

**В этом ноутбуке** пробуются модель на основе LSTM (так как у нас тут что-то похожее на временные ряды, то есть гипотеза, что LSTM хорошо здесь зайдет), обученная на предсказаниях леса. Ноутбук вдохновлен вот этим <https://www.kaggle.com/khalildmk/simple-two-layer-bidirectional-lstm-with-pytorch> (<https://www.kaggle.com/khalildmk/simple-two-layer-bidirectional-lstm-with-pytorch>) - оттуда взято большинство кода, но была переделана работа с данными и размерности в самой модели.

```
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
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
```

## 1. Параметры и модель

```
In [0]: 1 n_epochs = 100
2 lr = 0.01
3 n_folds = 5
4 lstm_input_size = 12
5 hidden_state_size = 30
6 batch_size = 30
7 num_sequence_layers = 2
8 output_dim = 11
9 num_time_steps = 4000
10 rnn_type = 'LSTM'
```

```

In [0]: 1 class Bi_RNN(nn.Module):
2
3     def __init__(self, input_dim, hidden_dim, batch_size, output_dim=11, num
4         super(Bi_RNN, self).__init__()
5         self.input_dim = input_dim
6         self.hidden_dim = hidden_dim
7         self.batch_size = batch_size
8         self.num_layers = num_layers
9
10        #Define the initial linear hidden layer
11        self.init_linear = nn.Linear(self.input_dim, self.input_dim)
12
13        # Define the LSTM layer
14        self.lstm = eval('nn.' + rnn_type)(self.input_dim, self.hidden_dim,
15
16        # Define the output layer
17        self.linear = nn.Linear(self.hidden_dim * 2, output_dim)
18
19    def init_hidden(self):
20        # This is what we'll initialise our hidden state as
21        return (torch.zeros(self.num_layers, self.batch_size, self.hidden_dim),
22                torch.zeros(self.num_layers, self.batch_size, self.hidden_dim))
23
24    def forward(self, input):
25        #Forward pass through initial hidden layer
26        linear_input = self.init_linear(input)
27
28        # Forward pass through LSTM layer
29        # shape of lstm_out: [batch_size, input_size, hidden_dim]
30        # shape of self.hidden: (a, b), where a and b both
31        # have shape (batch_size, num_layers, hidden_dim).
32        lstm_out, self.hidden = self.lstm(linear_input)
33
34        # Can pass on the entirety of lstm_out to the next layer if it is a
35        y_pred = self.linear(lstm_out)
36        return y_pred

```

## 2. Даталоадеры

```

In [0]: 1 class ION_Dataset_Sequential(Dataset):
2         def __init__(self, input, output):
3             self.input = input
4             self.output = output
5
6         def __len__(self):
7             return len(self.input)
8
9         def __getitem__(self, idx):
10            x = self.input[idx]
11            y = self.output[idx]
12            x = torch.tensor(x, dtype=torch.float)
13            y = torch.tensor(y, dtype=torch.float)
14            return x, y
15
16 class ION_Dataset_Sequential_test(Dataset):
17     def __init__(self, input):
18         self.input = input
19
20     def __len__(self):
21         return len(self.input)
22
23     def __getitem__(self, idx):
24         x = self.input[idx]
25         x = torch.tensor(x, dtype=torch.float)
26         return x

```

### 3. Данные

```

In [5]: 1 from google.colab import drive
2         drive.mount('/content/gdrive')

```

Drive already mounted at /content/gdrive; to attempt to forcibly remount, call drive.mount("/content/gdrive", force\_remount=True).

```

In [0]: 1 train_df = pd.read_csv('/content/gdrive/My Drive/data/train_clean.csv')
2         test_df = pd.read_csv('/content/gdrive/My Drive/data/test_clean.csv')

```

```

In [0]: 1 train_probs = np.load('/content/gdrive/My Drive/data/Y_train_proba.npy')
2         test_probs = np.load('/content/gdrive/My Drive/data/Y_test_proba.npy')

```

```

In [0]: 1 train_df = pd.concat([train_df, pd.DataFrame(train_probs)], axis=1)
2         test_df = pd.concat([test_df, pd.DataFrame(test_probs)], axis=1)

```

```

In [0]: 1 train_df.columns = train_df.columns.astype(str)
2         test_df.columns = test_df.columns.astype(str)

```

```
In [0]: 1 X = train_df[['signal', '0', '1', '2', '3', '4', '5', '6', '7', '8', '9', '10', '11', '12', '13', '14', '15', '16', '17', '18', '19', '20', '21', '22', '23', '24', '25', '26', '27', '28', '29', '30', '31', '32', '33', '34', '35', '36', '37', '38', '39', '40', '41', '42', '43', '44', '45', '46', '47', '48', '49', '50', '51', '52', '53', '54', '55', '56', '57', '58', '59', '60', '61', '62', '63', '64', '65', '66', '67', '68', '69', '70', '71', '72', '73', '74', '75', '76', '77', '78', '79', '80', '81', '82', '83', '84', '85', '86', '87', '88', '89', '90', '91', '92', '93', '94', '95', '96', '97', '98', '99']]
2 y = pd.get_dummies(train_df['open_channels']).values.reshape(-1, num_time_steps)
3 test_input = test_df[['signal', '0', '1', '2', '3', '4', '5', '6', '7', '8', '9', '10', '11', '12', '13', '14', '15', '16', '17', '18', '19', '20', '21', '22', '23', '24', '25', '26', '27', '28', '29', '30', '31', '32', '33', '34', '35', '36', '37', '38', '39', '40', '41', '42', '43', '44', '45', '46', '47', '48', '49', '50', '51', '52', '53', '54', '55', '56', '57', '58', '59', '60', '61', '62', '63', '64', '65', '66', '67', '68', '69', '70', '71', '72', '73', '74', '75', '76', '77', '78', '79', '80', '81', '82', '83', '84', '85', '86', '87', '88', '89', '90', '91', '92', '93', '94', '95', '96', '97', '98', '99']]
4 train_input_mean = X.mean()
5 train_input_sigma = X.std()
6 test_input = (test_input - train_input_mean) / train_input_sigma
7 test_preds = np.zeros((int(test_input.shape[0] * test_input.shape[1])))
8 test = ION_Dataset_Sequential_test(test_input)
9 test_loader = DataLoader(test, batch_size=batch_size, shuffle=False)
```

◀ [REDACTED] ▶

#### 4. Обучение с разбиением на фолды

```

In [11]: 1 kfold = KFold(n_splits=n_folds, shuffle=True, random_state=42)
2 local_val_score = 0
3 models = {}
4
5 k=0 #initialize fold number
6 for tr_idx, val_idx in kfold.split(X, y):
7     test_p = np.zeros((int(test_input.shape[0] * test_input.shape[1])))
8
9     print('starting fold', k)
10    k += 1
11
12    print(6*'#', 'splitting and reshaping the data')
13    train_input = X[tr_idx]
14    print(train_input.shape)
15    train_target = y[tr_idx]
16    val_input = X[val_idx]
17    val_target = y[val_idx]
18    train_input_mean = train_input.mean()
19    train_input_sigma = train_input.std()
20    val_input = (val_input-train_input_mean)/train_input_sigma
21    train_input = (train_input-train_input_mean)/train_input_sigma
22
23    print(6*'#', 'Loading')
24    train = ION_Dataset_Sequential(train_input, train_target)
25    valid = ION_Dataset_Sequential(val_input, val_target)
26    train_loader = torch.utils.data.DataLoader(train, batch_size=batch_size,
27    valid_loader = torch.utils.data.DataLoader(valid, batch_size=batch_size,
28
29    #Build tensor data for torch
30    train_preds = np.zeros((int(train_input.shape[0] * train_input.shape[1]))
31    val_preds = np.zeros((int(val_input.shape[0] * val_input.shape[1])))
32    best_val_preds = np.zeros((int(val_input.shape[0] * val_input.shape[1]))
33    train_targets = np.zeros((int(train_input.shape[0] * train_input.shape[1]))
34    avg_losses_f = []
35    avg_val_losses_f = []
36
37    #Define Loss function
38    loss_fn = torch.nn.BCEWithLogitsLoss()
39
40    #Build model, initialize weights and define optimizer
41    model = Bi_RNN(lstm_input_size, hidden_state_size, batch_size=batch_size)
42    model = model.to(device)
43    optimizer = torch.optim.Adam(model.parameters(), lr=lr, weight_decay=1e-5)
44    scheduler = ReduceLROnPlateau(optimizer, 'min', patience=150, factor=0.1)
45    temp_val_loss = 999999999
46    reached_val_score = 0
47
48    #Iterate through epochs
49    for epoch in range(n_epochs):
50        start_time = time.time()
51
52        #Train
53        model.train()
54        avg_loss = 0.
55        for i, (x_batch, y_batch) in enumerate(train_loader):
56            x_batch = x_batch.view(-1, num_time_steps, lstm_input_size)

```

```

57         y_batch = y_batch.view(-1, num_time_steps, output_dim)
58         optimizer.zero_grad()
59         y_pred = model(x_batch.cuda())
60         loss = loss_fn(y_pred.cpu(), y_batch)
61         loss.backward()
62         optimizer.step()
63         avg_loss += loss.item() / len(train_loader)
64         pred = F.softmax(y_pred, 2).detach().cpu().numpy().argmax(axis=-1)
65         train_preds[i * batch_size * train_input.shape[1]:(i + 1) * batch_size * train_input.shape[1]] = pred
66         train_targets[i * batch_size * train_input.shape[1]:(i + 1) * batch_size * train_input.shape[1]] = y_batch.cpu().numpy()
67         del y_pred, loss, x_batch, y_batch, pred
68
69     #Evaluate
70     model.eval()
71     avg_val_loss = 0.
72     for i, (x_batch, y_batch) in enumerate(valid_loader):
73         x_batch = x_batch.view(-1, num_time_steps, lstm_input_size)
74         y_batch = y_batch.view(-1, num_time_steps, output_dim)
75         y_pred = model(x_batch.cuda()).detach()
76         avg_val_loss += loss_fn(y_pred.cpu(), y_batch).item() / len(valid_loader)
77         pred = F.softmax(y_pred, 2).detach().cpu().numpy().argmax(axis=-1)
78         val_preds[i * batch_size * val_input.shape[1]:(i + 1) * batch_size * val_input.shape[1]] = pred
79         del y_pred, x_batch, y_batch, pred
80     if avg_val_loss < temp_val_loss:
81         temp_val_loss = avg_val_loss
82
83     #Calculate F1-score
84     train_score = f1_score(train_targets, train_preds, average='macro')
85     val_score = f1_score(val_target.argmax(axis=2).reshape((-1)), val_preds, average='macro')
86
87     #Print output of epoch
88     elapsed_time = time.time() - start_time
89     scheduler.step(avg_val_loss)
90     if epoch%10 == 0:
91         print('Epoch {}/{} \t loss={:.4f} \t train_f1={:.4f} \t val_loss={:.4f} \t val_f1={:.4f}'.format(epoch, n_epochs, avg_loss, train_score, temp_val_loss, val_score))
92
93     if val_score > reached_val_score:
94         reached_val_score = val_score
95         best_model = copy.deepcopy(model.state_dict())
96         best_val_preds = copy.deepcopy(val_preds)
97
98     #Calculate F1-score of the fold
99     val_score_fold = f1_score(val_target.argmax(axis=2).reshape((-1)), best_val_preds, average='macro')
100
101     #Save the fold's model in a dictionary
102     models[k] = best_model
103
104     #Print F1-score of the fold
105     print("BEST VALIDATION SCORE (F1): ", val_score_fold)
106     local_val_score += (1/n_folds) * val_score_fold
107
108     #Print final average k-fold CV F1-score
109     print("Final Score ", local_val_score)

```

```

starting fold 0
##### splitting and reshaping the data

```

(1000, 4000, 12)

##### Loading

Epoch 1/100	loss=0.2928	train_f1=0.1054	val_loss=0.2055
val_f1=0.2158	time=41.63s		
Epoch 11/100	loss=0.0204	train_f1=0.8489	val_loss=0.0238
val_f1=0.8492	time=42.01s		
Epoch 21/100	loss=0.0165	train_f1=0.9382	val_loss=0.0193
val_f1=0.9373	time=41.90s		
Epoch 31/100	loss=0.0164	train_f1=0.9377	val_loss=0.0188
val_f1=0.9372	time=41.97s		
Epoch 41/100	loss=0.0156	train_f1=0.9384	val_loss=0.0185
val_f1=0.9372	time=42.19s		
Epoch 51/100	loss=0.0157	train_f1=0.9382	val_loss=0.0185
val_f1=0.9374	time=41.84s		
Epoch 61/100	loss=0.0162	train_f1=0.9379	val_loss=0.0190
val_f1=0.9371	time=42.07s		
Epoch 71/100	loss=0.0157	train_f1=0.9383	val_loss=0.0186
val_f1=0.9375	time=41.90s		
Epoch 81/100	loss=0.0156	train_f1=0.9383	val_loss=0.0184
val_f1=0.9376	time=42.10s		
Epoch 91/100	loss=0.0152	train_f1=0.9383	val_loss=0.0184
val_f1=0.9375	time=41.94s		

BEST VALIDATION SCORE (F1): 0.9377366523157712

starting fold 1

##### splitting and reshaping the data

(1000, 4000, 12)

##### Loading

Epoch 1/100	loss=0.2628	train_f1=0.1831	val_loss=0.1796
val_f1=0.2677	time=41.80s		
Epoch 11/100	loss=0.0209	train_f1=0.8492	val_loss=0.0259
val_f1=0.8491	time=41.80s		
Epoch 21/100	loss=0.0168	train_f1=0.9378	val_loss=0.0208
val_f1=0.9368	time=41.91s		
Epoch 31/100	loss=0.0156	train_f1=0.9381	val_loss=0.0198
val_f1=0.9371	time=41.90s		
Epoch 41/100	loss=0.0159	train_f1=0.9379	val_loss=0.0197
val_f1=0.9371	time=41.70s		
Epoch 51/100	loss=0.0154	train_f1=0.9382	val_loss=0.0199
val_f1=0.9359	time=41.60s		
Epoch 61/100	loss=0.0151	train_f1=0.9383	val_loss=0.0196
val_f1=0.9368	time=41.82s		
Epoch 71/100	loss=0.0165	train_f1=0.9370	val_loss=0.0200
val_f1=0.9369	time=41.43s		
Epoch 81/100	loss=0.0152	train_f1=0.9381	val_loss=0.0193
val_f1=0.9372	time=41.68s		
Epoch 91/100	loss=0.0153	train_f1=0.9382	val_loss=0.0194
val_f1=0.9371	time=41.30s		

BEST VALIDATION SCORE (F1): 0.9376231785584731

starting fold 2

##### splitting and reshaping the data

(1000, 4000, 12)

##### Loading

Epoch 1/100	loss=0.2652	train_f1=0.1949	val_loss=0.1760
val_f1=0.2983	time=41.88s		
Epoch 11/100	loss=0.0225	train_f1=0.8494	val_loss=0.0252
val_f1=0.8475	time=41.57s		
Epoch 21/100	loss=0.0170	train_f1=0.9378	val_loss=0.0192

val_f1=0.9370	time=40.97s		
Epoch 31/100	loss=0.0167	train_f1=0.9367	val_loss=0.0193
val_f1=0.9366	time=40.41s		
Epoch 41/100	loss=0.0162	train_f1=0.9379	val_loss=0.0182
val_f1=0.9375	time=40.79s		
Epoch 51/100	loss=0.0156	train_f1=0.9381	val_loss=0.0181
val_f1=0.9375	time=40.29s		
Epoch 61/100	loss=0.0157	train_f1=0.9381	val_loss=0.0179
val_f1=0.9377	time=40.46s		
Epoch 71/100	loss=0.0156	train_f1=0.9379	val_loss=0.0178
val_f1=0.9377	time=40.37s		
Epoch 81/100	loss=0.0158	train_f1=0.9374	val_loss=0.0182
val_f1=0.9369	time=40.44s		
Epoch 91/100	loss=0.0158	train_f1=0.9380	val_loss=0.0185
val_f1=0.9368	time=40.86s		
BEST VALIDATION SCORE (F1): 0.9380991640612855			
starting fold 3			
##### splitting and reshaping the data			
(1000, 4000, 12)			
##### Loading			
Epoch 1/100	loss=0.2821	train_f1=0.1583	val_loss=0.1878
val_f1=0.2488	time=40.29s		
Epoch 11/100	loss=0.0225	train_f1=0.8489	val_loss=0.0228
val_f1=0.8505	time=40.34s		
Epoch 21/100	loss=0.0174	train_f1=0.9378	val_loss=0.0177
val_f1=0.9395	time=40.54s		
Epoch 31/100	loss=0.0162	train_f1=0.9375	val_loss=0.0172
val_f1=0.9388	time=40.92s		
Epoch 41/100	loss=0.0161	train_f1=0.9375	val_loss=0.0175
val_f1=0.9390	time=40.51s		
Epoch 51/100	loss=0.0159	train_f1=0.9378	val_loss=0.0166
val_f1=0.9396	time=40.60s		
Epoch 61/100	loss=0.0160	train_f1=0.9378	val_loss=0.0165
val_f1=0.9396	time=40.51s		
Epoch 71/100	loss=0.0160	train_f1=0.9375	val_loss=0.0164
val_f1=0.9394	time=40.49s		
Epoch 81/100	loss=0.0157	train_f1=0.9378	val_loss=0.0165
val_f1=0.9396	time=40.96s		
Epoch 91/100	loss=0.0159	train_f1=0.9379	val_loss=0.0165
val_f1=0.9398	time=40.47s		
BEST VALIDATION SCORE (F1): 0.93980716248632			
starting fold 4			
##### splitting and reshaping the data			
(1000, 4000, 12)			
##### Loading			
Epoch 1/100	loss=0.2710	train_f1=0.1693	val_loss=0.1904
val_f1=0.2708	time=40.36s		
Epoch 11/100	loss=0.0208	train_f1=0.8489	val_loss=0.0254
val_f1=0.8496	time=40.46s		
Epoch 21/100	loss=0.0159	train_f1=0.9380	val_loss=0.0203
val_f1=0.9375	time=41.19s		
Epoch 31/100	loss=0.0155	train_f1=0.9381	val_loss=0.0199
val_f1=0.9379	time=40.58s		
Epoch 41/100	loss=0.0153	train_f1=0.9383	val_loss=0.0196
val_f1=0.9378	time=40.38s		
Epoch 51/100	loss=0.0154	train_f1=0.9379	val_loss=0.0194
val_f1=0.9380	time=40.27s		



```

Epoch 61/100      loss=0.0152      train_f1=0.9383      val_loss=0.0195
val_f1=0.9382      time=40.18s
Epoch 71/100      loss=0.0151      train_f1=0.9379      val_loss=0.0199
val_f1=0.9377      time=40.22s
Epoch 81/100      loss=0.0153      train_f1=0.9379      val_loss=0.0194
val_f1=0.9382      time=40.19s
Epoch 91/100      loss=0.0151      train_f1=0.9378      val_loss=0.0196
val_f1=0.9375      time=40.18s
BEST VALIDATION SCORE (F1): 0.9384948498276348
Final Score 0.938352201449897

```

## 5. Предсказание - усреднение предсказаний моделей, обученных на разных фолдах

```

In [0]: 1 for k in range(n_folds):
2         test_p = np.zeros((int(test_input.shape[0] * test_input.shape[1])))
3         k += 1
4
5         #Import model of fold k
6         model = Bi_RNN(lstm_input_size, hidden_state_size, batch_size=batch_size)
7         model = model.to(device)
8         model.load_state_dict(models[k])
9
10        #Make predictions on test data
11        model.eval()
12        for i, x_batch in enumerate(test_loader):
13            x_batch = x_batch.view(-1, num_time_steps, lstm_input_size)
14            y_pred = model(x_batch.cuda()).detach()
15            pred = F.softmax(y_pred, 2).detach().cpu().numpy().argmax(axis=-1)
16            test_p[i * batch_size * test_input.shape[1]:(i + 1) * batch_size * t
17                del y_pred, x_batch, pred
18            test_preds += (1/n_folds) * test_p

```

## 6. Сохраняем ответы в файл

```

In [0]: 1 df_sub = pd.read_csv("/content/gdrive/My Drive/data/sample_submission.csv",
2         df_sub.open_channels = np.array(test_preds, np.int)
3         df_sub.to_csv("/content/gdrive/My Drive/data/submission_bilstm.csv", index=F

```

**Результат:** 0.939 на public lb