

M5-forecasting, CatBoost.

In [1]:

```
1 !pip install catboost
```

Collecting catboost

Downloading https://files.pythonhosted.org/packages/b1/61/2b8106c8870601671d99ca94d8b8d180f2b740b7cdb95c930147508abcf9/catboost-0.23-cp36-none-manylinux1_x86_64.whl (https://files.pythonhosted.org/packages/b1/61/2b8106c8870601671d99ca94d8b8d180f2b740b7cdb95c930147508abcf9/catboost-0.23-cp36-none-manylinux1_x86_64.whl) (64.7MB)

64.8MB 45kB/s

```
Requirement already satisfied: graphviz in /usr/local/lib/python3.6/dist-packages (from catboost) (0.10.1)
```

Requirement already satisfied: matplotlib in /usr/local/lib/python3.6/dist-packages (from catboost) (3.2.1)

Requirement already satisfied: pandas>=0.24.0 in /usr/local/lib/python3.6/dist-packages (from catboost) (1.0.3)

```
Requirement already satisfied: six in /usr/local/lib/python3.6/dist-packages (from catboost) (1.12.0)
```

Requirement already satisfied: plotly in /usr/local/lib/python3.6/dist-packages (from catboost) (4.4.1)

Requirement already satisfied: numpy>=1.16.0 in /usr/local/lib/python3.6/dist-packages (from catboost) (1.18.3)

Requirement already satisfied: scipy in /usr/local/lib/python3.6/dist-packages (from catboost) (1.4.1)

```
Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.1 in /usr/local/lib/python3.6/dist-packages (from matplotlib>catboost) (2.4.7)
```

Requirement already satisfied: cyciler>=0.10 in /usr/local/lib/python3.6/dist-packages (from matplotlib->catboost) (0.10.0)

```
Requirement already satisfied: python-dateutil>=2.1 in /usr/local/lib/
python3.6/dist-packages (from matplotlib->catboost) (2.8.1)
```

Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.6/dist-packages (from matplotlib->catboost) (1.2.0)

Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.6/dist-packages (from pandas>=0.24.0->catboost) (2018.9)

```
Requirement already satisfied: retrying>=1.3.3 in /usr/local/lib/python3.6/dist-packages (from plotly->catboost) (1.3.3)
```

Installing collected packages: catboost

Successfully installed catboost-0.23

In [0]:

```
1 import gc
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import os
5 import pandas as pd
6
7 from catboost import Pool, CatBoostRegressor
8 from catboost.utils import get_gpu_device_count
9 from datetime import datetime, timedelta
10 from tqdm.notebook import tqdm
```

In [18]:

```
1 print('GPU devices CatBoost:', get_gpu_device_count())
```

GPU devices CatBoost: 1

In [3]:

```
1 from google.colab import drive
2 drive.mount('/gdrive')
```

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aob&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly (https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aob&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly)

Enter your authorization code:

.....

Mounted at /gdrive

In [19]:

```
1 DATA_DIR = '/gdrive/My Drive/M5-forecasting'
2 MODEL_VER = 'v0'
3 BACKWARD_LAGS = 60
4 END_D = 1913
5 CUT_D = END_D - int(365 * 1.2)
6 END_DATE = '2016-04-24'
7 print(datetime.strptime(END_DATE, '%Y-%m-%d'))
8 np.random.seed(0)
```

2016-04-24 00:00:00

Загрузка данных

In [0]:

```
▼ 1 CALENDAR_DTYPES = {
  2     'date': 'str',
  3     'wm_yr_wk': 'int16',
  4     'weekday': 'object',
  5     'wday': 'int16',
  6     'month': 'int16',
  7     'year': 'int16',
  8     'd': 'object',
  9     'event_name_1': 'object',
 10     'event_type_1': 'object',
 11     'event_name_2': 'object',
 12     'event_type_2': 'object',
 13     'snap_CA': 'int16',
 14     'snap_TX': 'int16',
 15     'snap_WI': 'int16'
 16 }
 17 PARSE_DATES = ['date']
▼ 18 SPRICES_DTYPES = {
 19     'store_id': 'object',
 20     'item_id': 'object',
 21     'wm_yr_wk': 'int16',
 22     'sell_price': 'float32'
 23 }
```

In [0]:

```
1 def get_df(is_train:bool=True,
2           backward_lags:int=None):
3     strain = pd.read_csv('{}sales_train_validation.csv'.format(DATA_DIR))
4     print('sales_train_validation.csv:', strain.shape)
5     cat_cols = ['id', 'item_id', 'dept_id', 'store_id', 'cat_id', 'state_id']
6     last_day = int(strain.columns[-1].replace('d_', ''))
7     print('First day:', CUT_D)
8     print('Last day:', last_day)
9     if not is_train:
10         for day in range(last_day + 1, last_day + 28 + 28 + 1):
11             strain['d_{}'.format(day)] = np.nan
12         value_vars = [col for col in strain.columns
13                       if (col.startswith('d_') and (int(col.replace('d_', ''))
14 else:
15     value_vars = [col for col in strain.columns
16                   if (col.startswith('d_') and (int(col.replace('d_', ''))
17 strain = pd.melt(
18     strain,
19     id_vars = cat_cols,
20     value_vars = value_vars,
21     var_name = 'd',
22     value_name = 'sales'
23 )
24 print('Melted sales:', strain.shape)
25 calendar = pd.read_csv('{}calendar.csv'.format(DATA_DIR), dtype=CALENDAR)
26 print('calendar.csv:', calendar.shape)
27 strain = strain.merge(calendar, on='d', copy=False)
28 del calendar
29 gc.collect()
30 print('Merged done')
31 sprices = pd.read_csv('{}sell_prices.csv'.format(DATA_DIR), dtype=SPRICES)
32 print('read prices:', sprices.shape)
33 strain = strain.merge(
34     sprices,
35     on=['store_id', 'item_id', 'wm_yr_wk'],
36     copy=False
37 )
38 del sprices
39 gc.collect()
40
41 if not is_train:
42     strain = strain.loc[
43         strain['date'] >= (datetime.strptime(END_DATE, '%Y-%m-%d') - timedelta(days=1))
44     ]
45 return strain
```

In [0]:

```
1 def make_features(strain):
2     print('make features dataframe:', strain.shape)
3     lags = [7, 28]
4     windows= [7, 28]
5     wnd_feats = ['id', 'item_id']
6     lag_cols = ['lag_{}'.format(lag) for lag in lags ]
7     for lag, lag_col in zip(lags, lag_cols):
8         strain[lag_col] = strain[['id', 'sales']].groupby('id')['sales'].shift(lag)
9     for wnd_feat in wnd_feats:
10         for wnd in windows:
11             for lag_col in lag_cols:
12                 wnd_col = '{}_{}_rmean_{}'.format(lag_col, wnd_feat, wnd)
13                 strain[wnd_col] = strain[[wnd_feat, lag_col]].groupby(wnd_feat).apply(
14                     lambda x: x.rolling(wnd).mean()
15                 )
16     date_features = {
17         'week_num': 'weekofyear',
18         'quarter': 'quarter',
19         'mday': 'day'
20     }
21     for date_feat_name, date_feat_func in date_features.items():
22         strain[date_feat_name] = getattr(strain['date'].dt, date_feat_func).astype(int)
23     strain['d'] = strain['date'].apply(lambda x: int(x.replace('d_', '')))
24     return strain
```

In []:

```
1 %%time
2 strain = get_df(is_train=True, backward_lags=None)
3 strain = make_features(strain)
```

In [25]:

```
1 strain.head(3)
```

Out[25]:

	id	item_id	dept_id	store_id	cat_id	state_id
0	HOBBIES_1_001_CA_1_validation	HOBBIES_1_001	HOBBIES_1	CA_1	HOBBIES	C
1	HOBBIES_1_001_CA_1_validation	HOBBIES_1_001	HOBBIES_1	CA_1	HOBBIES	C
2	HOBBIES_1_001_CA_1_validation	HOBBIES_1_001	HOBBIES_1	CA_1	HOBBIES	C

In [0]:

```
1 drop_cols = ['id', 'sales', 'date', 'wm_yr_wk', 'weekday']
2 train_cols = strain.columns[~strain.columns.isin(drop_cols)]
3 cat_cols = [
4     'item_id', 'dept_id', 'store_id', 'cat_id', 'state_id',
5     'event_name_1', 'event_type_1', 'event_name_2', 'event_type_2'
6 ]
7 strain[cat_cols] = strain[cat_cols].fillna(0)
```

CatBoost

In [10]:

```
1 %%time
2 val_size = int(strain.shape[0] * .15)
3 val_idx = np.random.choice(strain.index.values, val_size, replace=False)
4 train_idx = np.setdiff1d(strain.index.values, val_idx)
5 train_pool = Pool(
6     strain.loc[train_idx][train_cols],
7     strain.loc[train_idx]['sales'],
8     cat_features=cat_cols
9 )
10 val_pool = Pool(
11     strain.loc[val_idx][train_cols],
12     strain.loc[val_idx]['sales'],
13     cat_features=cat_cols
14 )
15 del strain
16 gc.collect()
```

CPU times: user 42.1 s, sys: 1.85 s, total: 44 s

Wall time: 41.1 s

In [11]:

```
1 model = CatBoostRegressor(  
2     iterations=1000,  
3     task_type='GPU',  
4     verbose=200,  
5     loss_function='RMSE',  
6     boosting_type='Plain',  
7     depth=8,  
8     gpu_cat_features_storage='CpuPinnedMemory',  
9     #max_ctr_complexity=2  
10 model.fit(  
11     train_pool,  
12     eval_set = val_pool  
13     #plot=True  
14 )  
15 del train_pool, val_pool  
16 gc.collect()
```

Learning rate set to 0.35204

```
0:      learn: 2.9428610      test: 2.9549983 best: 2.9549983 (0)  
total: 1.32s   remaining: 21m 59s  
200:    learn: 1.9990122      test: 2.0761268 best: 2.0761084 (199)  
total: 3m 36s   remaining: 14m 22s  
400:    learn: 1.9394297      test: 2.0586270 best: 2.0586270 (400)  
total: 7m 10s   remaining: 10m 42s  
600:    learn: 1.9001072      test: 2.0478947 best: 2.0478416 (598)  
total: 10m 50s  remaining: 7m 11s  
800:    learn: 1.8674147      test: 2.0435845 best: 2.0432455 (789)  
total: 14m 27s  remaining: 3m 35s  
999:    learn: 1.8398260      test: 2.0388751 best: 2.0387949 (998)  
total: 18m 2s   remaining: 0us  
bestTest = 2.038794938  
bestIteration = 998  
Shrink model to first 999 iterations.
```

Out[11]:

0

In [0]:

```
1 model.save_model('model_{}.cbm'.format(MODEL_VER))
```

Prediction Loop

In []:

```
1 %%time
2 spread = get_df(is_train=False, backward_lags=BACKWARD_LAGS)
3 for pred_day in tqdm(range(1, 28 + 28 + 1)):
4     pred_date = datetime.strptime(END_DATE, '%Y-%m-%d') + timedelta(days=pred_day)
5     pred_date_back = pred_date - timedelta(days=BACKWARD_LAGS + 1)
6     print('-' * 70)
7     print('forecast day forward:', pred_day, '| forecast date:', pred_date)
8     spread_data = spread[(spread['date'] >= pred_date_back) & (spread['date'] <= pred_date)]
9     spread_data = make_features(spread_data)
10    spread_data = spread_data.loc[spread['date'] == pred_date, train_cols]
11    spread_data[cat_cols] = spread_data[cat_cols].fillna(0)
12    spread.loc[spread['date'] == pred_date, 'sales'] = model.predict(spread_data)
13 del spread_data
14 gc.collect()
```

Submission

In [14]:

```
1 spread_subm = spread.loc[spread['date'] > END_DATE, ['id', 'd', 'sales']].copy()
2 last_d = int(spread.loc[spread['date'] == END_DATE, 'd'].unique()[0].replace('d', ''))
3 print('last d num:', last_d)
4 spread_subm['d'] = spread_subm['d'].apply(lambda x: 'F{}'.format(int(x.replace('d', ''))))
5 spread_subm.loc[spread_subm['sales'] < 0, 'sales'] = 0
```

last d num: 1913

In [0]:

```
1 f_cols = ['F{}'.format(x) for x in range(1, 28 + 28 + 1)]
2 spread_subm = spread_subm.set_index(['id', 'd']).unstack()['sales'][f_cols].reset_index()
3 spread_subm.fillna(0, inplace=True)
4 spread_subm.sort_values('id', inplace=True)
5 spread_subm.reset_index(drop=True, inplace=True)
```

In [16]:

```
1 f_cols_val = ['F{}'.format(x) for x in range(1, 28 + 1)]
2 f_cols_eval = ['F{}'.format(x) for x in range(28 + 1, 28 + 28 + 1)]
3 spread_subm_eval = spread_subm.copy()
4 spread_subm.drop(columns=f_cols_eval, inplace=True)
5 spread_subm_eval.drop(columns=f_cols_val, inplace=True)
6 spread_subm_eval.columns = spread_subm.columns
7 spread_subm_eval['id'] = spread_subm_eval['id'].str.replace('validation', 'eval')
8 spread_subm = pd.concat([spread_subm, spread_subm_eval], axis=0, sort=False)
9 spread_subm.reset_index(drop=True, inplace=True)
10 spread_subm.to_csv('submission.csv', index=False)
11 print('submission saved:', spread_subm.shape)
```

submission saved: (60980, 29)

