

# One Feature Model Scores LB 0.930!

In this notebook, we will explore the Kaggle Ion Comp data and explore a one feature model. The LB result of 0.930 is enlightening.

Here we manually remove signal drift. Note that it is better to use machine learning to remove drift, but doing it by hand once allows us to understand its nature and build better models later.

## Load Libraries and Data

```
In [1]: import numpy as np
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import matplotlib.pyplot as plt
import os
from sklearn.metrics import f1_score
import graphviz
from sklearn import tree

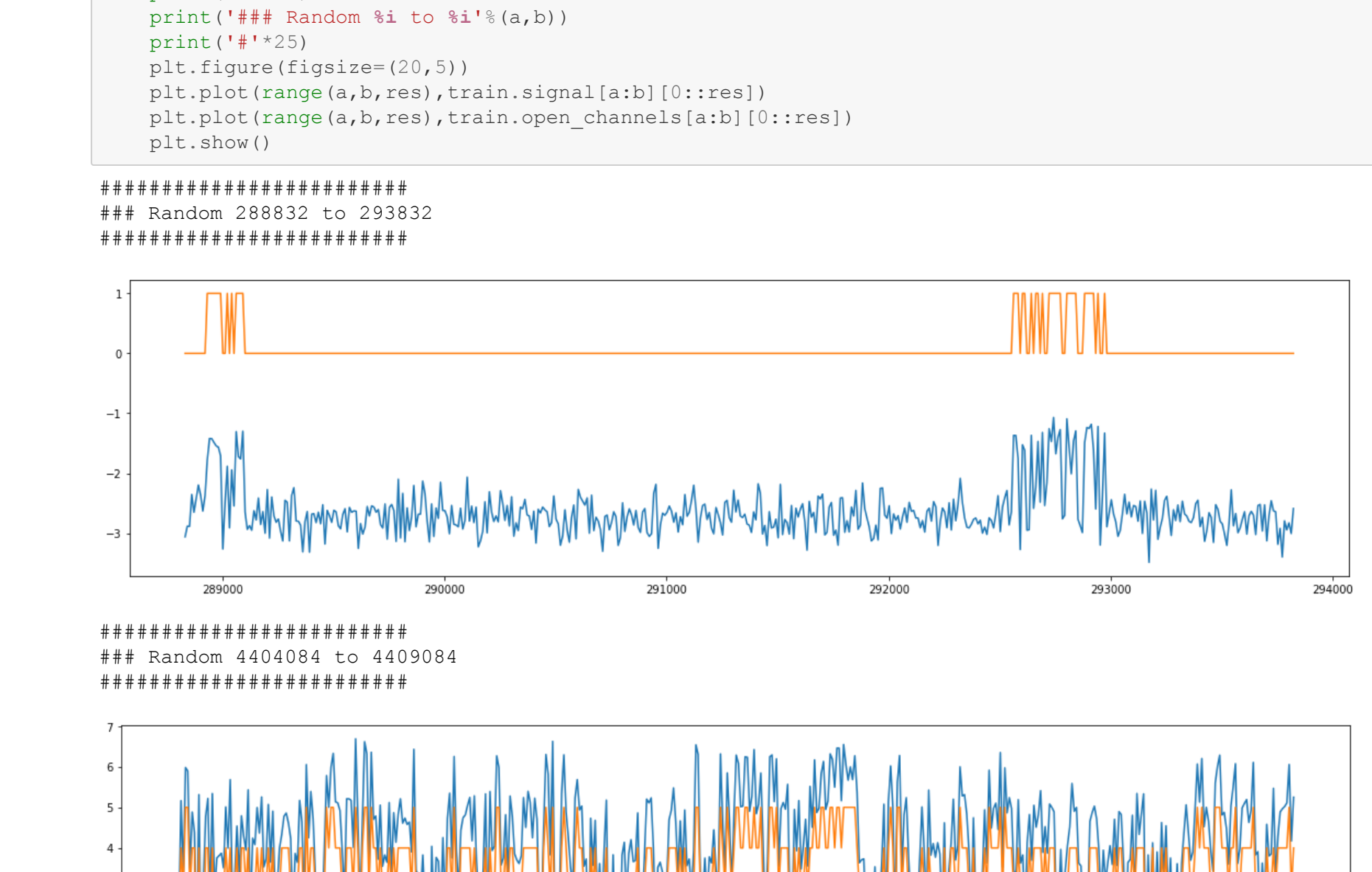
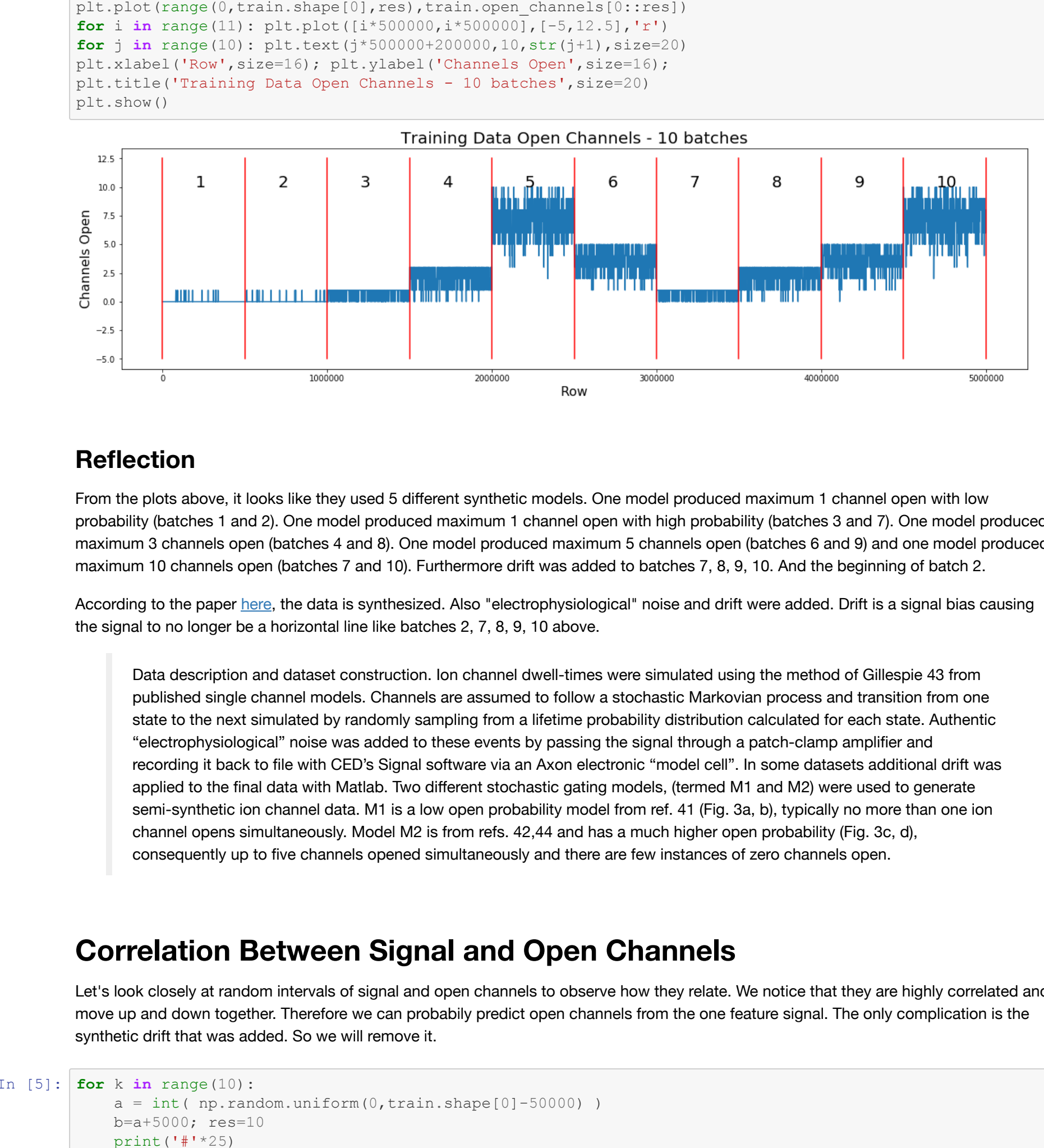
In [2]: test = pd.read_csv('../input/liverpool-ion-switching/test.csv')
train = pd.read_csv('../input/liverpool-ion-switching/train.csv')
train.head()
```

```
Out [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

## Description of Data

The training data is recordings in time. At each 10,000th of a second, the strength of the signal was recorded and the number of ion channels open was modelled. It is our task to build a model that predicts the number of open channels from signal at each time step. Furthermore we are told that the data was recorded in batches of 50 seconds. Therefore each 500,000 rows is one batch. The training data contains 10 batches and the test data contains 4 batches. Let's display the number of open channels and signal strength together for each training batch.



## Reflection

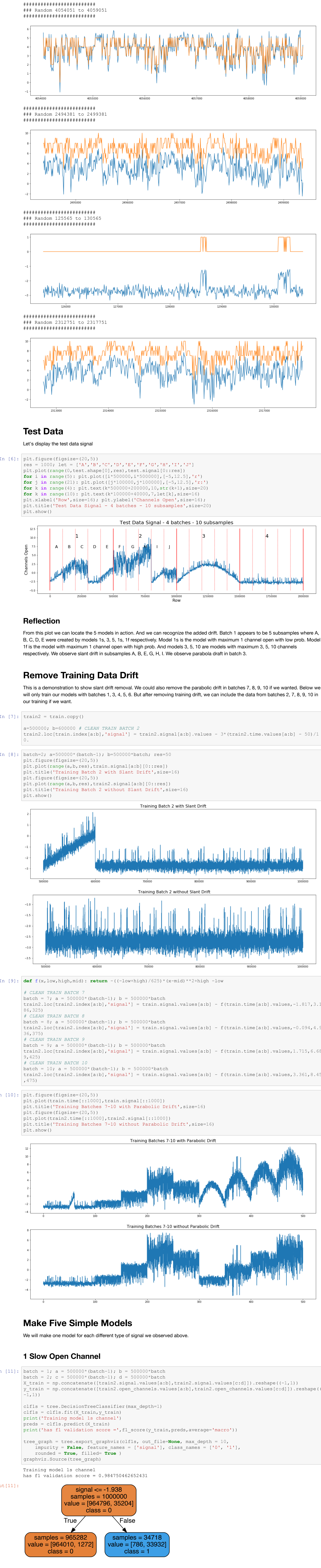
From the plots above, it looks like they used 5 different synthetic models. One model produced maximum 1 channel open with low probability (batches 1 and 2). One model produced maximum 1 channel open with high probability (batches 3 and 7). One model produced maximum 3 channels open (batches 4 and 8). One model produced maximum 5 channels open (batches 6 and 9) and one model produced maximum 10 channels open (batches 7 and 10). Furthermore drift was added to batches 7, 8, 9, 10. And the beginning of batch 2.

According to the paper [here](#), the data is synthesized. Also "electrophysiological" noise and drift were added. Drift is a signal bias causing the signal to no longer be a horizontal line at batches 2, 7, 8, 9, 10 above.

Data description and dataset construction. Ion channel dwell-times were simulated using the method of Gillespie 43 from published single channel models. Channels are assumed to follow a stochastic Markovian process and transition from one state to the next simulated by randomly sampling from a lifetime probability distribution calculated for each state. Authentic "electrophysiological" noise was added to these events by passing the signal through a patch-clamp amplifier and recording it back to file with CED's Signal software via an Axon electronic "model cell". In some datasets additional drift was applied to the final data with Matlab. Two different stochastic gating models, (termed M1 and M2) were used to generate semi-synthetic ion channel data. M1 is a low open probability model from ref. 41 (Fig. 3a, b), typically no more than one ion channel opens simultaneously. M2 is from refs. 42,44 and has a much higher open probability (Fig. 3c, d), consequently up to five channels open simultaneously and there are few instances of zero channels open.

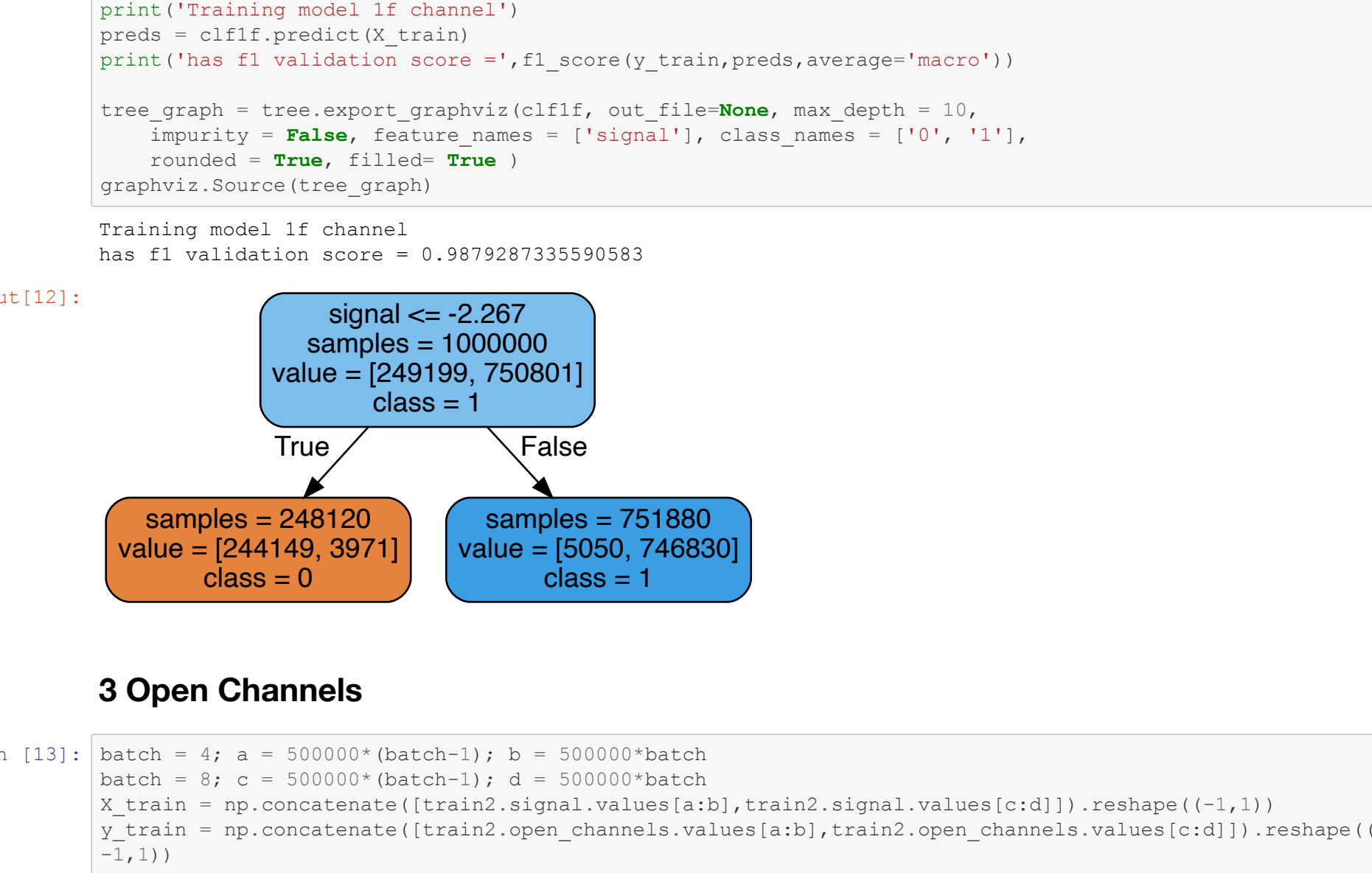
## Correlation Between Signal and Open Channels

Let's look closely at random intervals of signal and open channels to observe how they relate. We notice that they are highly correlated and move up and down together. Therefore we can probably predict open channels from the one feature signal. The only complication is the synthetic drift that was added. So we will remove it.



## Test Data

Let's display the test data signal

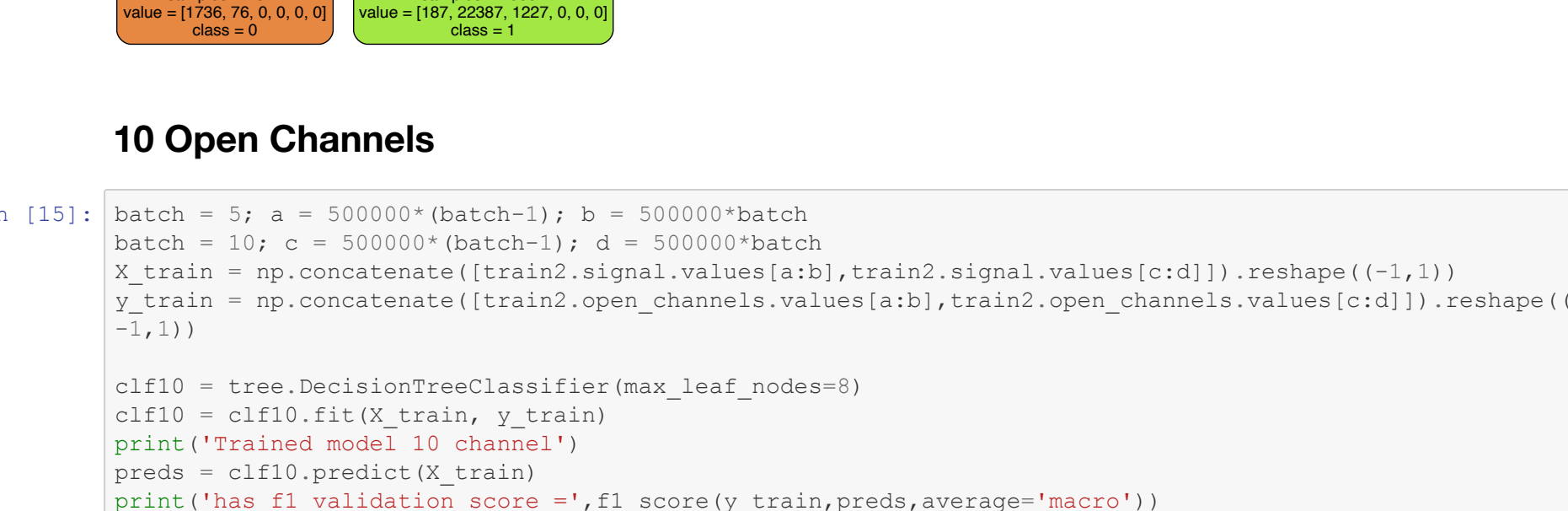
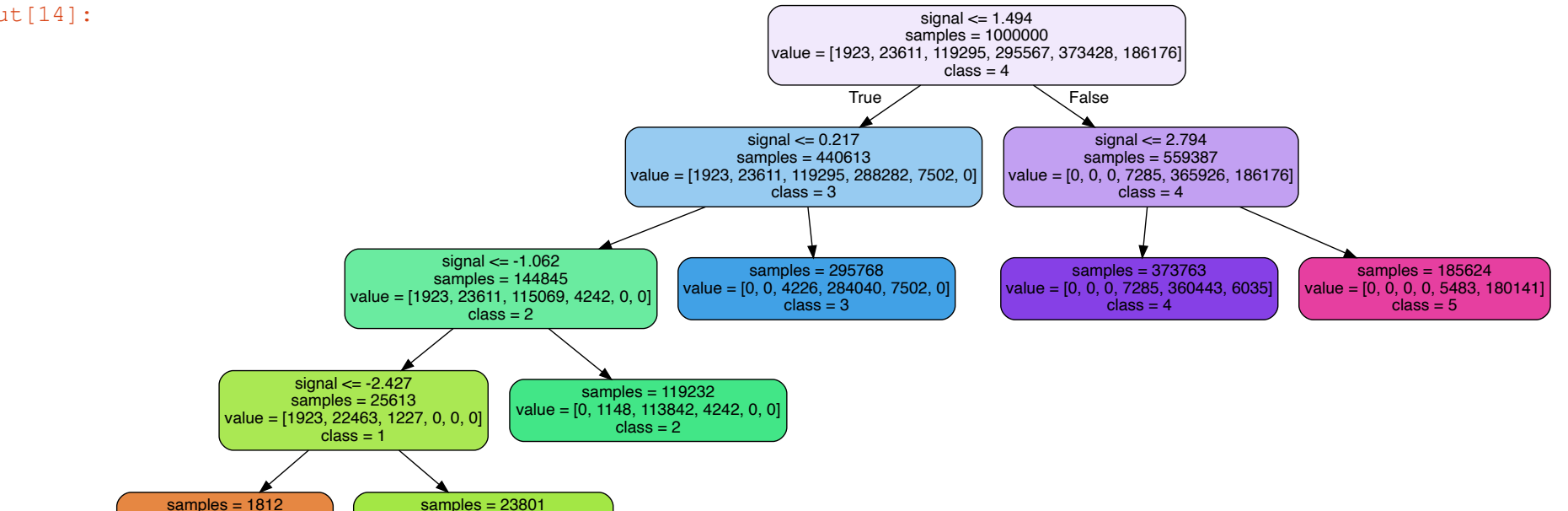


## Reflection

From this plot, we can locate the 5 models in action. And we can recognize the added drift. Batch 1 appears to be 5 subsamples where A, B, C, D, E were created by models 1, 3, 5, 1s, 1r respectively. Model 1s is the model with maximum 1 channel open with low prob. Model 1r is the model with maximum 1 channel open with high prob. And models 3, 5, 10 are models with maximum 3, 5, 10 channels respectively. We observe slant drift in subsamples A, B, E, G, H, I. We observe parabola drift in batch 3.

## Remove Training Data Drift

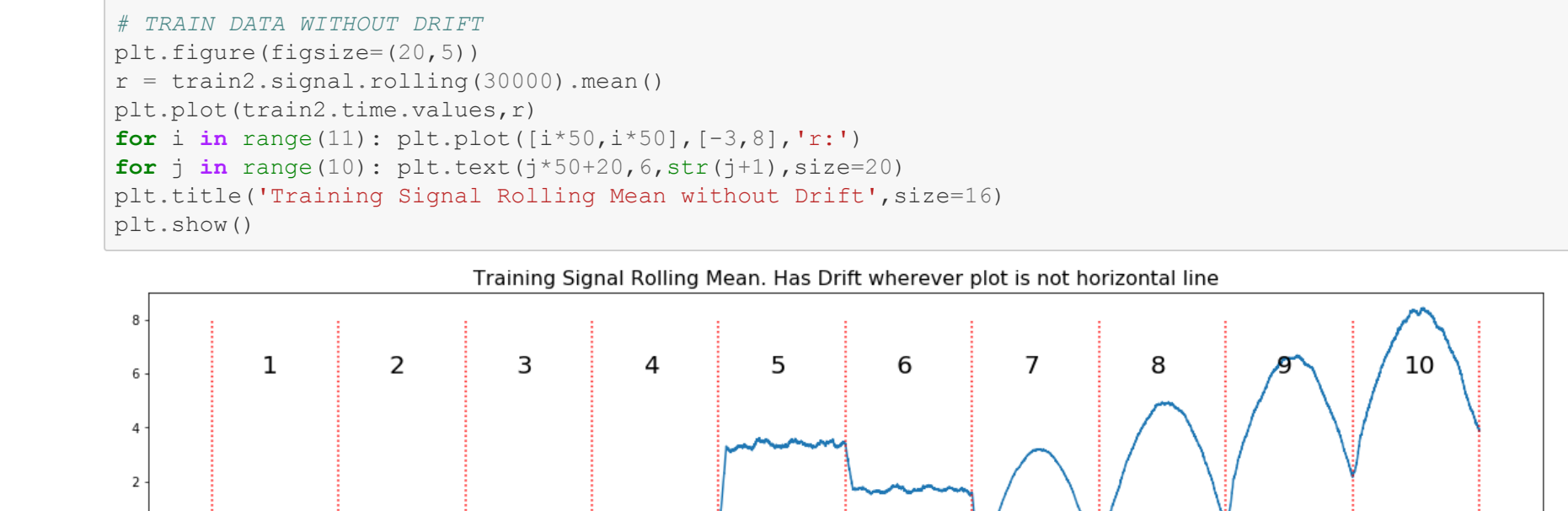
This is a demonstration to show slant drift removal. We could also remove the parabolic drift in batches 7, 8, 9, 10 if we wanted. Below we remove train our models with batches 1, 3, 4, 5, 6. But after removing training drift, we can include the data from batches 2, 7, 8, 9, 10 in our training if we want.



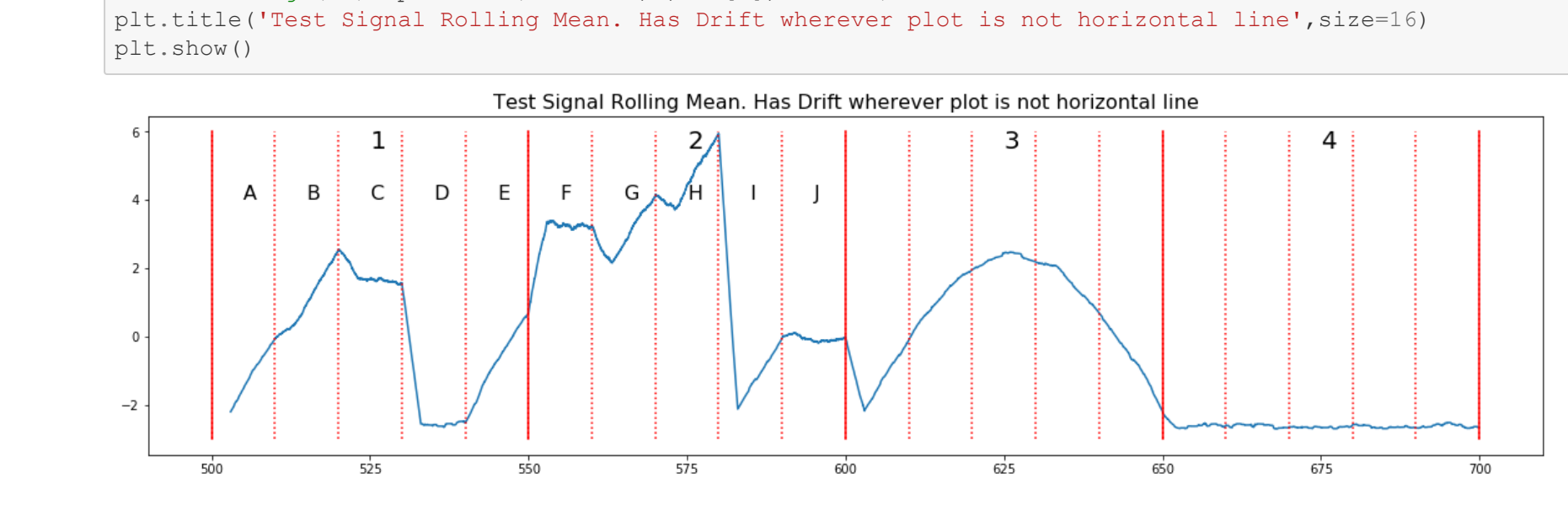
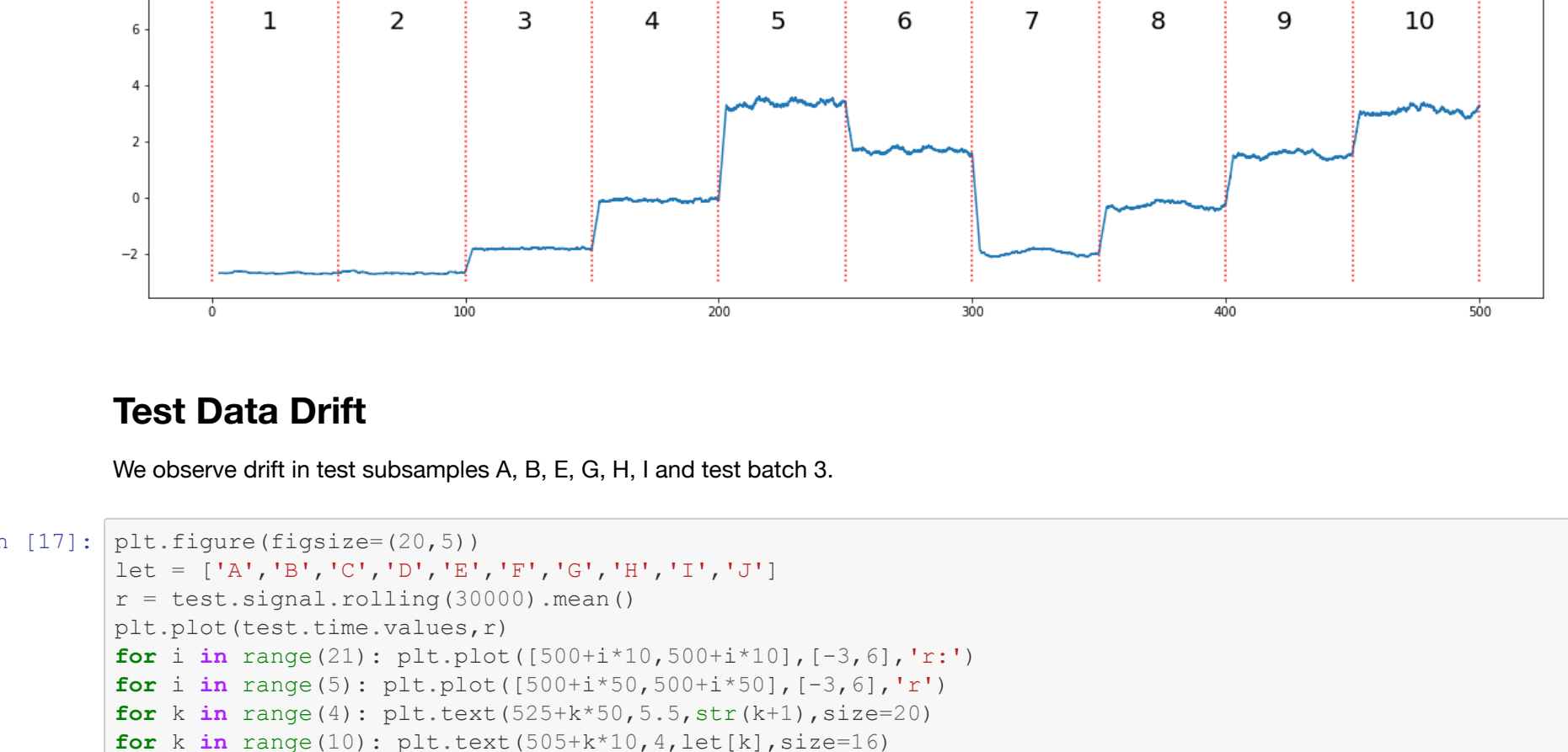
## Make Five Simple Models

We will make one model for each different type of signal we observed above.

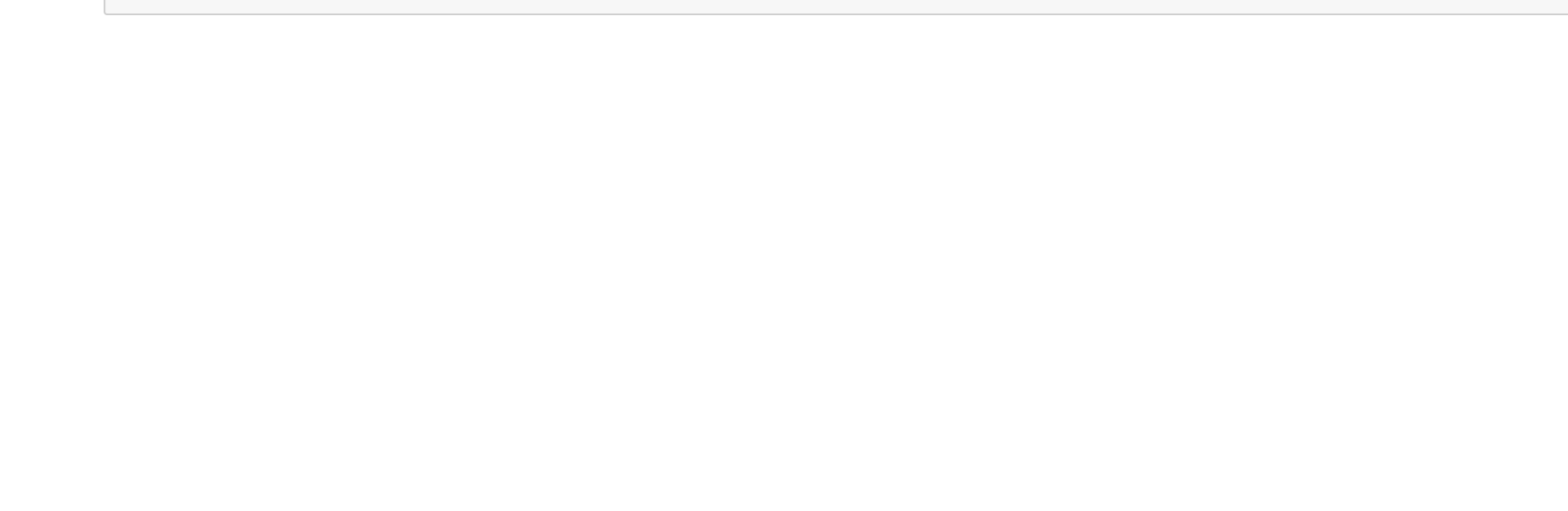
### 1 Slow Open Channel



### 1 Fast Open Channel



### 3 Open Channels



### 5 Open Channels



### 10 Open Channels

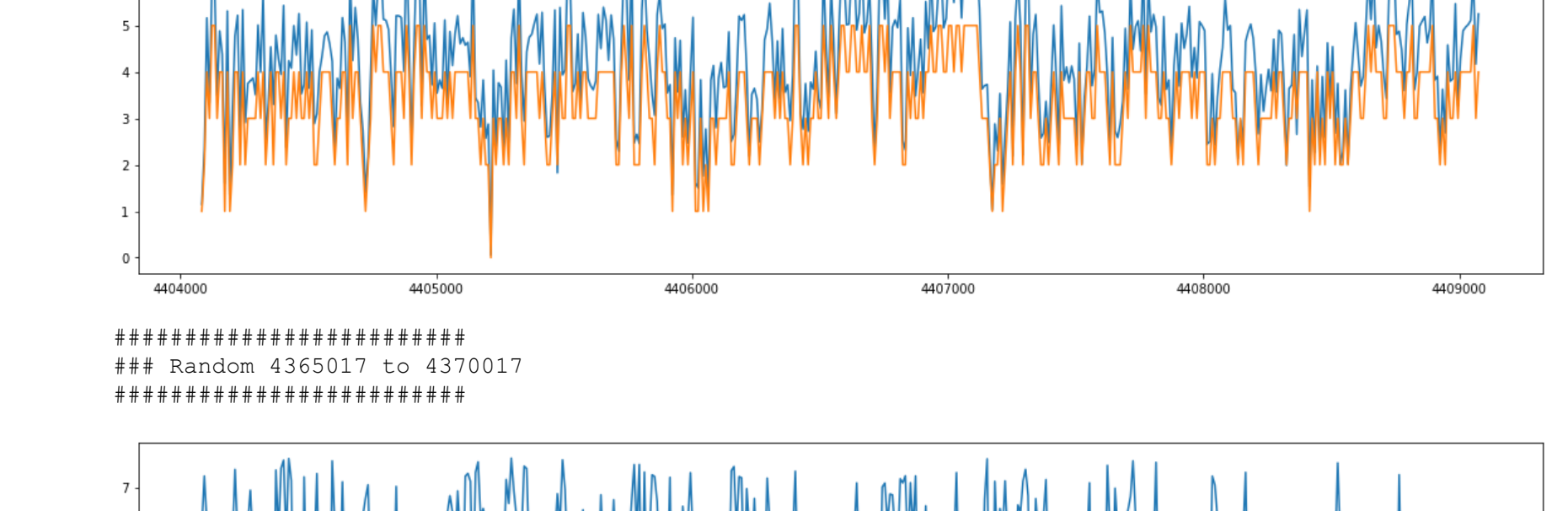
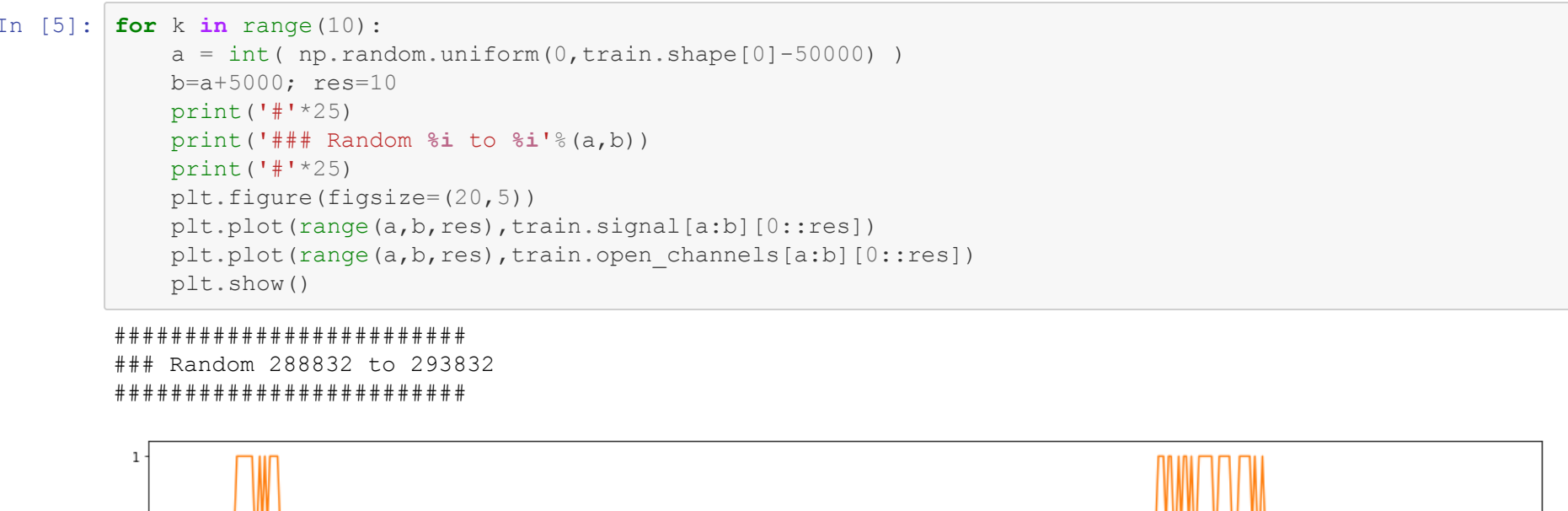
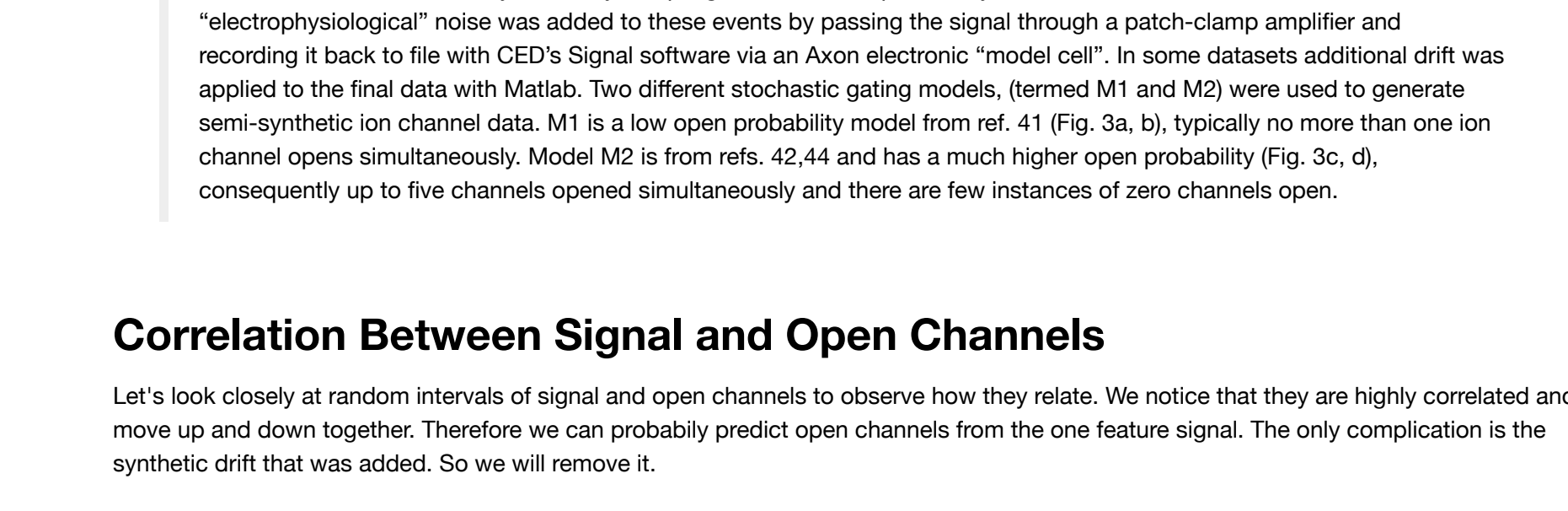


## Analyze Test Data Drift

Let's plot the drift in the training and test data

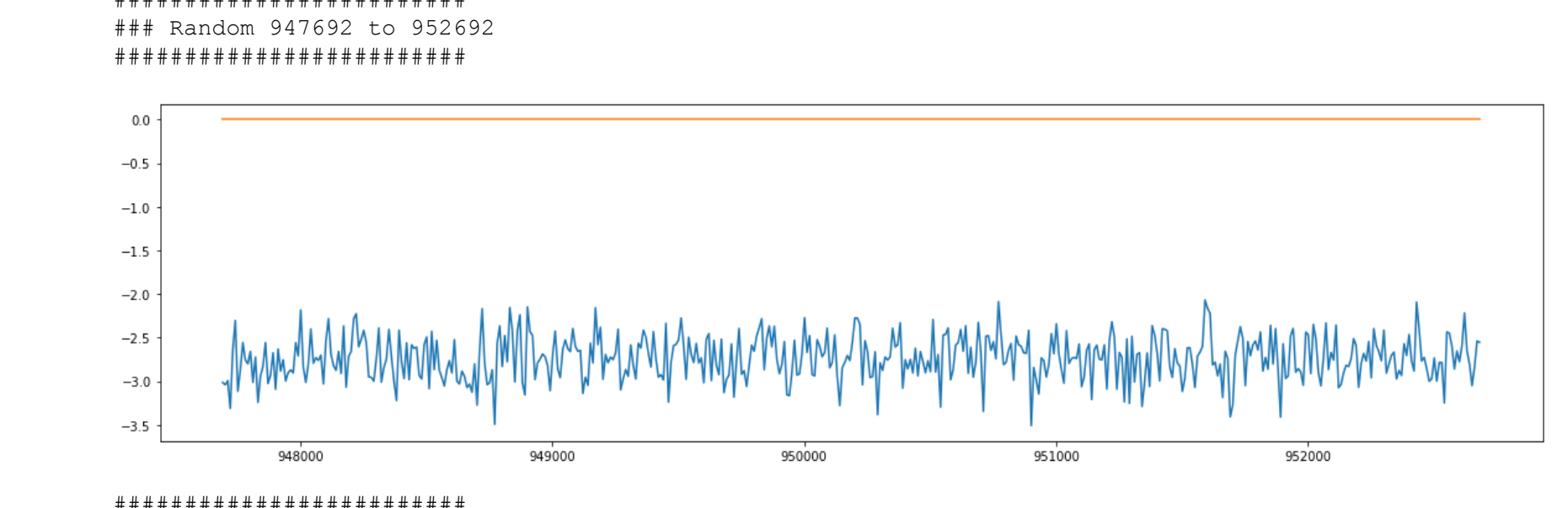
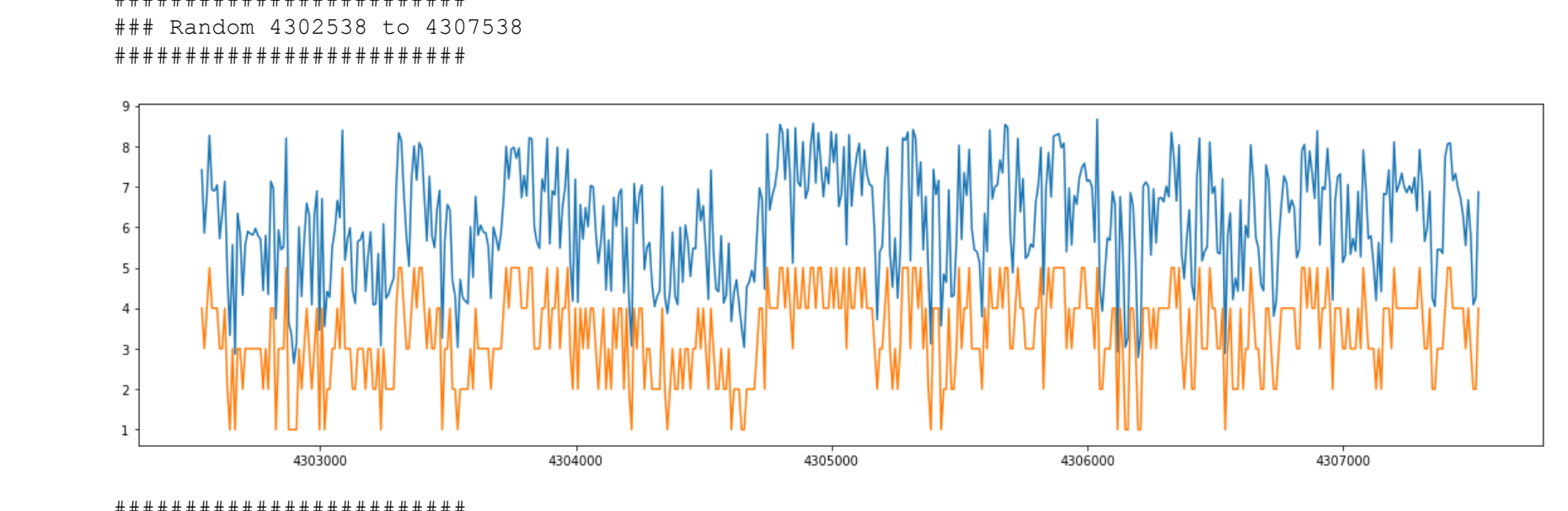
### Training Data Drift

We observe drift wherever the following plot is not a horizontal line. We see drift in batches 2, 7, 8, 9, 10.

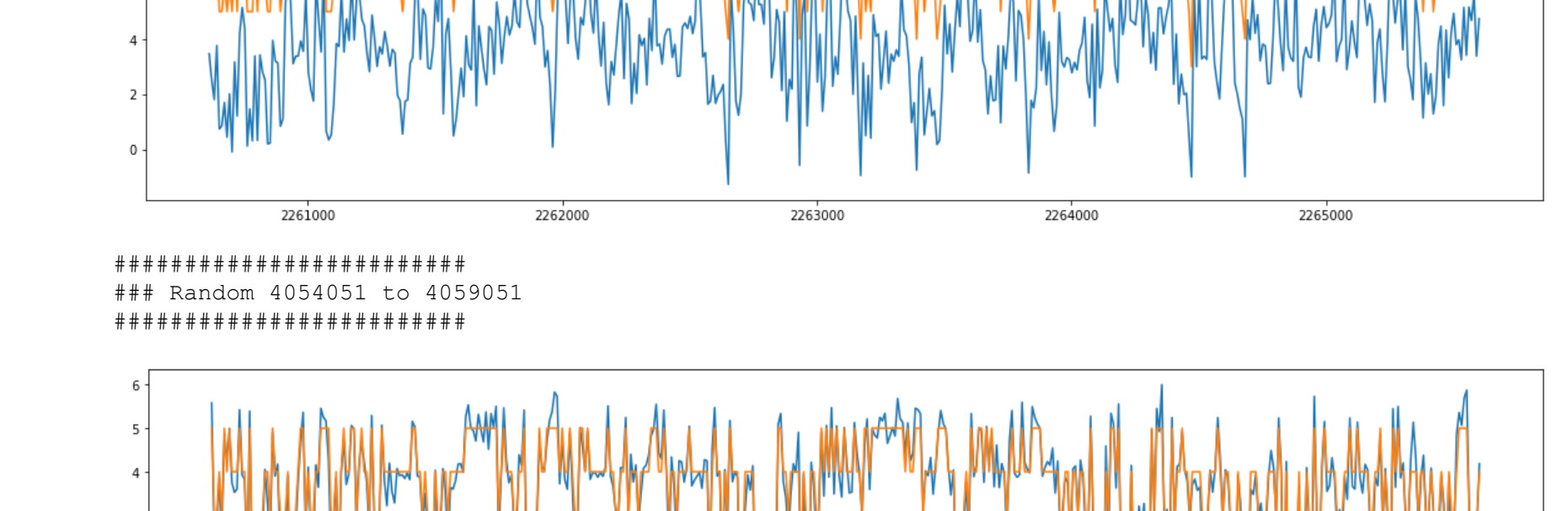


### Test Data Drift

We observe drift in test subsamples A, B, E, G, H, I and test batch 3.



## Remove Test Data Drift





```
[22]: plt.figure(figsize=(20,5))
res = 1000
plt.plot(range(0,test2.shape[0],res),res,test2.signal[0:res])
for i in range(5): plt.plot([(i*500000,i*500000)],[-5,12.5], 'r')
for i in range(21): plt.plot([(i*100000,i*100000)],[-5,12.5], 'r')
for k in range(4): plt.text(k*500000+250000,10, str(k+1),size=20)
for k in range(10): plt.text(k*100000+40000,7.5, str(k),size=16)
plt.title("Test Signal without Drift",size=16)
plt.show()

plt.figure(figsize=(20,5))
r = test2.signal.rolling(30000).mean()
plt.plot(test2.time.values,r)
for i in range(21): plt.plot([500+i*10,500+i*10],[-2,6], 'r')
for i in range(5): plt.plot([500+i*50,500+i*50],[-2,6], 'r')
for k in range(4): plt.text(525+k*50,5.5, str(k+1),size=20)
for k in range(10): plt.text(505+k*10,4, str(k),size=16)
plt.title("Test Signal Rolling Mean without Drift",size=16)
plt.show()
```

In [23]:

```
sub = pd.read_csv('.../input/liverpool-ion-switching/sample_submission.csv')

a = 0 # SUBSAMPLE A, Model 1s
sub.iloc[100000*a:100000*(a+1),1] = clf1s.predict(test2.signal.values[100000*a:100000*(a+1)].reshape((-1,1)))

a = 1 # SUBSAMPLE B, Model 3s
sub.iloc[100000*a:100000*(a+1),1] = clf3s.predict(test2.signal.values[100000*a:100000*(a+1)].reshape((-1,1)))

a = 2 # SUBSAMPLE C, Model 5s
sub.iloc[100000*a:100000*(a+1),1] = clf5s.predict(test2.signal.values[100000*a:100000*(a+1)].reshape((-1,1)))

a = 3 # SUBSAMPLE D, Model 1s
sub.iloc[100000*a:100000*(a+1),1] = clf1s.predict(test2.signal.values[100000*a:100000*(a+1)].reshape((-1,1)))

a = 4 # SUBSAMPLE E, Model 1f
sub.iloc[100000*a:100000*(a+1),1] = clf1f.predict(test2.signal.values[100000*a:100000*(a+1)].reshape((-1,1)))

a = 5 # SUBSAMPLE F, Model 10
sub.iloc[100000*a:100000*(a+1),1] = clf10.predict(test2.signal.values[100000*a:100000*(a+1)].reshape((-1,1)))

a = 6 # SUBSAMPLE G, Model 5
sub.iloc[100000*a:100000*(a+1),1] = clf5s.predict(test2.signal.values[100000*a:100000*(a+1)].reshape((-1,1)))

a = 7 # SUBSAMPLE H, Model 10
sub.iloc[100000*a:100000*(a+1),1] = clf10.predict(test2.signal.values[100000*a:100000*(a+1)].reshape((-1,1)))

a = 8 # SUBSAMPLE I, Model 1s
sub.iloc[100000*a:100000*(a+1),1] = clf1s.predict(test2.signal.values[100000*a:100000*(a+1)].reshape((-1,1)))

a = 9 # SUBSAMPLE J, Model 3
sub.iloc[100000*a:100000*(a+1),1] = clf3s.predict(test2.signal.values[100000*a:100000*(a+1)].reshape((-1,1)))

# BATCHES 3 AND 4, Model 1s
sub.iloc[1000000:2000000,1] = clf1s.predict(test2.signal.values[1000000:2000000]).reshape((-1,1)))
```

## Predict Test

In [24]:

```
plt.figure(figsize=(20,5))
res = 1000
plt.plot(range(0,test.shape[0],res),sub.open_channels[0:res])
for i in range(5): plt.plot([(i*500000,i*500000)],[-5,12.5], 'r')
for i in range(21): plt.plot([(i*100000,i*100000)],[-5,12.5], 'r')
for k in range(4): plt.text(k*500000+250000,10, str(k+1),size=20)
for k in range(10): plt.text(k*100000+40000,7.5, str(k),size=16)
plt.title("Test Data Predictions",size=16)
plt.show()
```

In [25]:

```
sub.to_csv('submission.csv',index=False,float_format='%4f')
```