

In [2]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import statsmodels.api as sm
5 import seaborn as sns
6 from tqdm import tqdm_notebook
7 import warnings
8 from statsmodels.tsa.stattools import kpss
9 from statsmodels.stats.multitest import multipletests
10 from statsmodels.tsa.holtwinters import ExponentialSmoothing
11
12 warnings.filterwarnings('ignore')
13 sns.set_style('whitegrid')
```

## 0. Загрузка данных

In [3]:

```
1 train = pd.read_csv('data/train.csv')
2 test = pd.read_csv('data/test.csv')
```

In [3]:

```
1 print(train.shape)
2 print(test.shape)
3 print(train.head())
```

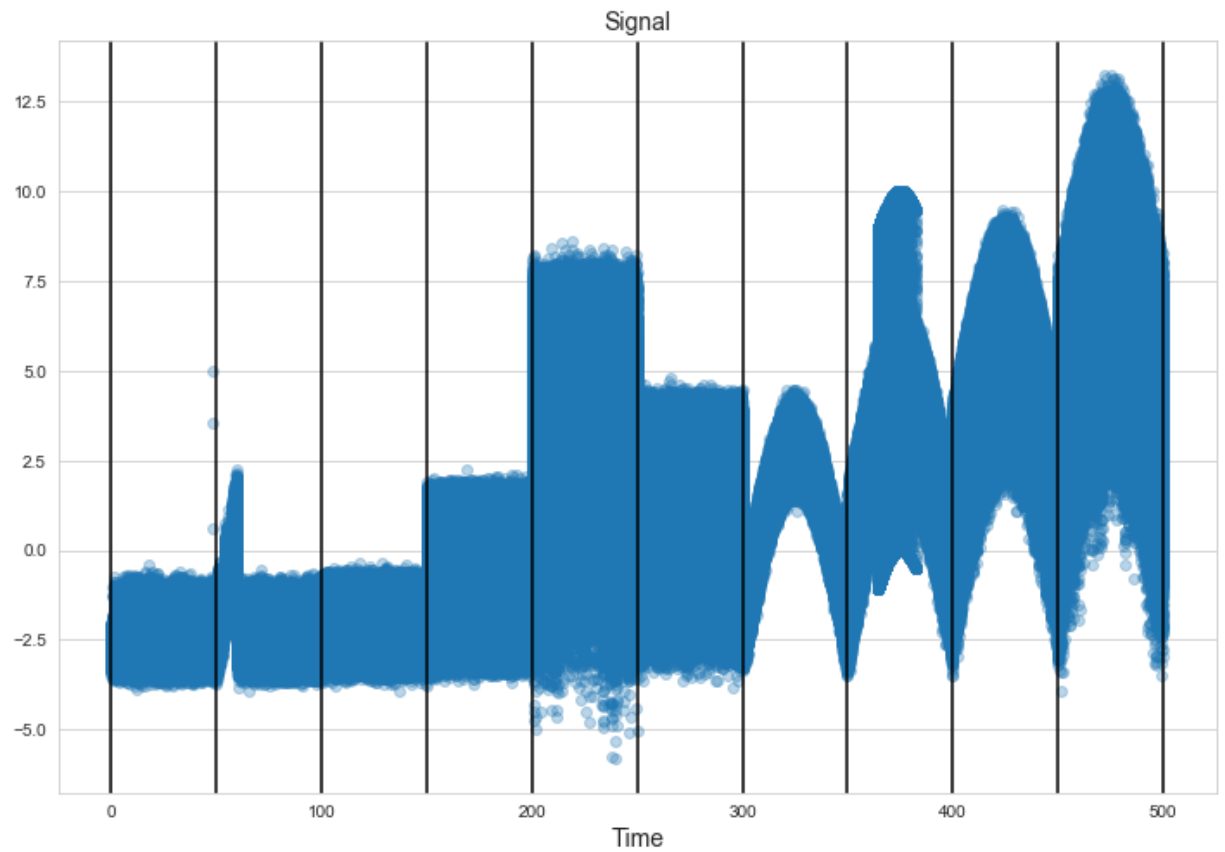
```
(5000000, 3)
(2000000, 2)
   time  signal  open_channels
0  0.0001 -2.7600             0
1  0.0002 -2.8557             0
2  0.0003 -2.4074             0
3  0.0004 -3.1404             0
4  0.0005 -3.1525             0
```

## 1. Сигнал и кол-во каналов в зависимости от времени

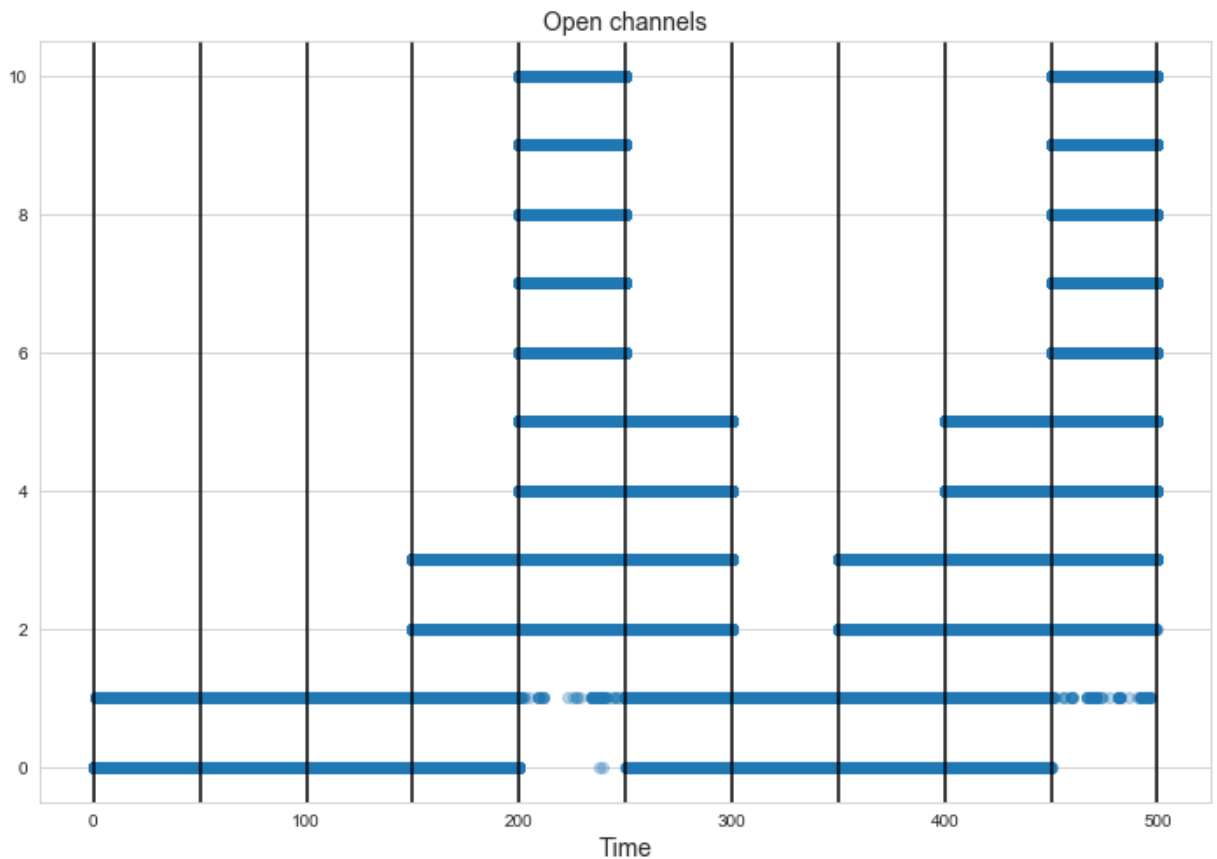
In [79]:

```
1 def draw_time_series(time, data, title, plot_type='plot',
2                       multibatch=True):
3     plt.figure(figsize=(12,8))
4     plt.title(title, fontsize=14)
5
6     if plot_type == 'plot':
7         plt.plot(time, data, alpha=0.9)
8     if plot_type == 'scatter':
9         plt.scatter(time, data, alpha=0.3)
10
11     if multibatch:
12         for i in range(11):
13             plt.axvline(50*i, color='black')
14
15     plt.xlabel('Time', fontsize=14)
16     plt.show()
17
18
19 def draw_signal_with_channels(time, signal, channels, title):
20     plt.figure(figsize=(12, 8))
21     plt.title(title, fontsize=14)
22     plt.plot(time, signal, label='signal')
23     plt.plot(time, channels, label='channels')
24     plt.legend(fontsize=14)
25     plt.xlabel('Time', fontsize=14)
26     plt.show()
```

```
In [24]: 1 draw_time_series(train.time, train.signal,  
2                      'Signal', plot_type='scatter')
```



```
In [25]: 1 draw_time_series(train.time, train.open_channels,
2          'Open channels', plot_type='scatter')
```



### Наблюдения:

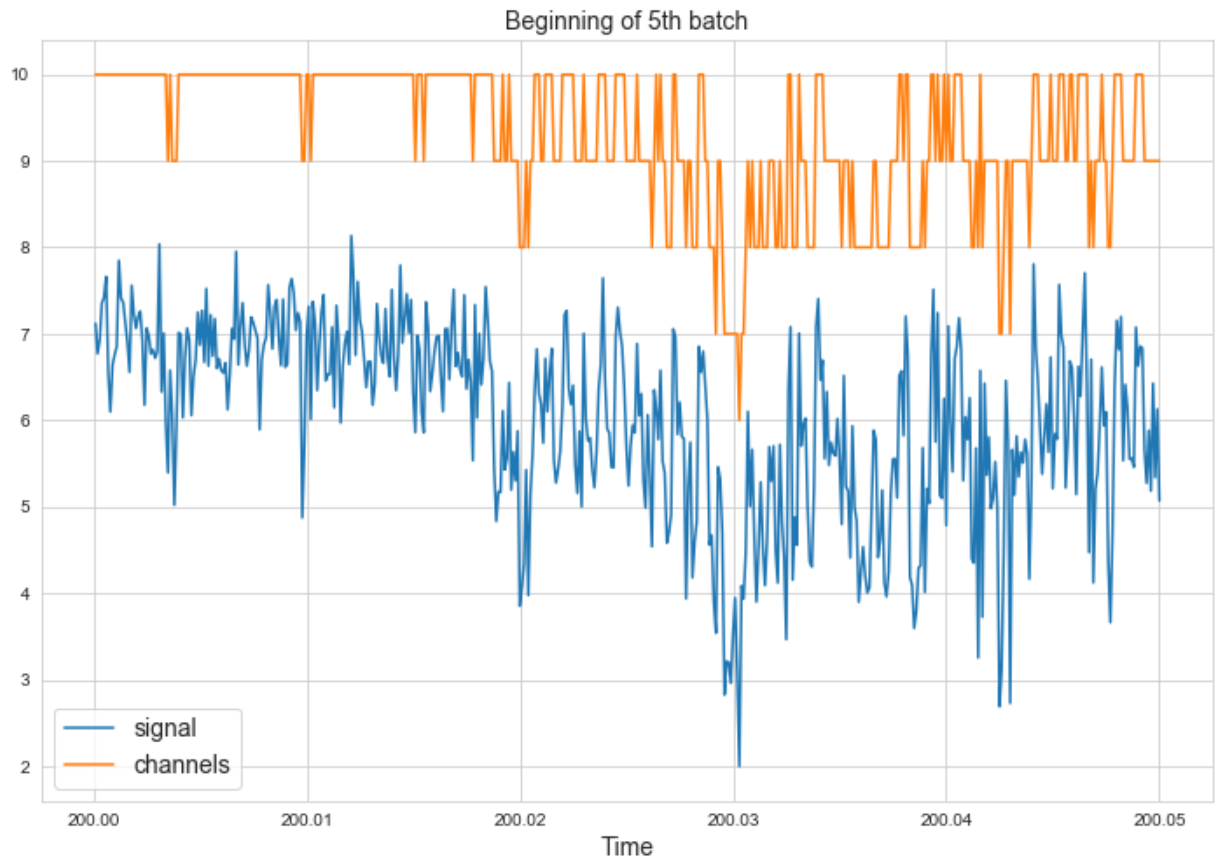
1) 6 из 10 временных батчей по 50 секунд сигнал колеблется с постоянной дисперсией, причем за короткий промежуток времени успевает принять весь диапазон своих значений в пределах этого батча. Большое количество открытых каналов получается в батчах с большей дисперсией, и соответственно, с большим значением тока, что предсказуемо. Позже проверим эту зависимость статистическими методами.

2) Последние 4 батча напоминают синусоиду, и так же в батчах с меньшими значениями сигнала количество открытых каналов тоже меньше.

Пятый батч в приближении:

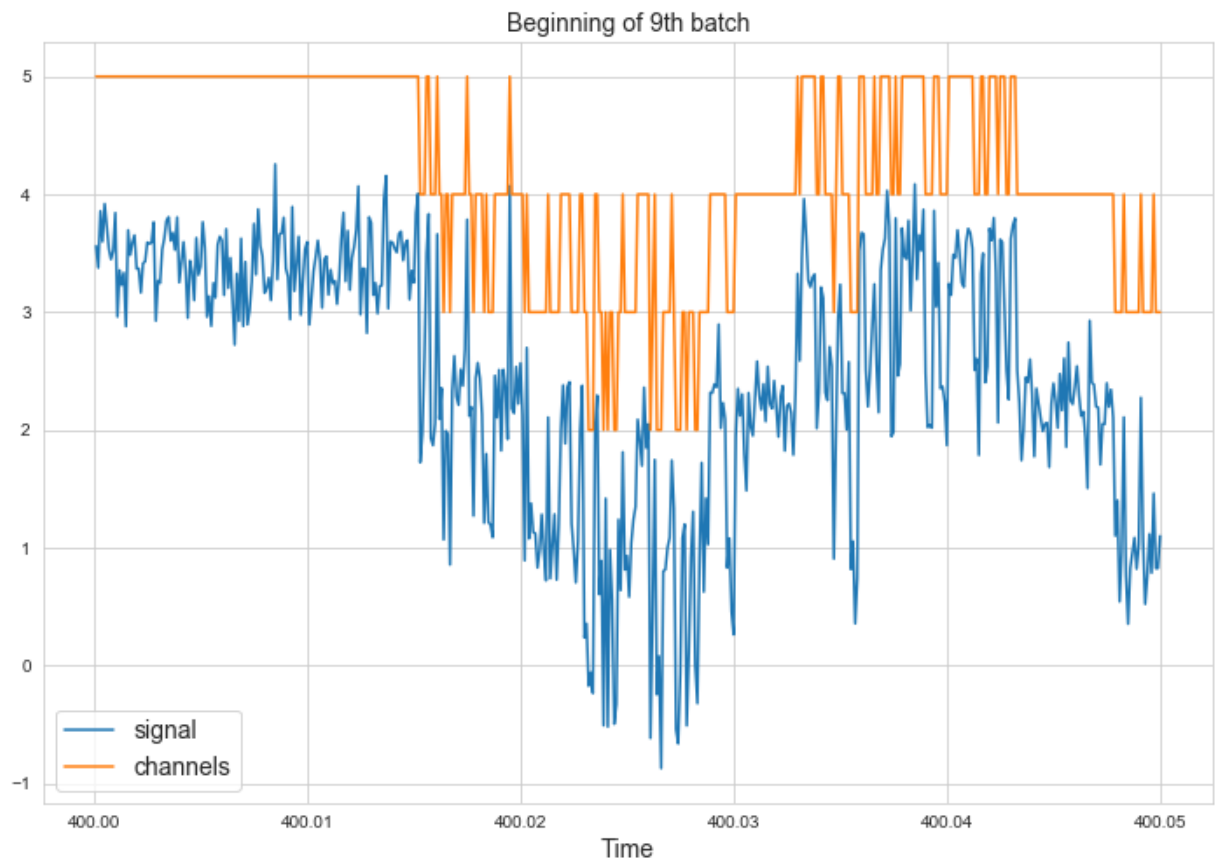
In [35]:

```
1 left = 500*(10**3)*4
2 num_steps = 500
3 right = left + num_steps
4 draw_signal_with_channels(train.time[left:right],
5                           train.signal[left:right],
6                           train.open_channels[left:right],
7                           'Beginning of 5th batch')
```



Девятый батч в приближении:

```
In [36]: 1 left = 500*(10**3)*8
2 num_steps = 500
3 right = left + num_steps
4 draw_signal_with_channels(train.time[left:right],
5                           train.signal[left:right],
6                           train.open_channels[left:right],
7                           'Beginning of 9th batch')
```



**Наблюдения:** когда количество открытых каналов постоянно, ток может колебаться в некотором диапазоне шириной около двух единиц, закрытие одного-двух каналов ведёт к мгновенному падению тока. "Волатильность" тока при постоянном количестве открытых каналов наводит на мысль, что полезно будет посмотреть на сглаженный график тока, в дальнейшем это может пригодиться для генерации новых признаков (со сглаженным сигналам моделям будет работать явно проще).

## 2. Стационарность временных рядов

Прежде чем что-то думать про тренд и сезонность, проверим, а не являются ли наши временные ряды стационарными. Это поможет понять, в каком направлении стоит двигаться в дальнейшем анализе.

Проверять стационарность будем критерием KPSS, а затем к результатам применим множественную проверку гипотез:

```
In [68]: 1 pvalues = []
2 batch_size = 500*(10**3)
3
4 for i in range(10):
5     current_batch = train.signal[i*batch_size:(i+1)*batch_size]
6     current_batch = np.array(current_batch)
7     pvalues.append(kpss(current_batch)[1])
8
9 multipletests(pvalues)[0]
```

```
Out[68]: array([False, False, False, False, False, False, False, False, False,
False])
```

**Результат:** гипотеза о стационарности не отвергается ни для одного временного ряда. Значит не стоит особо возиться с трендом и сезонностью - скорее всего их нет.

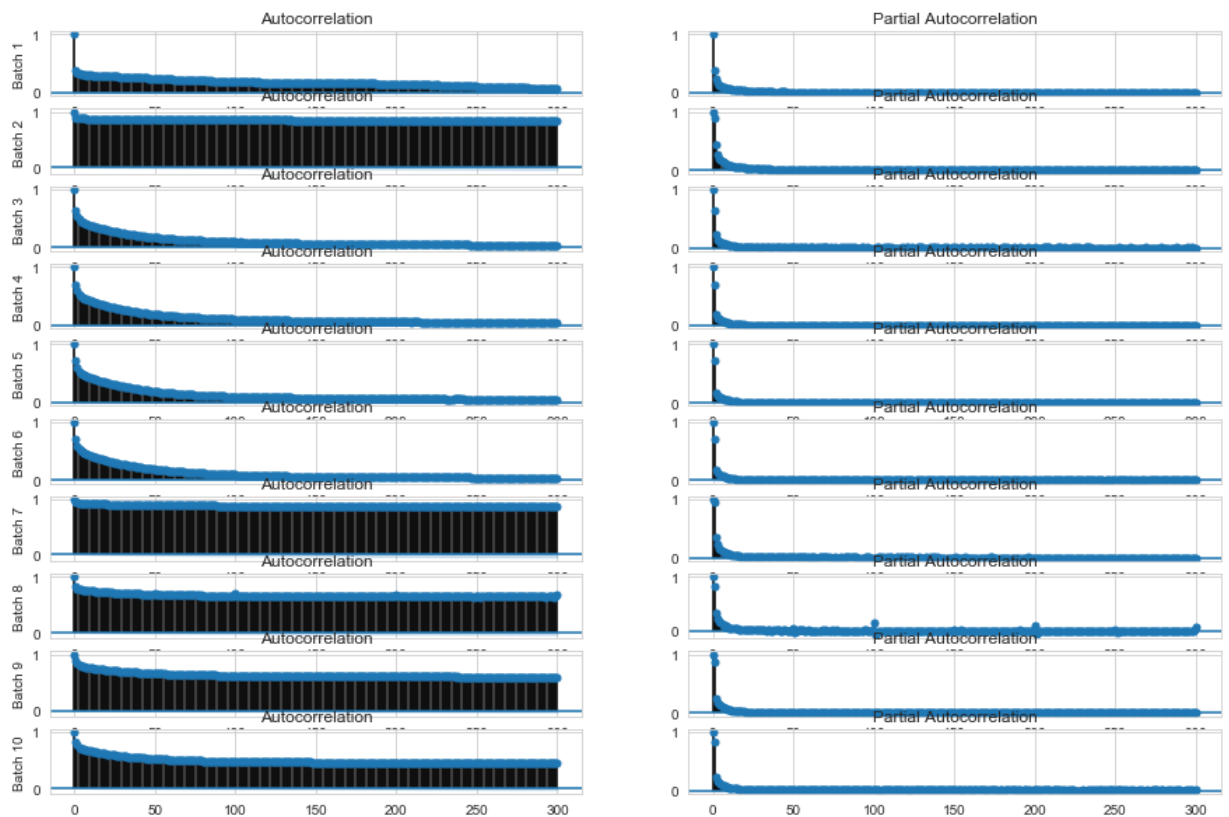
### 3. Выявление сезонности

Тем не менее, нарисовать коррелограммы - никогда не бывает лишним, поэтому все-таки попробуем найти какую-нибудь сезонность в рядах. В конце концов, результат проверки гипотезы о стационарности получился статистически незначимым.

```
In [83]: 1 def draw_corr_grams(time, signal):
2     batch_size = 500*(10**3)
3     num_batches = 10
4
5     plt.subplots(num_batches, 2, figsize=(15,10))
6
7     for i in tqdm_notebook(range(num_batches)):
8         current_batch = signal[i*batch_size:(i+1)*batch_size]
9
10        ax = plt.subplot(num_batches, 2, (i*2)+1)
11        ax.set_ylabel(f'Batch {i+1}')
12        sm.graphics.tsa.plot_acf(current_batch, lags=300, ax=ax)
13
14        ax = plt.subplot(num_batches, 2, (i*2)+2)
15        sm.graphics.tsa.plot_pacf(current_batch, lags=300, ax=ax)
16
17    plt.show()
```

```
In [84]: 1 draw_corr_grams(train.time, train.signal)
```

```
HBox(children=(IntProgress(value=0, max=10), HTML(value='')))
```



**Наблюдения:** сезонности сигнала нет ни для одного батча.

## 4. Сглаживание - скользящее среднее

Исследуем способ сглаживания сигнала путём взятия скользящего среднего с разной шириной окна.



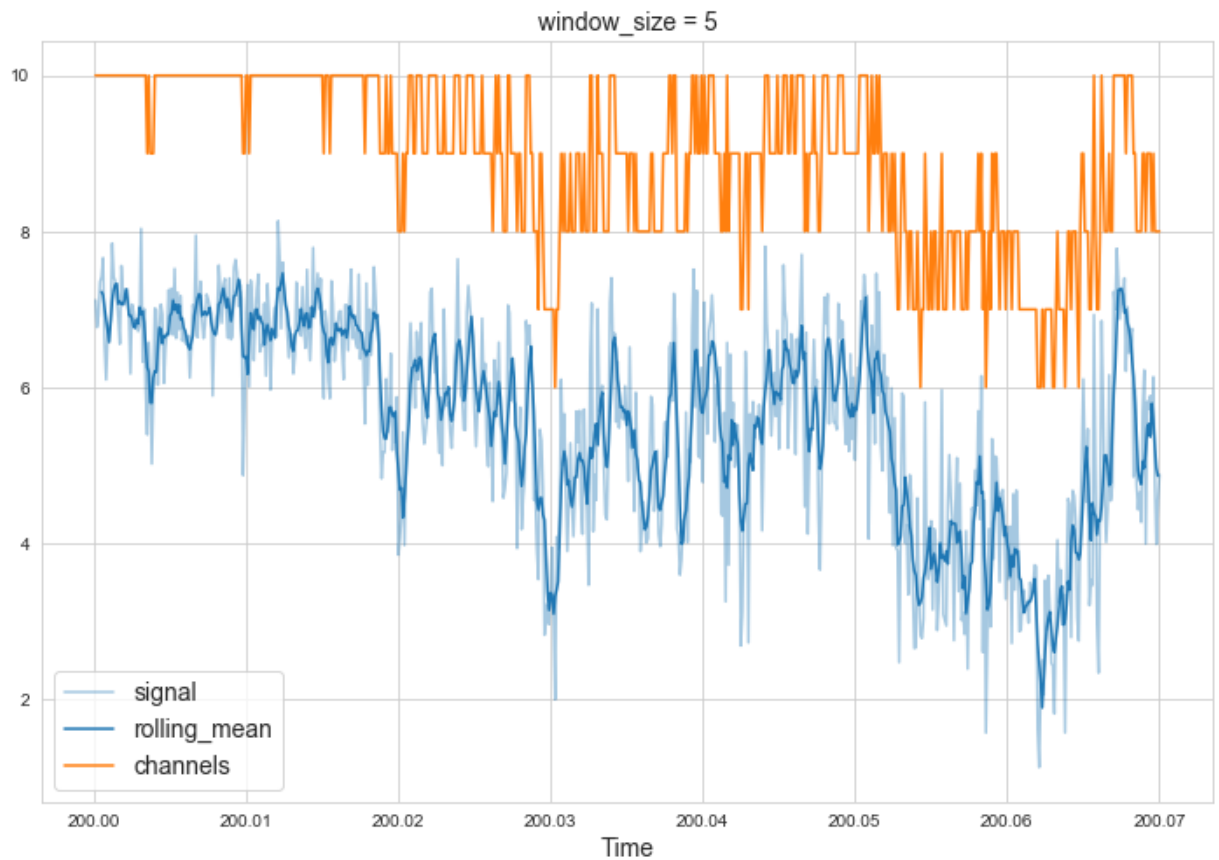
```
In [52]: 1 def signal_channels_rolling(time, signal, channels, title,
2         rolling, window_size):
3         rolling_time = np.array(rolling.dropna().index).astype('float')
4         rolling_time = rolling_time / 10000
5         rolling_value = np.array(rolling.dropna())
6
7         plt.figure(figsize=(12, 8))
8         plt.title(title, fontsize=14)
9         plt.plot(time, signal, label='signal', alpha=0.4,
10                color='C0')
11
12         plt.plot(rolling_time[window_size:],
13                rolling_value[window_size:],
14                label='rolling_mean', color='C0')
15         plt.plot(time, channels, label='channels', color='C1')
16         plt.legend(fontsize=14)
17         plt.xlabel('Time', fontsize=14)
18         plt.show()
```

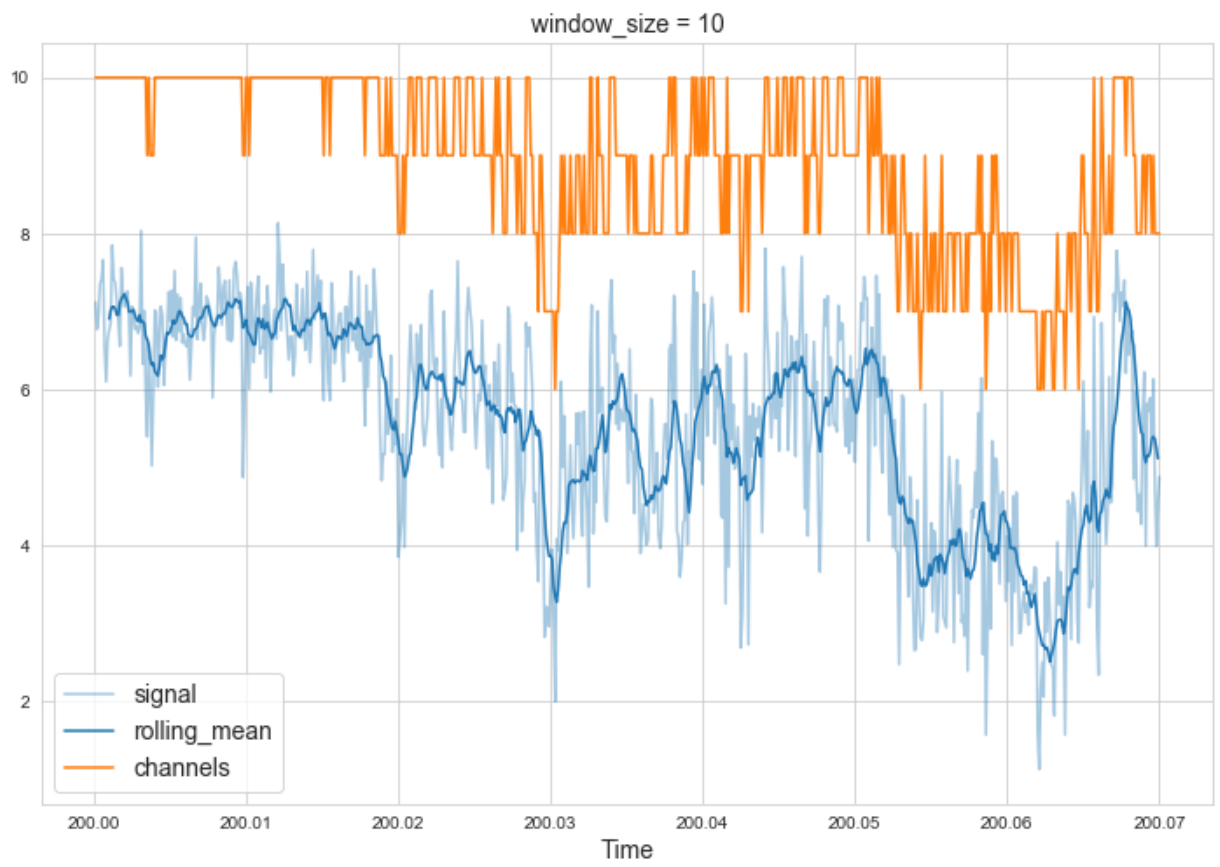
```
In [59]: 1 window_sizes = [5, 10, 20, 30, 60, 100]
2         rollings = []
3
4         for window_size in tqdm_notebook(window_sizes):
5             rollings.append(train.signal
6                             .rolling(window=window_size)
7                             .mean())
```

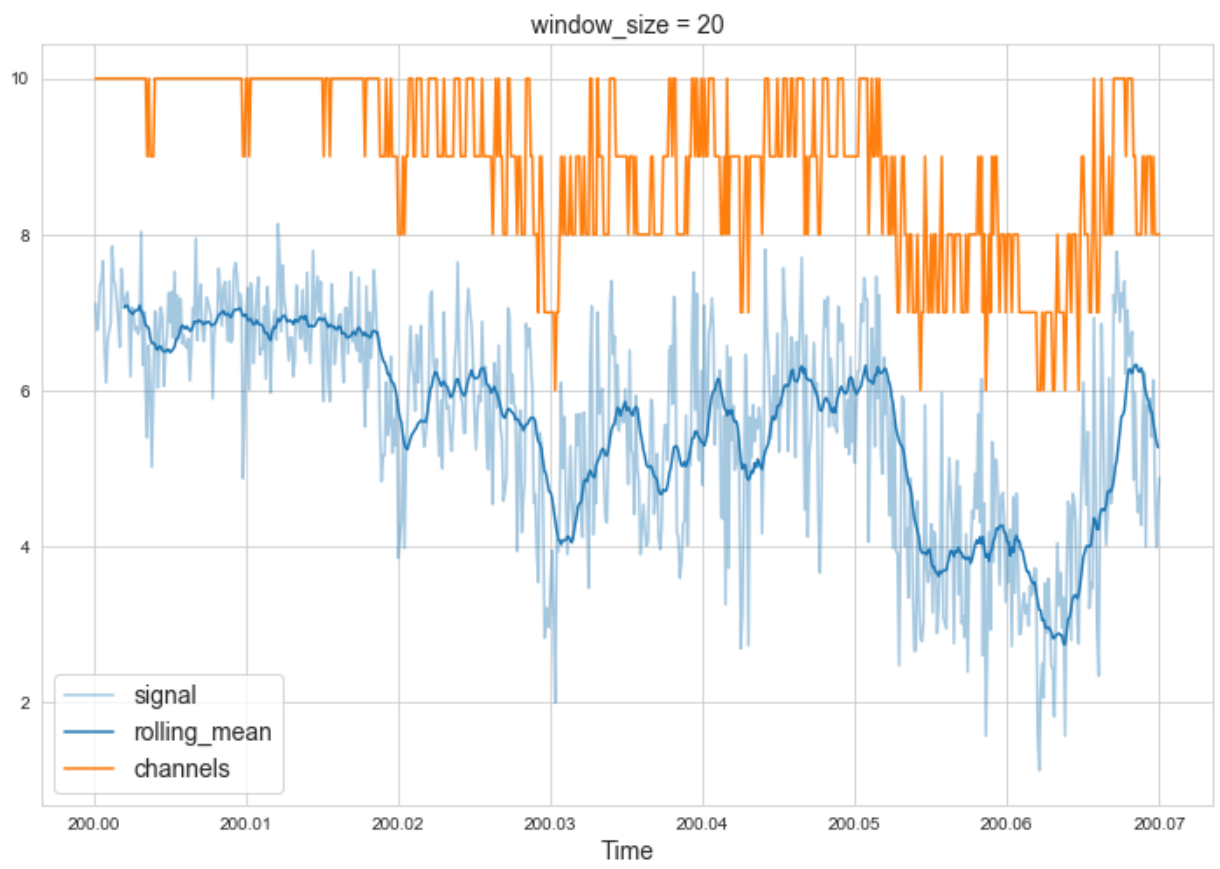
```
HBox(children=(IntProgress(value=0, max=6), HTML(value='')))
```

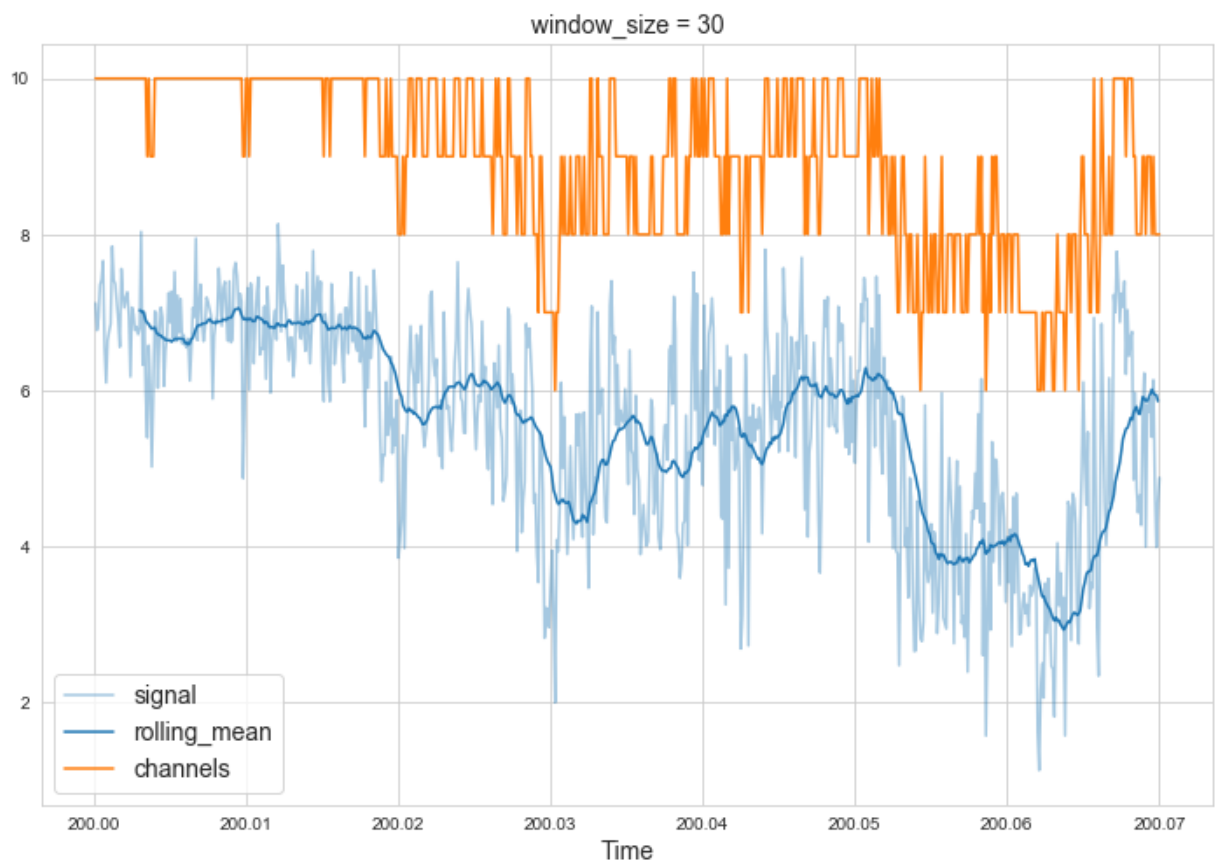
In [60]:

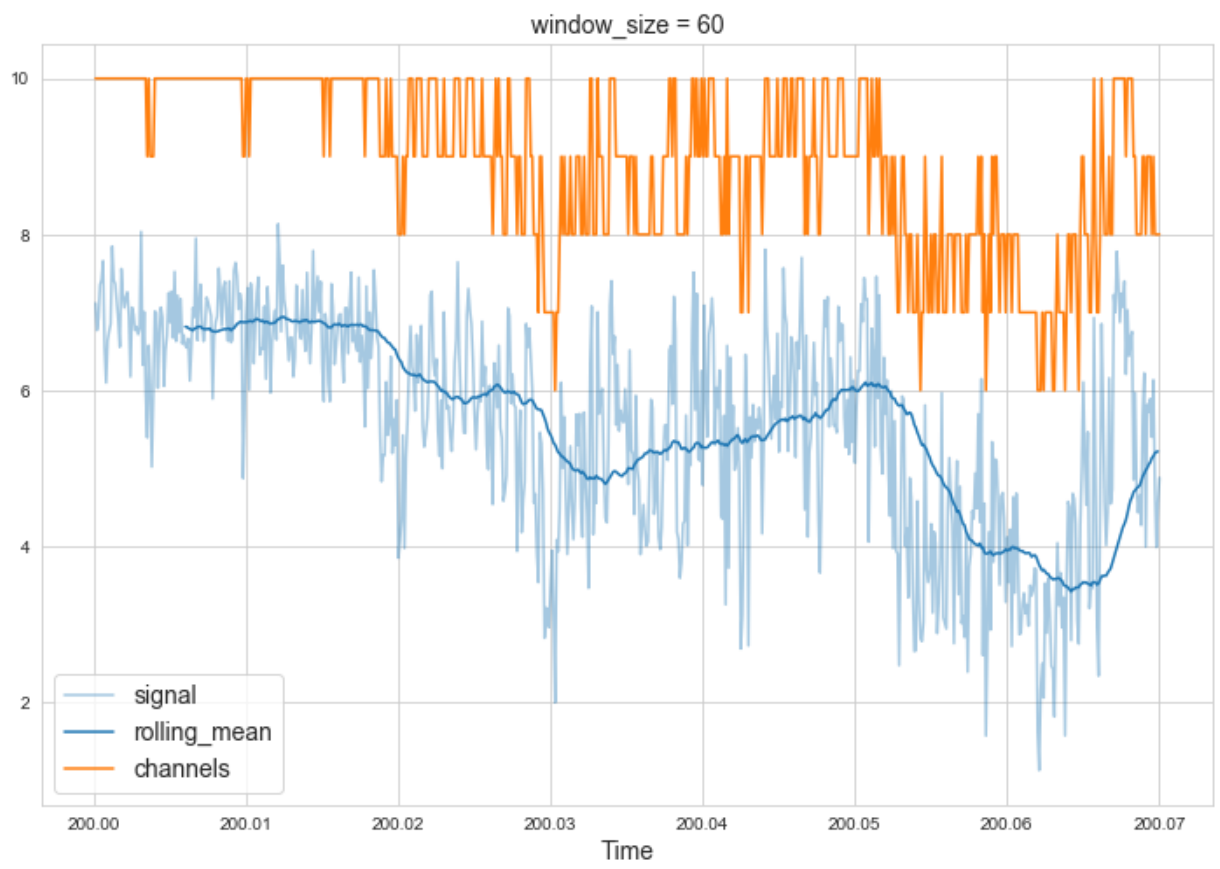
```
1 left = 500*(10**3)*4
2 num_steps = 700
3 right = left + num_steps
4
5 for i, window_size in enumerate(window_sizes):
6     signal_channels_rolling(train.time[left:right],
7                             train.signal[left:right],
8                             train.open_channels[left:right],
9                             f'window_size = {window_size}',
10                            rollings[i][left:right],
11                            window_size=window_size)
```

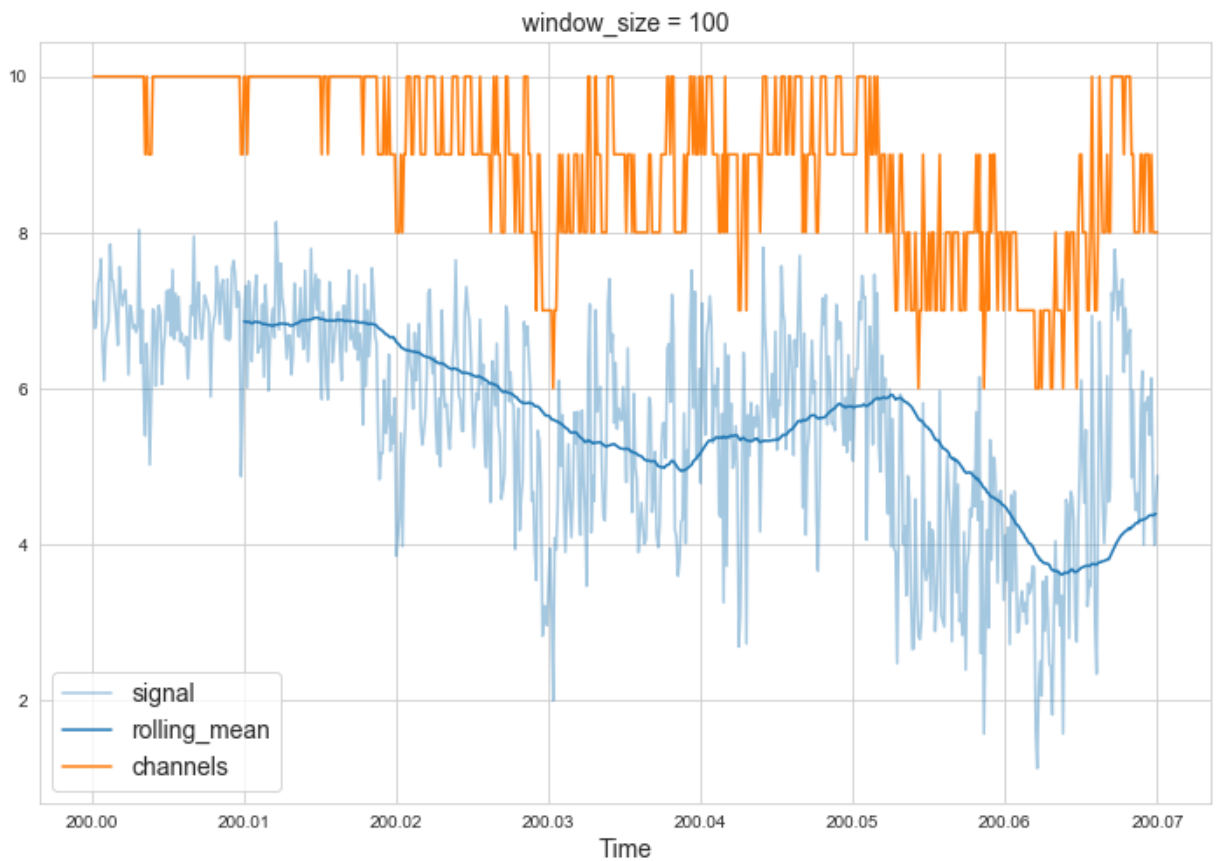












**Наблюдения:** рисунок скользящего среднего хорошо повторяет рисунок изменения количества открытых каналов, а значит при feature engineering нужно будет добавить побольше различных скользящих статистик. Особенно хорошо смотрится скользящее среднее по окну размера 20, там и дисперсия хорошо сглаживается, и мелкие изменения количества открытых каналов все равно остаются заметными.

## 5. Экспоненциальное сглаживание

Рассмотрим также вариант экспоненциального сглаживания.

In [74]:

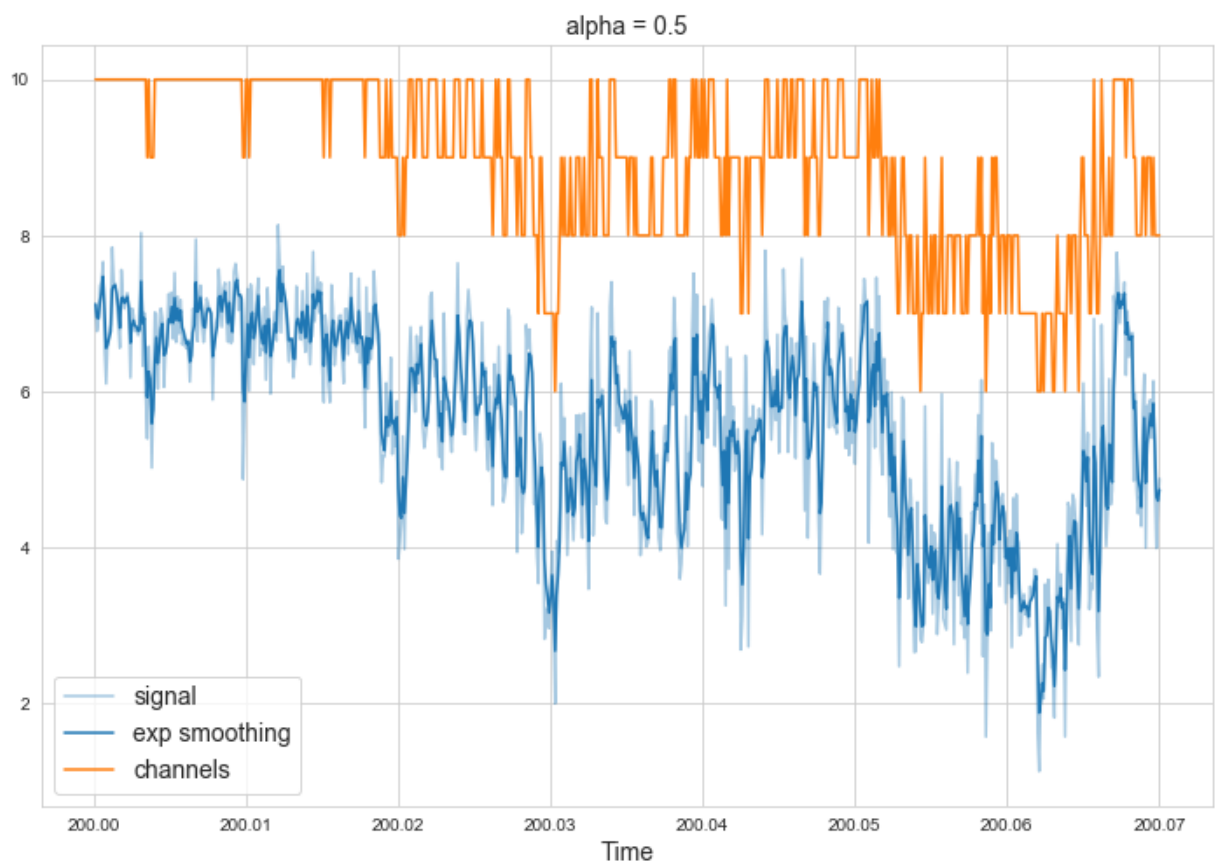
```
1 def exponential_smoothing(y, alpha):
2     res = np.zeros(len(y))
3     res[0] = y[0]
4
5     for i in range(1, len(y)):
6         res[i] = res[i-1] + alpha*(y[i] - res[i-1])
7
8     return res
9
10
11 def signal_channels_exp(time, signal, channels, title, exp):
12     plt.figure(figsize=(12, 8))
13     plt.title(title, fontsize=14)
14     plt.plot(time, signal, label='signal', alpha=0.4,
15              color='C0')
16     plt.plot(time, exp,
17              label='exp smoothing', color='C0')
18     plt.plot(time, channels, label='channels', color='C1')
19     plt.legend(fontsize=14)
20     plt.xlabel('Time', fontsize=14)
21     plt.show()
```

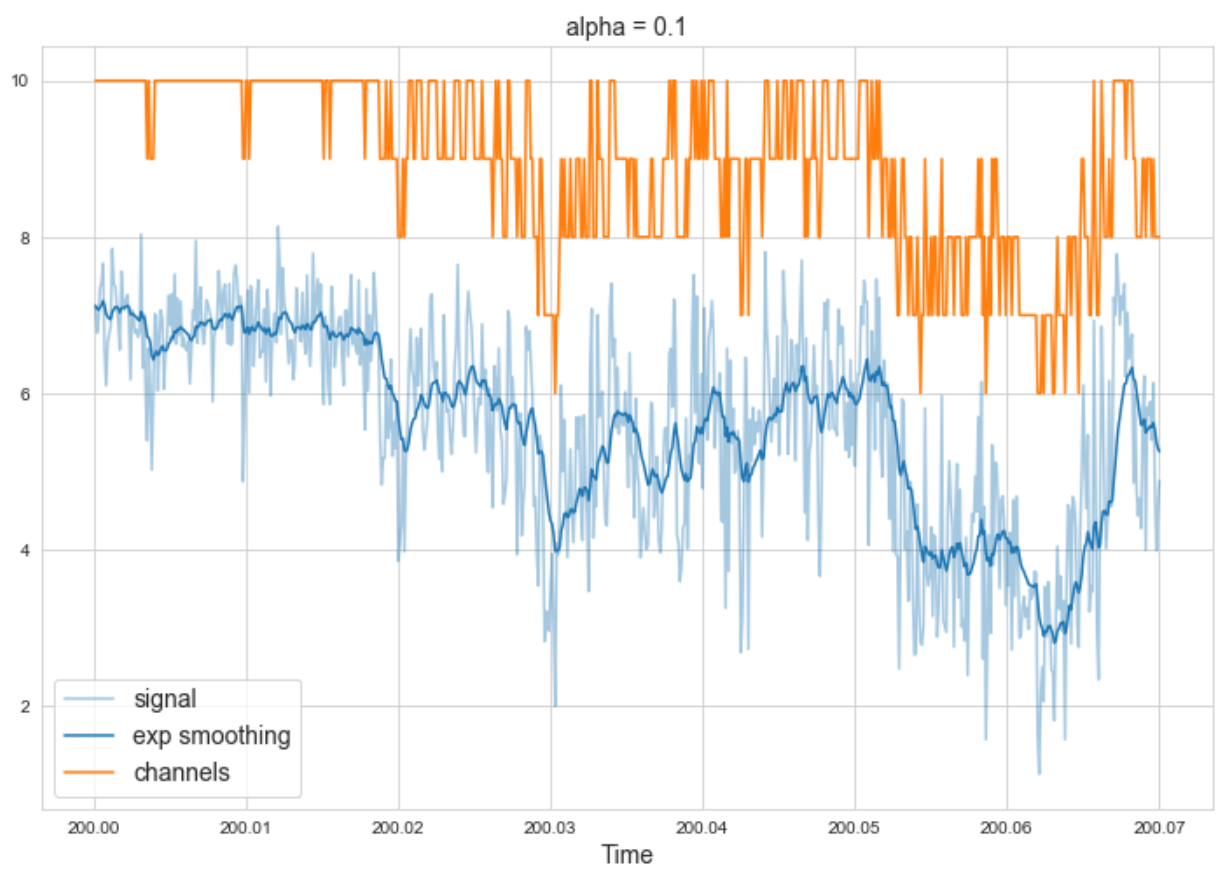
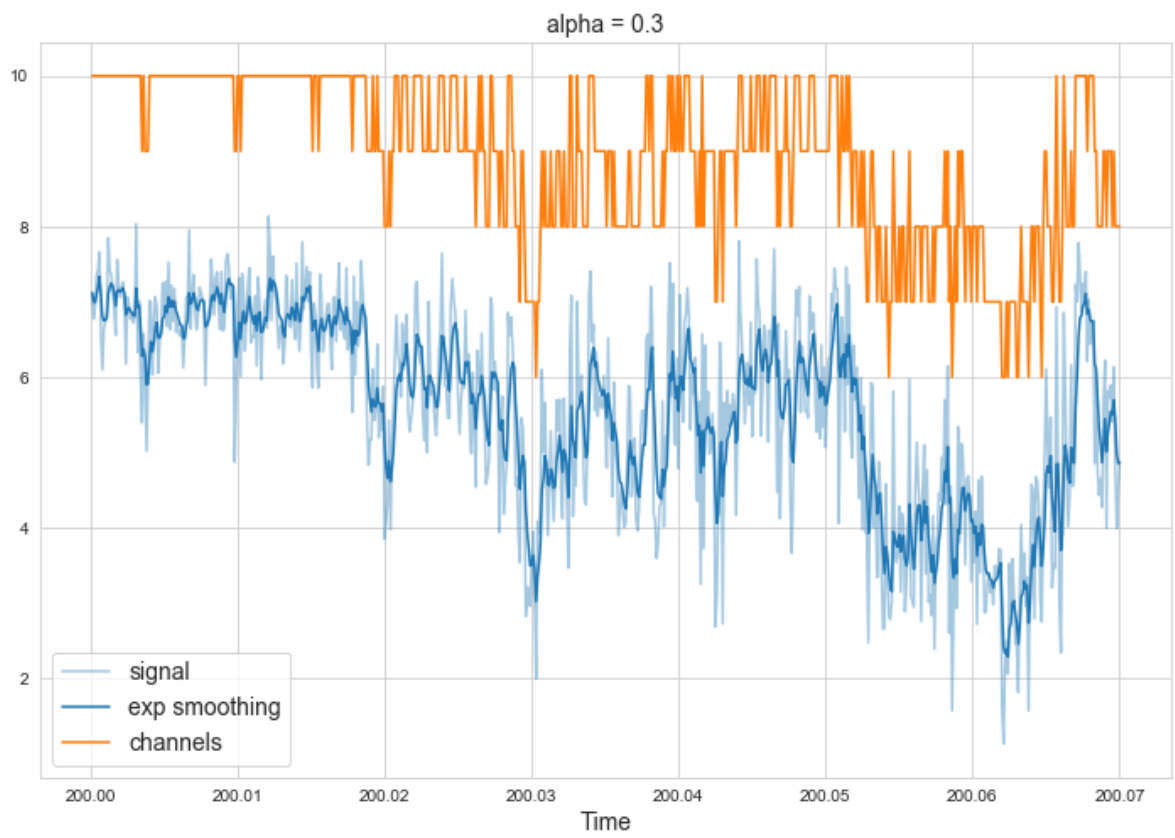


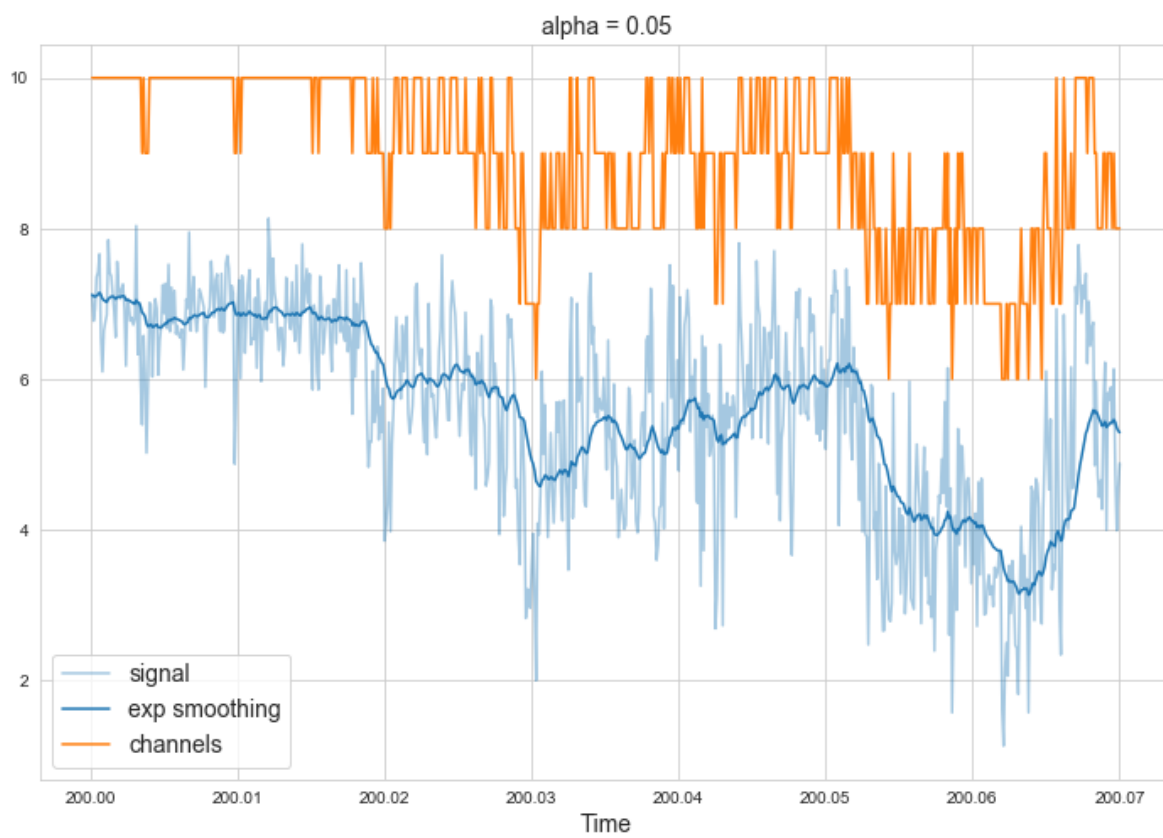
In [78]:

```
1 alphas = [0.5, 0.3, 0.1, 0.05, 0.01]
2
3 left = 500*(10**3)*4
4 num_steps = 700
5 right = left + num_steps
6
7 for i, alpha in tqdm_notebook(enumerate(alphas)):
8     exp = exponential_smoothing(np.array(train.signal[left:right]),
9                                alpha=alpha)
10    signal_channels_exp(train.time[left:right],
11                       train.signal[left:right],
12                       train.open_channels[left:right],
13                       f'alpha = {alpha}',
14                       exp)
```

HBox(children=(IntProgress(value=1, bar\_style='info', max=1), HTML(value='')))





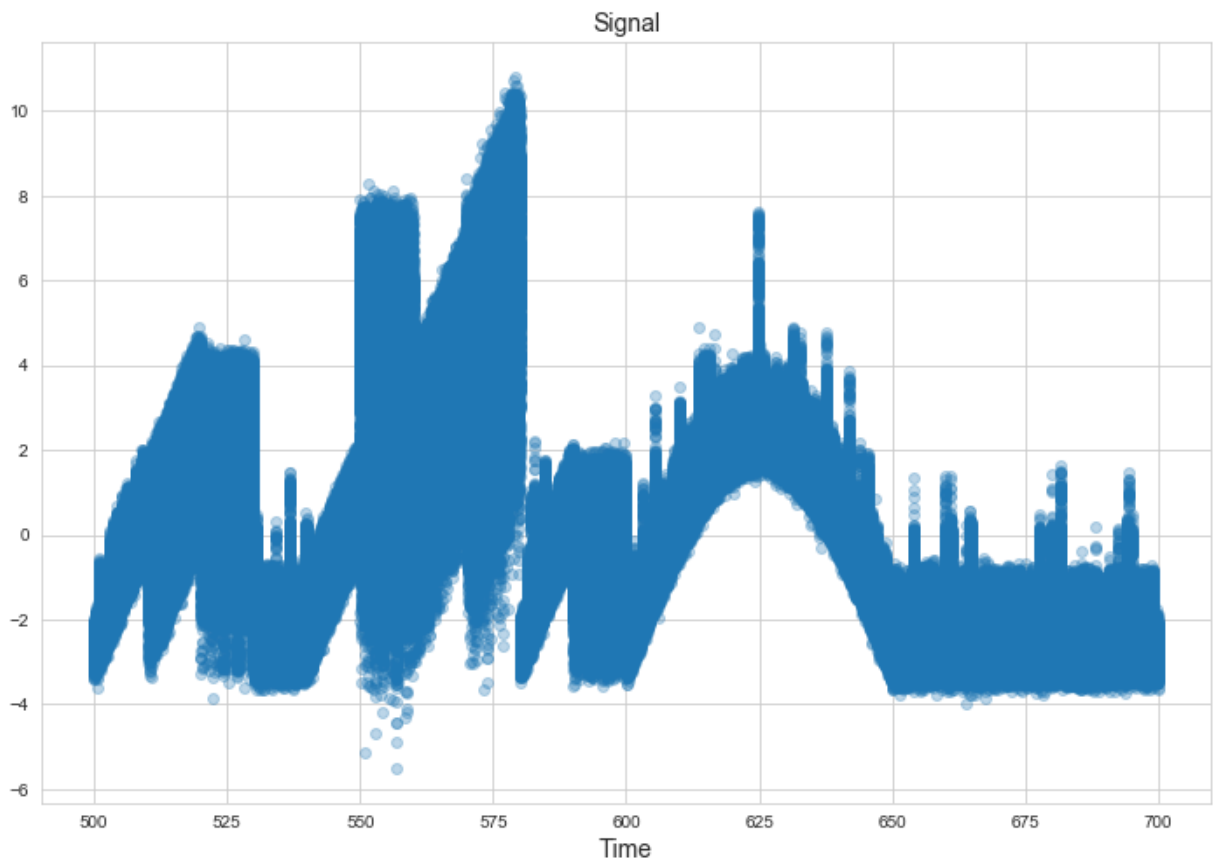




**Наблюдения:** при  $\alpha = 0.05$ , сглаженный график учитывает и мелкие колебания количества открытых каналов, и хорошо снижает волатильность. Учтем это также при генерации новых признаков.

## 6. График тестовой выборки

```
In [80]: 1 draw_time_series(test.time, test.signal,
          2 'Signal', plot_type='scatter', multibatch=False)
```

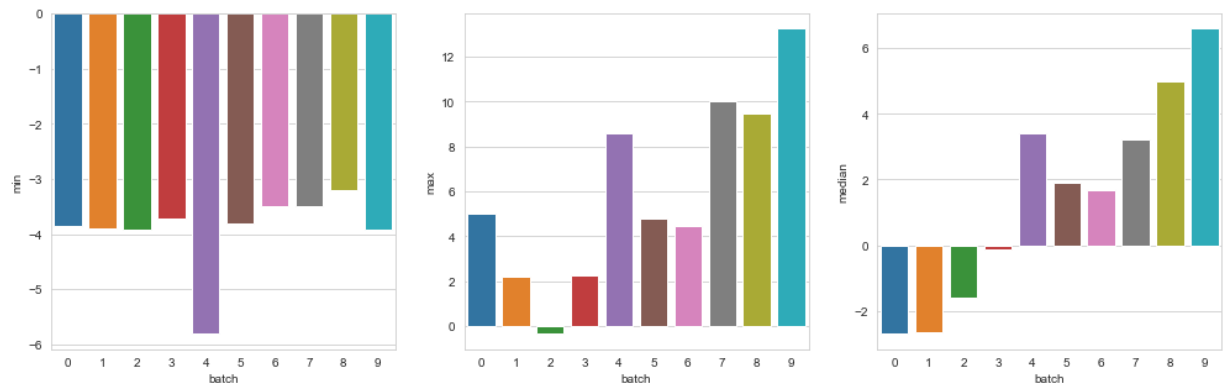


**Наблюдения:** видно, что к батчам на тесте тоже были применены некоторые преобразования, похожие на синус и линейную функцию. При обучении моделей нужно будет попробовать избавиться от этого, а еще стоит попробовать применять разные модели для разных рисунков сигнала (например, предсказывать третий батч на тесте моделью, обучающейся на последних 4-х батчах из теста, т.к. рисунок сигнала у них похож).

## 7. Статистики, агрегированные по батчам

In [4]:

```
1 def plot_batch_stats(df, batch_size=50 * 10**4):
2     df = df.copy()
3
4     stats = ['min', 'max', 'median']
5     df['batch'] = df.index // batch_size
6     by_batch = df.groupby(['batch'])['signal'].agg(stats).reset_index()
7
8
9     plt.figure(figsize=(17, 5))
10    for i, curr_stat in enumerate(stats):
11        plt.subplot(1, 3, i+1)
12
13        sns.barplot(y=curr_stat, x='batch', data=by_batch)
14
15    plot_batch_stats(train)
```



**Наблюдение:** если собирать статистики по данным в пределах конкретных батчей, то с помощью них можно четко различить батчи друг от друга, а значит это однозначно надо внедрить в feature\_engineering.