[1000] valid_0's MacroF1Metric: 0.937717
Early stopping, best iteration is:
[869] valid_0's MacroF1Metric: 0.937803
Wall time: 24min 52s

```

Out[32]: <lightgbm.basic.Booster at 0x26383974e08>



```
In [37]: 1 %%time
        2 scores, oof, prediction = lgb_cv_loop(X_train, y_train, X_test)
```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

```
score: 0.9378635292772155
score: 0.9380152351255461
score: 0.9375839819218609
score: 0.9393842052037662
score: 0.9384854338633025
```

Wall time: 40min 46s

```
In [38]: 1 np.mean(scores)
```

Out[38]: 0.9382664770783382

```
In [39]: 1 np.save('lgb13_oof_preds.npy', oof)
```

```
In [40]: 1 %%time
        2
        3 train_set = lgb.Dataset(X_train, y_train)
        4
        5 del X_train, y_train
        6
        7 gc.collect()
        8
        9 num_iterations = 900
       10 model = lgb.train(params, train_set,
       11                     num_iterations)
       12
       13 y_lgb_pred = model.predict(X_test)
       14
       15 np.save('lgb13_test_preds.npy', y_lgb_pred)
```

Wall time: 9min 41s

```
In [41]: 1 y_lgb_pred = pred_proc(y_lgb_pred)
        2 sample_df = pd.read_csv("data/sample_submission.csv", dtype={'time':str})
        3 sample_df['open_channels'] = y_lgb_pred
        4 sample_df.to_csv("lgb13.csv", index=False, float_format='%.4f')
```

**Результат:** 0.94 на public lb

