

В этом ноутбуке пробуются модель wavenet. Большинство кода взято отсюда <https://www.kaggle.com/mobassir/understanding-ion-switching-with-modeling> (<https://www.kaggle.com/mobassir/understanding-ion-switching-with-modeling>) - просто отчаянная попытка хоть как-то повысить скор за счет ноутбуков на катгле.

```
In [0]: 1 import numpy as np
2 import pandas as pd
3 import os
4 import torch
5 import torch.nn as nn
6 import time
7 import copy
8 from torch.utils.data import Dataset, DataLoader
9 import torch.nn.functional as F
10 from sklearn.metrics import f1_score
11 from sklearn.model_selection import KFold, GroupKFold
12 device = torch.device("cuda:0") if torch.cuda.is_available() else torch.device("cpu")
13 from torch.optim.lr_scheduler import ReduceLROnPlateau
14 import gc
15 import torchcontrib
16 from tqdm import tqdm
17 from torchcontrib.optim import SWA
18 import torchcontrib
```

1. Параметры

```
In [0]: 1 EPOCHS = 80 #150
2 NN_BATCH_SIZE = 16
3 GROUP_BATCH_SIZE = 4000
4 SEED = 321
5 LR = 0.001
6 SPLITS = 5
7
8 outdir = 'wavenet_models'
9 flip = False
10 noise = False
11
12
13 if not os.path.exists(outdir):
14     os.makedirs(outdir)
15
16
17
18 def seed_everything(seed):
19     random.seed(seed)
20     np.random.seed(seed)
21     os.environ['PYTHONHASHSEED'] = str(seed)
22     tf.random.set_seed(seed)
```

2. Данные

```
In [5]: 1 from google.colab import drive
        2 drive.mount('/content/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%3aoob&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%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%3aoob&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%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 /content/gdrive

In [0]:

```
1 def read_data():
2     train = pd.read_csv('/content/gdrive/My Drive/data/train_clean_kalman.csv')
3     test = pd.read_csv('/content/gdrive/My Drive/data/test_clean_kalman.csv')
4     sub = pd.read_csv('/content/gdrive/My Drive/data/sample_submission.csv')
5     return train, test, sub
6
7 def batching(df, batch_size):
8     #print(df)
9     df['group'] = df.groupby(df.index//batch_size, sort=False)['signal'].agg('mean')
10    df['group'] = df['group'].astype(np.uint16)
11    return df
12
13 # normalize the data (standard scaler). We can also try other scalers for a
14 def normalize(train, test):
15     train_input_mean = train.signal.mean()
16     train_input_sigma = train.signal.std()
17     train['signal'] = (train.signal - train_input_mean) / train_input_sigma
18     test['signal'] = (test.signal - train_input_mean) / train_input_sigma
19     return train, test
20
21 # get lead and lags features
22 def lag_with_pct_change(df, windows):
23     for window in windows:
24         df['signal_shift_pos_' + str(window)] = df.groupby('group')['signal'].shift(-window)
25         df['signal_shift_neg_' + str(window)] = df.groupby('group')['signal'].shift(window)
26     return df
27
28 # main module to run feature engineering. Here you may want to try and add o
29 def run_feat_engineering(df, batch_size):
30     # create batches
31     df = batching(df, batch_size = batch_size)
32     # create leads and lags (1, 2, 3 making them 6 features)
33     df = lag_with_pct_change(df, [1, 2, 3])
34     # create signal ** 2 (this is the new feature)
35     df['signal_2'] = df['signal'] ** 2
36     return df
37
38 # fillna with the mean and select features for training
39 def feature_selection(train, test):
40     features = [col for col in train.columns if col not in ['index', 'group']]
41     train = train.replace([np.inf, -np.inf], np.nan)
42     test = test.replace([np.inf, -np.inf], np.nan)
43     for feature in features:
44         feature_mean = pd.concat([train[feature], test[feature]], axis = 0).mean()
45         train[feature] = train[feature].fillna(feature_mean)
46         test[feature] = test[feature].fillna(feature_mean)
47     return train, test, features
48
49
50 def split(GROUP_BATCH_SIZE=4000, SPLITS=5):
51     print('Reading Data Started...')
52     train, test, sample_submission = read_data()
53     train, test = normalize(train, test)
54     print('Reading and Normalizing Data Completed')
55     print('Creating Features')
56     print('Feature Engineering Started...')
```

```

57 train = run_feat_engineering(train, batch_size=GROUP_BATCH_SIZE)
58 test = run_feat_engineering(test, batch_size=GROUP_BATCH_SIZE)
59 train, test, features = feature_selection(train, test)
60 print(train.head())
61 print('Feature Engineering Completed...')
62
63 target = ['open_channels']
64 group = train['group']
65 kf = GroupKFold(n_splits=SPLITS)
66 splits = [x for x in kf.split(train, train[target], group)]
67 new_splits = []
68 for sp in splits:
69     new_split = []
70     new_split.append(np.unique(group[sp[0]]))
71     new_split.append(np.unique(group[sp[1]]))
72     new_split.append(sp[1])
73     new_splits.append(new_split)
74 target_cols = ['open_channels']
75 print(train.head(), train.shape)
76 train_tr = np.array(list(train.groupby('group').apply(lambda x: x[target]
77 train = np.array(list(train.groupby('group').apply(lambda x: x[features]
78 test = np.array(list(test.groupby('group').apply(lambda x: x[features].v
79 print(train.shape, test.shape, train_tr.shape)
80 return train, test, train_tr, new_splits

```

3. Модель

In [0]:

```
1  class Wave_Block(nn.Module):
2
3      def __init__(self,in_channels,out_channels,dilation_rates):
4          super(Wave_Block,self).__init__()
5          self.num_rates = dilation_rates
6          self.convs = nn.ModuleList()
7          self.filter_convs = nn.ModuleList()
8          self.gate_convs = nn.ModuleList()
9
10         self.convs.append(nn.Conv1d(in_channels,out_channels,kernel_size=1))
11         dilation_rates = [2**i for i in range(dilation_rates)]
12         for dilation_rate in dilation_rates:
13             self.filter_convs.append(nn.Conv1d(out_channels,out_channels,kerne
14             self.gate_convs.append(nn.Conv1d(out_channels,out_channels,kerne
15             self.convs.append(nn.Conv1d(out_channels,out_channels,kernel_siz
16
17     def forward(self,x):
18         x = self.convs[0](x)
19         res = x
20         for i in range(self.num_rates):
21             x = F.tanh(self.filter_convs[i](x))*F.sigmoid(self.gate_convs[i]
22             x = self.convs[i+1](x)
23             x += res
24         return x
25
26
27
28
29 class Classifier(nn.Module):
30     def __init__(self):
31         super().__init__()
32         input_size = 128
33
34         self.LSTM = nn.GRU(input_size=input_size,hidden_size=64,num_layers=2
35
36         self.wave_block1 = Wave_Block(8,16,12)
37         self.wave_block2 = Wave_Block(16,32,8)
38         self.wave_block3 = Wave_Block(32,64,4)
39         self.wave_block4 = Wave_Block(64, 128, 1)
40         self.fc = nn.Linear(128, 11)
41
42     def forward(self,x):
43         x = x.permute(0, 2, 1)
44
45         x = self.wave_block1(x)
46         x = self.wave_block2(x)
47         x = self.wave_block3(x)
48
49         #x,_ = self.LSTM(x)
50         x = self.wave_block4(x)
51         x = x.permute(0, 2, 1)
52         x,_ = self.LSTM(x)
53         #x = self.conv1(x)
54         #x = self.attention(x)
55         x = self.fc(x)
56         return x
```

```

57
58
59
60 class EarlyStopping:
61     def __init__(self, patience=5, delta=0, checkpoint_path='checkpoint.pt',
62                 self.patience, self.delta, self.checkpoint_path = patience, delta, c
63                 self.counter, self.best_score = 0, None
64                 self.is_maximize = is_maximize
65
66
67     def load_best_weights(self, model):
68         model.load_state_dict(torch.load(self.checkpoint_path))
69
70     def __call__(self, score, model):
71         if self.best_score is None or \
72            (score > self.best_score + self.delta if self.is_maximize el
73            torch.save(model.state_dict(), self.checkpoint_path)
74            self.best_score, self.counter = score, 0
75            return 1
76         else:
77             self.counter += 1
78             if self.counter >= self.patience:
79                 return 2
80         return 0

```

4. Даталоадеры

```

In [0]: 1 from torch.utils.data import Dataset, DataLoader
2 class IronDataset(Dataset):
3     def __init__(self, data, labels, training=True, transform=None, seq_len=
4         self.data = data
5         self.labels = labels
6         self.transform = transform
7         self.training = training
8         self.flip = flip
9         self.noise_level = noise_level
10        self.class_split = class_split
11        self.seq_len = seq_len
12
13    def __len__(self):
14        return len(self.data)
15
16    def __getitem__(self, idx):
17        if torch.is_tensor(idx):
18            idx = idx.tolist()
19
20        data = self.data[idx]
21        labels = self.labels[idx]
22
23        return [data.astype(np.float32), labels.astype(int)]

```

```
In [9]: 1 train, test, train_tr, new_splits = split()
```

Reading Data Started...

Reading and Normalizing Data Completed

Creating Features

Feature Engineering Started...

	time	signal	...	signal_shift_neg_3	signal_2
0	0.0001	-1.148772	...	-1.298012	1.319677
1	0.0002	-1.184075	...	-1.303999	1.402034
2	0.0003	-1.012891	...	-1.104036	1.025949
3	0.0004	-1.298012	...	-1.123085	1.684836
4	0.0005	-1.303999	...	-1.082287	1.700413

[5 rows x 11 columns]

Feature Engineering Completed...

	time	signal	...	signal_shift_neg_3	signal_2
0	0.0001	-1.148772	...	-1.298012	1.319677
1	0.0002	-1.184075	...	-1.303999	1.402034
2	0.0003	-1.012891	...	-1.104036	1.025949
3	0.0004	-1.298012	...	-1.123085	1.684836
4	0.0005	-1.303999	...	-1.082287	1.700413

[5 rows x 11 columns] (5000000, 11)

(1250, 4000, 8) (500, 4000, 8) (1250, 4000, 1)

5. Обучение

```
In [10]: 1 model = Classifier()  
2 model
```

```
Out[10]: Classifier(  
  (LSTM): GRU(128, 64, num_layers=2, batch_first=True, bidirectional=True)  
  (wave_block1): Wave_Block(  
    (convs): ModuleList(  
      (0): Conv1d(8, 16, kernel_size=(1,), stride=(1,))  
      (1): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (2): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (3): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (4): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (5): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (6): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (7): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (8): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (9): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (10): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (11): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (12): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
    )  
    (filter_convs): ModuleList(  
      (0): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (1): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (2): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (3): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (4): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (5): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (6): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (7): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (8): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (9): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (10): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (11): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
      (12): Conv1d(16, 16, kernel_size=(1,), stride=(1,))  
    )  
  )  
)
```

In [11]:

```
1 %%time
2 test_y = np.zeros([int(2000000/GROUP_BATCH_SIZE), GROUP_BATCH_SIZE, 1])
3 test_dataset = IronDataset(test, test_y, flip=False)
4 test_dataloader = DataLoader(test_dataset, NN_BATCHSIZE, shuffle=False, num_w
5 test_preds_all = np.zeros((2000000, 11))
6
7
8 oof_score = []
9 for index, (train_index, val_index, _) in enumerate(new_splits[0:], start=0)
10     print("Fold : {}".format(index))
11     train_dataset = IronDataset(train[train_index], train_tr[train_index], s
12     train_dataloader = DataLoader(train_dataset, NN_BATCHSIZE, shuffle=True,
13
14     valid_dataset = IronDataset(train[val_index], train_tr[val_index], seq_
15     valid_dataloader = DataLoader(valid_dataset, NN_BATCHSIZE, shuffle=False,
16
17     it = 0
18     model = Classifier()
19     model = model.cuda()
20
21     early_stopping = EarlyStopping(patience=40, is_maximize=True,
22                                   checkpoint_path=os.path.join(outdir, "gru
23
24
25     weight = None#cal_weights()
26     criterion = nn.CrossEntropyLoss(weight=weight)
27     optimizer = torch.optim.Adam(model.parameters(), lr=LR)
28
29     optimizer = torchcontrib.optim.SWA(optimizer, swa_start=10, swa_freq=5,
30
31
32     scheduler = torch.optim.lr_scheduler.ReduceLROnPlateau(optimizer, mode=
33     avg_train_losses, avg_valid_losses = [], [])
34
35
36
37     for epoch in range(EPOCHS):
38         print('*****')
39         print("Folder : {} Epoch : {}".format(index, epoch))
40         print("Curr learning_rate: {:.9f}".format(optimizer.param_groups[0]
41         train_losses, valid_losses = [], []
42         tr_loss_cls_item, val_loss_cls_item = [], []
43
44         model.train() # prep model for training
45         train_preds, train_true = torch.Tensor([]).cuda(), torch.LongTensor(
46
47         for x, y in tqdm(train_dataloader):
48             x = x.cuda()
49             y = y.cuda()
50             #print(x.shape)
51
52
53             optimizer.zero_grad()
54             #loss_fn(model(input), target).backward()
55
56
```



```

57
58     #optimizer.zero_grad()
59     predictions = model(x)
60
61     predictions_ = predictions.view(-1, predictions.shape[-1])
62     y_ = y.view(-1)
63
64     loss = criterion(predictions_, y_)
65
66     # backward pass: compute gradient of the loss with respect to model parameters
67     loss.backward()
68     # perform a single optimization step (parameter update)
69     optimizer.step()
70
71     #scheduler.step()
72     # record training loss
73     train_losses.append(loss.item())
74     train_true = torch.cat([train_true, y_], 0)
75     train_preds = torch.cat([train_preds, predictions_], 0)
76
77     #model.eval() # prep model for evaluation
78     optimizer.swap_swa_sgd()
79     val_preds, val_true = torch.Tensor([]).cuda(), torch.LongTensor([]).cuda()
80     print('EVALUATION')
81     with torch.no_grad():
82         for x, y in tqdm(valid_dataloader):
83             x = x.cuda()#.to(device)
84             y = y.cuda()#.to(device)
85
86             predictions = model(x)
87             predictions_ = predictions.view(-1, predictions.shape[-1])
88             y_ = y.view(-1)
89
90             loss = criterion(predictions_, y_)
91
92             valid_losses.append(loss.item())
93
94
95             val_true = torch.cat([val_true, y_], 0)
96             val_preds = torch.cat([val_preds, predictions_], 0)
97
98     # calculate average loss over an epoch
99     train_loss = np.average(train_losses)
100    valid_loss = np.average(valid_losses)
101    avg_train_losses.append(train_loss)
102    avg_valid_losses.append(valid_loss)
103    print("train_loss: {:.0.6f}, valid_loss: {:.0.6f}".format(train_loss, valid_loss))
104
105    train_score = f1_score(train_true.cpu().detach().numpy(), train_preds.cpu().detach().numpy(),
106                          labels=list(range(11)), average='macro')
107
108    val_score = f1_score(val_true.cpu().detach().numpy(), val_preds.cpu().detach().numpy(),
109                        labels=list(range(11)), average='macro')
110
111    scheduler.step(val_score)
112    print("train_f1: {:.0.6f}, valid_f1: {:.0.6f}".format(train_score, val_score))
113    res = early_stopping(val_score, model)

```

```

114         #print('fres:', res)
115         if res == 2:
116             print("Early Stopping")
117             print('folder %d global best val max f1 model score %f' % (index,
118                             best_val_max_f1_model_score))
119             break
120         elif res == 1:
121             print('save folder %d global val max f1 model score %f' % (index,
122                             best_val_max_f1_model_score))
123             print('Folder {} finally best global max f1 score is {}'.format(index,
124                                     best_val_max_f1_model_score))
125             oof_score.append(round(early_stopping.best_score, 6))
126
127         model.eval()
128         pred_list = []
129         with torch.no_grad():
130             for x, y in tqdm(test_dataloader):
131
132                 x = x.cuda()
133                 y = y.cuda()
134
135                 predictions = model(x)
136                 predictions_ = predictions.view(-1, predictions.shape[-1]) # shape [2000000, 11]
137                 #print(predictions.shape, F.softmax(predictions_, dim=1).cpu().numpy())
138                 pred_list.append(F.softmax(predictions_, dim=1).cpu().numpy())
139                 #a = input()
140             test_preds = np.vstack(pred_list) # shape [2000000, 11]
141             test_preds_all += test_preds

```

curr_learning_rate: 0.000100000

```

100%|██████████| 63/63 [00:37<00:00, 1.67it/s]
 0%|          | 0/16 [00:00<?, ?it/s]

```

EVALUATION

```

100%|██████████| 16/16 [00:04<00:00, 3.97it/s]

```

train_loss: 0.079804, valid_loss: 0.087705

```

 0%|          | 0/32 [00:00<?, ?it/s]

```

train_f1: 0.938050, valid_f1: 0.937883

Folder 4 finally best global max f1 score is 0.9384072723611167

```

100%|██████████| 32/32 [00:08<00:00, 3.92it/s]

```

CPU times: user 4h 9min 29s, sys: 30min 1s, total: 4h 39min 30s

Wall time: 4h 41min 49s

```
In [12]: 1 print('all folder score is:%s'%str(oof_score))
2 print('OOF mean score is: %f'% (sum(oof_score)/len(oof_score)))
3 print('Generate submission.....')
4 submission_csv_path = '/content/gdrive/My Drive/data/sample_submission.csv'
5 ss = pd.read_csv(submission_csv_path, dtype={'time': str})
6 test_preds_all = test_preds_all / np.sum(test_preds_all, axis=1)[:, None]
7 test_pred_frame = pd.DataFrame({'time': ss['time'].astype(str),
8                                'open_channels': np.argmax(test_preds_all, a
9 test_pred_frame.to_csv("/content/gdrive/My Drive/data/gru_submission.csv", i
10 print('over')
```

all folder score is:[0.938465, 0.939936, 0.93835, 0.93827, 0.938407]

OOF mean score is: 0.938686

Generate submission.....

over

Результат: 0.941 на public lb