

## Алгоритм Витерби для скрытого марковского процесса.

Не секрет, что данные в соревновании получены искусственно.

Утверждается, они были сгенерированы с помощью некоего марковского процесса, затем к ним был добавлен шум (сгенерированный, потом прогнанный через физический усилитель и записанный на микрофон), затем сдвиг (drift).

Следовательно, возникает задача получения скрытых состояний(количества открытых ионных каналов) цепи по видимым(сигналу).

В данном ноутбуке проведён эксперимент с использованием алгоритма Витерби поиска наиболее подходящего списка состояний.

[illegible]

In [3]:

```
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 from sklearn.metrics import f1_score, accuracy_score, confusion_matrix
10 from sklearn.model_selection import StratifiedKFold, KFold
11
12 from tqdm.notebook import tqdm
13
14 from sklearn.metrics import f1_score
15
16 from scipy.stats import norm
```

In [4]:

```
1 train = pd.read_csv('data-without-drift/train_clean.csv')
2 test  = pd.read_csv('data-without-drift/test_clean.csv')
```

Прежде чем реализовывать алгоритм самостоятельно, посмотрим, осмысленно ли это делать вообще.

Получим предсказания с использованием kfold кросс-валидации, и посмотрим на качество на валидации.

Реализацию алгоритма на данном этапе возьмем из паблик ноутбуков соревнования.

<https://www.kaggle.com/friedchips/the-viterbi-algorithm-a-complete-solution>  
(<https://www.kaggle.com/friedchips/the-viterbi-algorithm-a-complete-solution>)

In [5]:

```
1 class ViterbiClassifier:
2     def __init__(self):
3         self._p_trans = None
4         self._p_signal = None
5         self._p_in = None
6
7     def fit(self, x, y):
8         self._p_trans = self.markov_p_trans(y)
9         self._dists = []
10        self._states = len(np.unique(y))
11        for s in np.arange(y.min(), y.max() + 1):
12            self._dists.append((np.mean(x[y == s]), np.std(x[y == s])))
13
14        return self
15
16    def predict(self, x):
17        p_signal = self.markov_p_signal(x)
18        return self.viterbi(self._p_trans, p_signal, x)
19
20    def markov_p_signal(self, signal):
21        p_signal = np.zeros((self._states, len(signal)))
22        for k, dist in enumerate(self._dists):
23            p_signal[k, :] = norm.pdf(signal, *dist)
24
25        return p_signal
26
27    def markov_p_trans(self, states):
28        max_state = np.max(states)
29        states_next = np.roll(states, -1)
30        matrix = []
31        for i in tqdm(range(max_state + 1)):
32            current_row = np.histogram(states_next[states == i], bins=np.arange(
33                if np.sum(current_row) == 0:
34                    current_row = np.ones(max_state + 1) / (max_state + 1) # ...us
35                else:
36                    current_row = current_row / np.sum(current_row) # normalize to
37            matrix.append(current_row)
38        return np.array(matrix)
39
40    def viterbi(self, p_trans, p_signal, signal):
41        offset = 10**(-20)
42
43        p_trans_tlog = np.transpose(np.log2(p_trans + offset))
44        p_signal_tlog = np.transpose(np.log2(p_signal + offset))
45
46        T1 = np.zeros(p_signal.shape)
47        T2 = np.zeros(p_signal.shape)
48
49        T1[:, 0] = p_signal_tlog[0, :]
50        T2[:, 0] = 0
51
52        for j in tqdm(range(1, p_signal.shape[1])):
53            for i in range(len(p_trans)):
54                T1[i, j] = np.max(T1[:, j - 1] + p_trans_tlog[:, i] + p_signal
55                T2[i, j] = np.argmax(T1[:, j - 1] + p_trans_tlog[:, i] + p_sig
56
57        x = np.empty(p_signal.shape[1], 'B')
58        x[-1] = np.argmax(T1[:, p_signal.shape[1] - 1])
59        for i in reversed(range(1, p_signal.shape[1])):
```

```
60         x[i - 1] = T2[x[i], i]
61
62     return x
```

Проведём кросс-валидацию и получим out of fold предсказания для данного алгоритма.

In [7]:

```
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     model = ViterbiClassifier().fit(X_train_, y_train_)
21
22
23     # скор на валидации
24     preds = model.predict(X_valid)
25     oof[validation_index] = preds.reshape(-1,)
26
27     preds = np.round(np.clip(preds, 0, 10)).astype(int)
28     score = f1_score(y_valid, preds, average = 'macro')
29     scores.append(score)
30
31     # предсказание на тесте
32     preds = model.predict(X_test)
33     prediction += preds
34
35     print(f'score: {score}')
```

40% 2/5 [27:32<34:19, 686.34s/it]

100% 11/11 [00:02<00:00, 3.78it/s]

100% 999999/999999 [11:23<00:00, 1463.31it/s]

100% 1999999/1999999 [07:30<00:00, 4443.66it/s]

score: 0.7533647190608334

100% 11/11 [00:03<00:00, 3.09it/s]

100% 999999/999999 [03:45<00:00, 4434.41it/s]

100% 1999999/1999999 [11:34<00:00, 2877.98it/s]

score: 0.7530597098194463

100% 11/11 [03:53<00:00, 21.26s/it]

100% 999999/999999 [03:47<00:00, 4390.93it/s]

9% 184787/1999999 [00:52<06:54, 4380.43it/s]

```
-----
KeyboardInterrupt                                Traceback (most recent call
last)
<ipython-input-7-9ad60175c2dd> in <module>
    30
    31     # предсказание на тесте
----> 32     preds = model.predict(X_test)
    33     prediction += preds
    34

<ipython-input-5-245828350e11> in predict(self, x)
    16     def predict(self, x):
    17         p_signal = self.markov_p_signal(x)
----> 18         return self.viterbi(self._p_trans, p_signal, x)
    19
    20     def markov_p_signal(self, signal):

<ipython-input-5-245828350e11> in viterbi(self, p_trans, p_signal, sig
nal)
    55         for i in range(len(p_trans)):
    56             T1[i, j] = np.max(T1[:, j - 1] + p_trans_tlog[
:, i] + p_signal_tlog[j, i])
----> 57             T2[i, j] = np.argmax(T1[:, j - 1] + p_trans_tl
og[:, i] + p_signal_tlog[j, i])
    58
    59         x = np.empty(p_signal.shape[1], 'B')
```

KeyboardInterrupt:

Можно видеть, что качество на валидации очень низкое, поэтому не стоит ожидать высокого качества на LB, и использовать в стекинге, т.к. модель верхнего уровня в стекинге/блендинге, скорее всего, просто возьмёт ответы данной модели с коэффициентом 0.

В паблик ноутбуках с помощью этого алгоритма получают достаточно высокое качество, т.к. используют его для сложных батчей на тестовой и обучающей выборке (предполагается, что если значения сигналов в различных батчах схожи, то для генерации данных в них используются схожие макровские процессы).

Поэтому, использовать данную модель в стекинге не представляется возможным.

Стоит отметить, что в результате работы алгоритма на выходе мы получаем матрицу вероятностей переходов. Её можно использовать в качестве новой фичи. Также, данная модель будет полезна при обучении остальных рассмотренных моделей на отдельных батчах.

In [ ]:

1	
---	--