

Алгоритм Витерби без обучения на отдельных батчах

В данном ноутбуке рассмотрена реализация алгоритма Витерби без обучения на отдельных схожих сегментах данных.

В данном случае это будет нам полезно возможностью использования данного алгоритма при стекинге с бустингами и лесами.

В данном ноутбуке не так много оригинального кода. В основном использованы результаты отсюда: <https://www.kaggle.com/friedchips/the-viterbi-algorithm-a-complete-solution> (<https://www.kaggle.com/friedchips/the-viterbi-algorithm-a-complete-solution>)

И отсюда: <https://www.kaggle.com/miklgr500/viterbi-algorithm-without-segmentation-on-groups> (<https://www.kaggle.com/miklgr500/viterbi-algorithm-without-segmentation-on-groups>)

В данном случае пока достаточно понимания принципа работы данного алгоритма. Реализовывать его самостоятельно пока не имеет смысла в связи с тем, что есть ненулевая вероятность того, что данный алгоритм ничего не даст при стекинге, т.к. его качество (~ 0.93 - 0.934) недотягивает до качества стекингов, лесов и нейронных сетей, и может быть попросту незамечено метамоделью стекинга.

In [2]:

```
1 import os
2
3 import numpy as np
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 from tqdm.notebook import tqdm
9
10 from sklearn.metrics import f1_score, accuracy_score, confusion_matrix
11 from sklearn.model_selection import StratifiedKFold, KFold
12
13
14 plt.style.use('dark_background')
```

Попробуем сначала использовать предобработанный датасет с более интеллектуальной очисткой данных, рассмотренный в вышеперечисленных публичных ноутбуках.

In [7]:

```
1 train = pd.read_csv('clean-kalman/train_clean_kalman.csv')
2 test  = pd.read_csv('clean-kalman/test_clean_kalman.csv')
```

```

In [8]: 1 class ViterbiClassifier:
2     def __init__(self, num_bins=1000):
3         self._n_bins = num_bins
4         self._p_trans = None
5         self._p_signal = None
6         self._signal_bins = None
7         self._p_in = None
8
9     def fit(self, x, y):
10        self._p_trans = self.markov_p_trans(y)
11        self._p_signal, self._signal_bins = self.markov_p_signal(true_state, x, self._n_bins)
12
13        self._p_in = np.ones(len(self._p_trans)) / len(self._p_trans)
14        return self
15
16    def predict(self, x):
17        x_dig = self.digitize_signal(x, self._signal_bins)
18        return self.viterbi(self._p_trans, self._p_signal, self._p_in, x_dig)
19
20    @classmethod
21    def digitize_signal(cls, signal, signal_bins):
22        # https://www.kaggle.com/friedchips/the-viterbi-algorithm-a-complete-solution
23        signal_dig = np.digitize(signal, bins=signal_bins) - 1
24        signal_dig = np.minimum(signal_dig, len(signal_bins) - 2)
25        return signal_dig
26
27    @classmethod
28    def markov_p_signal(cls, state, signal, num_bins = 1000):
29        # https://www.kaggle.com/friedchips/the-viterbi-algorithm-a-complete-solution
30        states_range = np.arange(state.min(), state.max() + 1)
31        signal_bins = np.linspace(signal.min(),
32                                   signal.max(),
33                                   num_bins + 1)
34        p_signal = np.array([ np.histogram(signal[state == s],
35                                           bins=signal_bins)[0] for s in states_range ])
36        p_signal = np.array([ p / np.sum(p) if np.sum(p) != 0 else p for p in p_signal ])
37        return p_signal, signal_bins
38
39    @classmethod
40    def markov_p_trans(cls, states):
41        # https://www.kaggle.com/friedchips/the-viterbi-algorithm-a-complete-solution

```

```

42     max_state = np.max(states)
43     states_next = np.roll(states, -1)
44     matrix = []
45     for i in tqdm(range(max_state + 1)):
46         current_row = np.histogram(states_next[states == i],
47                                     bins=np.arange(max_state + 2))[0]
48         if np.sum(current_row) == 0:
49             current_row = np.ones(max_state + 1) / (max_state + 1)
50         else:
51             current_row = current_row / np.sum(current_row)
52         matrix.append(current_row)
53     return np.array(matrix)
54
55     @classmethod
56     def viterbi(cls, p_trans, p_signal, p_in, signal):
57         # https://www.kaggle.com/friedchips/the-viterbi-algorithm-a-complete-solution
58         offset = 10**(-20)
59
60         p_trans_tlog = np.transpose(np.log2(p_trans + offset))
61         p_signal_tlog = np.transpose(np.log2(p_signal + offset))
62         p_in_log = np.log2(p_in + offset)
63
64         p_state_log = [ p_in_log + p_signal_tlog[signal[0]] ]
65
66         for s in tqdm(signal[1:]):
67             p_state_log.append(np.max(p_state_log[-1] + p_trans_tlog, axis=1)
68                                 + p_signal_tlog[s])
69
70         states = np.argmax(p_state_log, axis=1)
71
72     return states

```

```

In [9]: 1 true_state = train.open_channels.values
        2 signal = train.signal.values

```

Обучение

In [10]:

```
1 viterbi = ViterbiClassifier().fit(signal, true_state)
2 train_prediction = viterbi.predict(signal)
```

```
HBox(children=(FloatProgress(value=0.0, max=11.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=4818797.0), HTML(value='')))
```

In [11]:

```
1 print("Accuracy =", accuracy_score(y_pred=train_prediction, y_true=true_state))
2 print("F1 macro =", f1_score(y_pred=train_prediction, y_true=true_state, average='macro'))
```

```
Accuracy = 0.9641219656852186
```

```
F1 macro = 0.9320404278331105
```

In [12]:

```
1 X_train = train.signal
2 y_train = train.open_channels
3
4 X_test = test.signal
5
6 n_fold = 5
7 folds = KFold(n_splits=n_fold, shuffle=True, random_state=17)
8
9 oof = np.zeros(len(X_train))
10 prediction = np.zeros(len(X_test))
11 scores = []
12
13 for training_index, validation_index in tqdm(folds.split(X_train), total=n_fold):
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     true_state = y_train_.values
21     signal = X_train_.values
22
23     model = ViterbiClassifier().fit(signal, true_state)
24
25     # скор на валидации
26     preds = model.predict(X_valid.values)
27     oof[validation_index] = preds.reshape(-1,)
28
29     preds = np.round(np.clip(preds, 0, 10)).astype(int)
30     score = f1_score(y_valid, preds, average = 'macro')
31     scores.append(score)
32
33     # предсказание на тесте
34     preds = model.predict(X_test)
35     prediction += preds
36
37     print(f'score: {score}')
38
39 prediction /= n_fold
```

HBox(children=(FloatProgress(value=0.0, max=5.0), HTML(value='')))

```
HBox(children=(FloatProgress(value=0.0, max=11.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=963759.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=1999999.0), HTML(value='')))
```

```
score: 0.9293983141212808
```

```
HBox(children=(FloatProgress(value=0.0, max=11.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=963759.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=1999999.0), HTML(value='')))
```

```
score: 0.9294417188919251
```

```
HBox(children=(FloatProgress(value=0.0, max=11.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=963759.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=1999999.0), HTML(value='')))
```

```
score: 0.9298128243081841
```

```
HBox(children=(FloatProgress(value=0.0, max=11.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=963758.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=1999999.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=1999999.0), HTML(value='')))
```

```
score: 0.9299081397497022
```

```
HBox(children=(FloatProgress(value=0.0, max=11.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=963758.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=1999999.0), HTML(value='')))
```

```
score: 0.930220770278687
```

Алгоритм показал сравнительно высокий скор на валидации. Учитывая, что никаких дополнительных фичей здесь не используется, только непосредственное приближение восстановления скрытых состояний по видимым, это очень хороший результат.

Сохраним результат работы и oof предсказания для стекинга.

```
In [69]: 1 def pred_proc(pred):  
2         pred = np.round(np.clip(pred, 0, 10))  
3         return pred.astype(int)
```

```
In [70]: 1 y_catboost_pred = pred_proc(prediction)  
2  
3 sample_df = pd.read_csv("data/sample_submission.csv", dtype={'time':str})  
4 sample_df['open_channels'] = y_catboost_pred  
5 sample_df.to_csv("viterbi_best.csv", index=False, float_format='%0.4f')
```

```
In [71]: 1 np.save('preds_viterbi_best', prediction)  
2         np.save('oof_viterbi_best', oof)
```

Посмотрим теперь, какое качество алгоритм даёт на наших данных.

```
In [15]: 1 train = pd.read_csv('data-without-drift/train_clean.csv')
2 test  = pd.read_csv('data-without-drift/test_clean.csv')
3
4 true_state = train.open_channels.values
5 signal = train.signal.values
6
7 viterbi = ViterbiClassifier().fit(signal, true_state)
8 train_prediction = viterbi.predict(signal)
```

```
HBox(children=(FloatProgress(value=0.0, max=11.0), HTML(value='')))
```

```
HBox(children=(FloatProgress(value=0.0, max=4999999.0), HTML(value='')))
```

```
In [16]: 1 print("Accuracy =", accuracy_score(y_pred=train_prediction, y_true=true_state))
2 print("F1 macro =", f1_score(y_pred=train_prediction, y_true=true_state, average='macro'))
```

```
Accuracy = 0.9117018
```

```
F1 macro = 0.8518252491783009
```

Качество даже на обучающей выборке совсем далеко от оптимального, поэтому будем использовать предыдущий датасет.

```
In [ ]: 1
```