# Spike Sorting: Tutorial

Christophe Pouzat

LASCON, January 23 2018, Tutorial 3

## Contents

# 1 Downloading the data

The data are available and can be downloaded with (watch out, you must use slightly different commands if you're using `Python 2`):

```python
from urllib.request import urlretrieve # Python 3
# from urllib import urlretrieve # Python 2
data_names = ['Locust_' + str(i) + '.dat.gz' for i in range(1,5)]
data_src = ['http://xtof.disque.math.cnrs.fr/data/' + n
            for n in data_names]
[urlretrieve(data_src[i],data_names[i]) for i in range(4)]
```

They were stored as floats coded on 64 bits and compressed with `gnuzip`. So we decompress it:

```python
import gzip
import shutil
data_snames = ['Locust_' + str(i) + '.dat' for i in range(1,5)]
```

```
4   for in_name, out_name in zip(data_names,data_snames):
5       with gzip.open(in_name,'rb') as f_in:
6           with open(out_name,'wb') as f_out:
7               shutil.copyfileobj(f_in, f_out)
```

20 seconds of data sampled at 15 kHz are contained in these files (see PouzatEtAl_2002 for details). Four files corresponding to the four electrodes or recording sites of a *tetrode* (see Sec. why-tetrode) are used.

## 2 Importing the required modules and loading the data

The individual functions developed for this kind of analysis are defined at the end of this document (Sec. 12). They can also be downloaded as a single file sorting_with_python.py which must then be imported with for instance:

```
1   import sorting_with_python as swp
```

where it is assumed that the working directory of your python session is the directory where the file sorting_with_python.py can be found. We are going to use numpy and pylab (we will also use pandas later on, but to generate only one figure so you can do the analysis without it). We are also going to use the interactive mode of the latter:

```
1   import numpy as np
2   import matplotlib.pylab as plt
3   plt.ion()
```

Python 3 was used to perform this analysis but everything also works with Python 2. We load the data with:

```
1   # Create a list with the file names
2   data_files_names = ['Locust_' + str(i) + '.dat' for i in range(1,5)]
3   # Get the lenght of the data in the files
4   data_len = np.unique(list(map(len, map(lambda n:
5                                   np.fromfile(n,np.double),
6                                   data_files_names))))[0]
7   # Load the data in a list of numpy arrays
8   data = [np.fromfile(n,np.double) for n in data_files_names]
```

## 3 Preliminary analysis

We are going to start our analysis by some "sanity checks" to make sure that nothing "weird" happened during the recording.

### 3.1 Five number summary

We should start by getting an overall picture of the data like the one provided by the mquantiles method of module scipy.stats.mstats using it to output a five-number

summary. The five numbers are the `minimum`, the `first quartile`, the `median`, the `third quartile` and the `maximum`:

```
1  from scipy.stats.mstats import mquantiles
2  np.set_printoptions(precision=3)
3  [mquantiles(x,prob=[0,0.25,0.5,0.75,1]) for x in data]
```

```
[array([-9.074, -0.371, -0.029,  0.326, 10.626]),
 array([-8.229, -0.45 , -0.036,  0.396, 11.742]),
 array([-6.89 , -0.53 , -0.042,  0.469,  9.849]),
 array([-7.347, -0.492, -0.04 ,  0.431, 10.564])]
```

In the above result, each row corresponds to a recording channel, the first column contains the minimal value; the second, the first quartile; the third, the median; the fourth, the third quartile; the fifth, the maximal value. We see that the data range (`maximum - minimum`) is similar (close to 20) on the four recording sites. The inter-quartiles ranges are also similar.

### 3.2 Were the data normalized?

We can check next if some processing like a division by the *standard deviation* (SD) has been applied:

```
1  [np.std(x) for x in data]
```

```
[0.9999983333319417, 0.9999983333319362, 0.9999983333319479, 0.9999983333317428]
```

We see that SD normalization was indeed applied to these data...

### 3.3 Discretization step amplitude

We can easily obtain the size of the digitization set:

```
1  [np.min(np.diff(np.sort(np.unique(x)))) for x in data]
```

```
[0.006709845078411547,
 0.009194500187932775,
 0.011888432902217971,
 0.009614042128660572]
```

## 4 Plot the data

Plotting the data for interactive exploration is trivial. The only trick is to add (or subtract) a proper offest (that we get here using the maximal value of each channel from our five-number summary), this is automatically implemented in our `plot_data_list` function:

```
1  tt = np.arange(0,data_len)/1.5e4
2  swp.plot_data_list(data,tt,0.1)
```

The first channel is drawn as is, the second is offset downward by the sum of its maximal value and of the absolute value of the minimal value of the first, etc. We then get something like Fig. **??**.
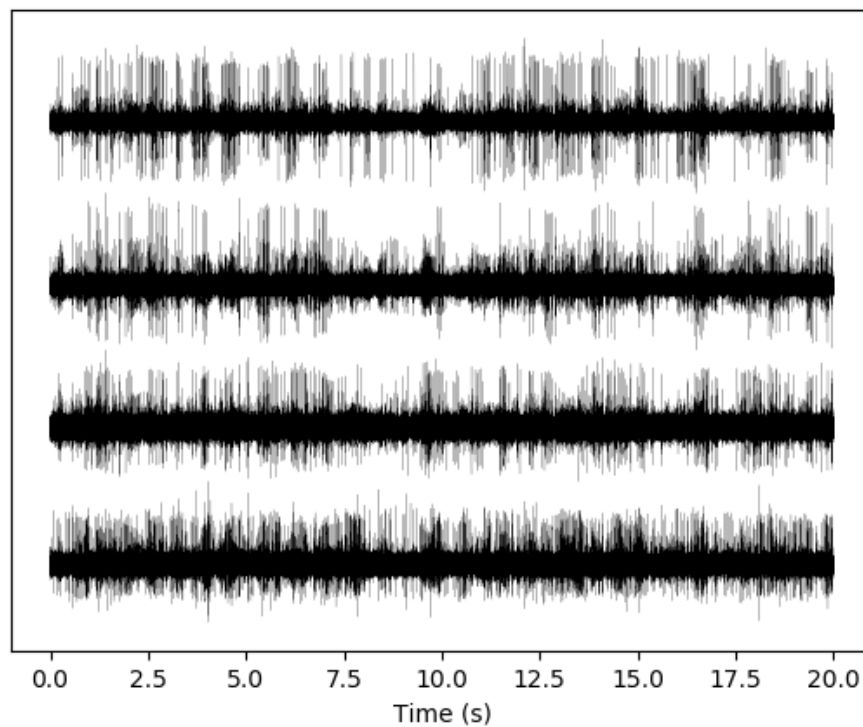


Figure 1: The whole (20 s) Locust antennal lobe data set.

It is also good to "zoom in" and look at the data with a finer time scale (Fig. **??**) with:
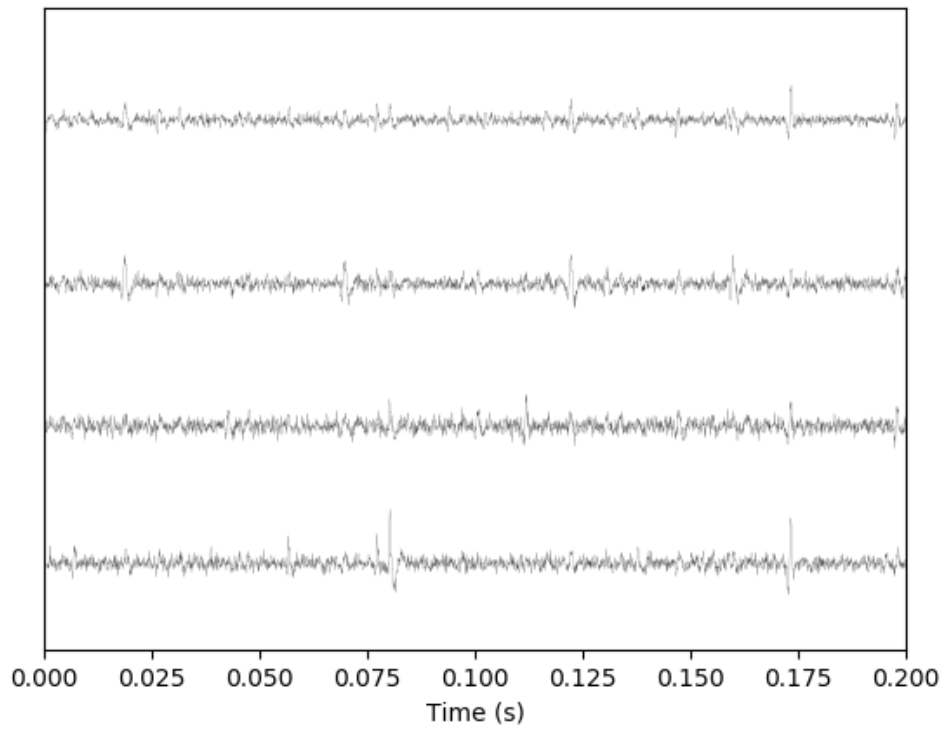
```
1  plt.xlim([0,0.2])
```

Figure 2: First 200 ms of the Locust data set.

# 5 Data renormalization

We are going to use a median absolute deviation (MAD) based renormalization. The goal of the procedure is to scale the raw data such that the *noise SD* is approximately 1. Since it is not straightforward to obtain a noise SD on data where both signal (*i.e.*, spikes) and noise are present, we use this robust type of statistic for the SD:

```
1  data_mad = list(map(swp.mad,data))
2  data_mad
```

```
[0.5172968482892563, 0.6270612350170097, 0.7402832060747951, 0.6841813852777244]
```

And we normalize accordingly (we also subtract the `median` which is not exactly 0):

```
1  data = list(map(lambda x: (x−np.median(x))/swp.mad(x), data))
```

We can check on a plot (Fig. **??**) how `MAD` and `SD` compare:

```
1  plt.plot(tt,data[0],color="black")
2  plt.xlim([0,0.2])
3  plt.ylim([−17,13])
4  plt.axhline(y=1,color="red")
5  plt.axhline(y=−1,color="red")
6  plt.axhline(y=np.std(data[0]),color="blue",linestyle="dashed")
7  plt.axhline(y=−np.std(data[0]),color="blue",linestyle="dashed")
8  plt.xlabel('Time␣(s)')
9  plt.ylim([−5,10])
```

## 5.1 A quick check that the `MAD` **"does its job"**

We can check that the `MAD` does its job as a robust estimate of the *noise* standard deviation by looking at Q-Q plots of the whole traces normalized with the `MAD` and normalized with the "classical" `SD` (Fig. **??**):

```
1  dataQ = map(lambda x:
2               mquantiles(x, prob=np.arange(0.01,0.99,0.001)),data)
3  dataQsd = map(lambda x:
4               mquantiles(x/np.std(x), prob=np.arange(0.01,0.99,0.001)),
5               data)
6  from scipy.stats import norm
7  qq = norm.ppf(np.arange(0.01,0.99,0.001))
8  plt.plot(np.linspace(−3,3,num=100),np.linspace(−3,3,num=100),
9           color='grey')
10 colors = ['black', 'orange', 'blue', 'red']
11 for i,y in enumerate(dataQ):
12     plt.plt.plot(qq,y,color=colors[i])
13
14 for i,y in enumerate(dataQsd):
15     plt.plot(qq,y,color=colors[i],linestyle="dashed")
```
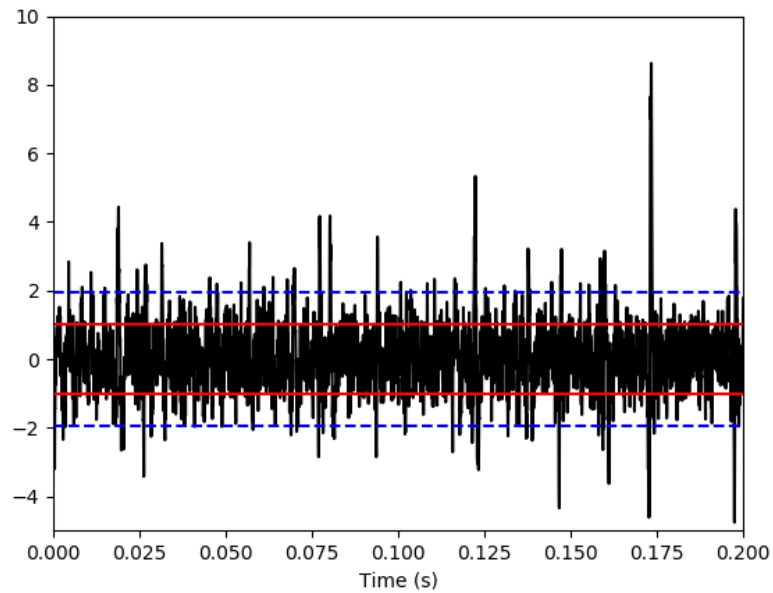
Figure 3: First 200 ms on site 1 of the Locust data set. In red: +/- the `MAD`; in dashed blue +/- the `SD`.

```
16
17  plt.xlabel('Normal␣quantiles')
18  plt.ylabel('Empirical␣quantiles')
```

We see that the behavior of the "away from normal" fraction is much more homogeneous for small, as well as for large in fact, quantile values with the `MAD` normalized traces than with the `SD` normalized ones. If we consider automatic rules like the three sigmas we are going to reject fewer events (*i.e.*, get fewer putative spikes) with the `SD` based normalization than with the `MAD` based one.

## 6 Detect peaks

We are going to filter the data slightly using a "box" filter of length 5. That is, the data points of the original trace are going to be replaced by the average of themselves with their four nearest neighbors. We will then scale the filtered traces such that the `MAD` is one on each recording sites and keep only the parts of the signal above 4:

```
1  from scipy.signal import fftconvolve
2  from numpy import apply_along_axis as apply
3  data_filtered = apply(lambda x:
4                        fftconvolve(x,np.array([1,1,1,1,1])/5.,'same'),
5                        1,np.array(data))
6  data_filtered = (data_filtered.transpose() / \
```
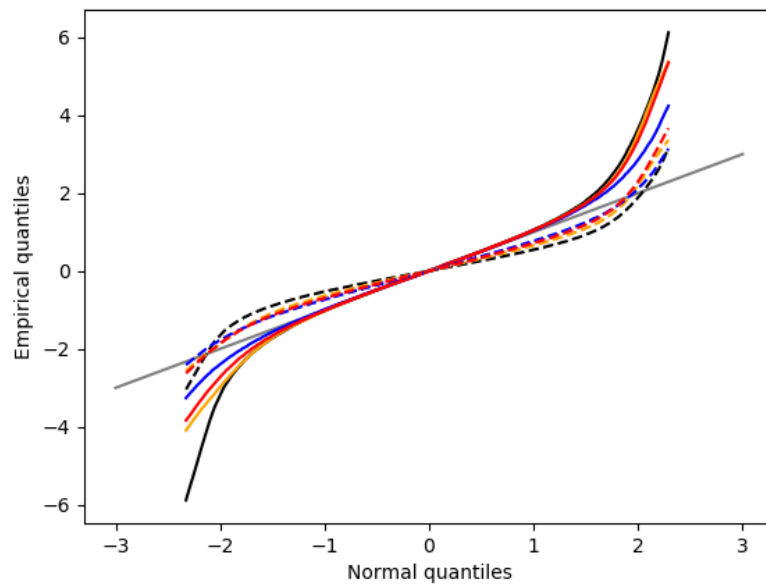
8

Figure 4: Performances of `MAD` based vs `SD` based normalizations. After normalizing the data of each recording site by its `MAD` (plain colored curves) or its `SD` (dashed colored curves), Q-Q plot against a standard normal distribution were constructed. Colors: site 1, black; site 2, orange; site 3, blue; site 4, red.

```
7                          apply(swp.mad,1,data_filtered)).transpose()
8   data_filtered[data_filtered < 4] = 0
```

We can see the difference between the *raw* trace and the *filtered and rectified* one (Fig. **??**) on which spikes are going to be detected with:

```
1   plt.plot(tt, data[0],color='black')
2   plt.axhline(y=4,color="blue",linestyle="dashed")
3   plt.plot(tt, data_filtered[0,],color='red')
4   plt.xlim([0,0.2])
5   plt.ylim([−5,10])
6   plt.xlabel('Time (s)')
```



Figure 5: First 200 ms on site 1 of data set `data`. The raw data are shown in black, the detection threshold appears in dashed blue and the filtered and rectified trace on which spike detection is going to be preformed appears in red.

We now use function `peak` on the sum of the rows of our filtered and rectified version of the data:

```
1   sp0 = swp.peak(data_filtered.sum(0))
```

Giving 1795 spikes, a mean inter-event interval of 167.0 sampling points, a standard deviation of 144.0 sampling points, a smallest inter-event interval of 16 sampling points and a largest of 1157 sampling points.

## 6.1 Interactive spike detection check

We can then check the detection quality with:

```
1  swp.plot_data_list_and_detection(data,tt,sp0)
2  plt.xlim([0,0.2])
```



Figure 6: First 200 ms of data set `data`. The raw data are shown in black, the detected events are signaled by red dots (a dot is put on each recording site at the amplitude on that site at that time).

## 6.2 Split the data set in two parts

As explained in the text, we want to "emulate" a long data set analysis where the model is estimated on the early part before doing template matching on what follows. We therefore get an "early" and a "late" part by splitting the data set in two:

```
1  sp0E = sp0[sp0 <= data_len/2.]
2  sp0L = sp0[sp0 > data_len/2.]
```

In `sp0E`, the number of detected events is: 908 ; the mean inter-event interval is: 165.0; the standard deviation of the inter-event intervals is: 139.0; the smallest inter-event interval is: 16 sampling points long; the largest inter-event interval is: 931 sampling points long.

In `sp0L`, the number of detected events is: 887; the mean inter-event interval is: 169.0; the standard deviation of the inter-event intervals is: 149.0; the smallest inter-event in-

11

terval is: 16 sampling points long; the largest inter-event interval is: 1157 sampling points long.

# 7 Cuts

After detecting our spikes, we must make our cuts in order to create our events' sample. The obvious question we must first address is: How long should our cuts be? The pragmatic way to get an answer is:

- Make cuts much longer than what we think is necessary, like 50 sampling points on both sides of the detected event's time.

- Compute robust estimates of the "central" event (with the `median`) and of the dispersion of the sample around this central event (with the `MAD`).

- Plot the two together and check when does the `MAD` trace reach the background noise level (at 1 since we have normalized the data).

- Having the central event allows us to see if it outlasts significantly the region where the `MAD` is above the background noise level.

Clearly cutting beyond the time at which the `MAD` hits back the noise level should not bring any useful information as far a classifying the spikes is concerned. So here we perform this task as follows:

```
1  evtsE = swp.mk_events(sp0E,np.array(data),49,50)
2  evtsE_median=apply(np.median,0,evtsE)
3  evtsE_mad=apply(swp.mad,0,evtsE)
```

```
1   plt.plot(evtsE_median, color='red', lw=2)
2   plt.axhline(y=0, color='black')
3   for i in np.arange(0,400,100):
4       plt.axvline(x=i, color='black', lw=2)
5
6   for i in np.arange(0,400,10):
7       plt.axvline(x=i, color='grey')
8
9   plt.plot(evtsE_median, color='red', lw=2)
10  plt.plot(evtsE_mad, color='blue', lw=2)
```

Fig. **??** clearly shows that starting the cuts 15 points before the peak and ending them 30 points after should fulfill our goals. We also see that the central event slightly outlasts the window where the `MAD` is larger than 1.

## 7.1 Events

Once we are satisfied with our spike detection, at least in a provisory way, and that we have decided on the length of our cuts, we proceed by making `cuts` around the detected events. :
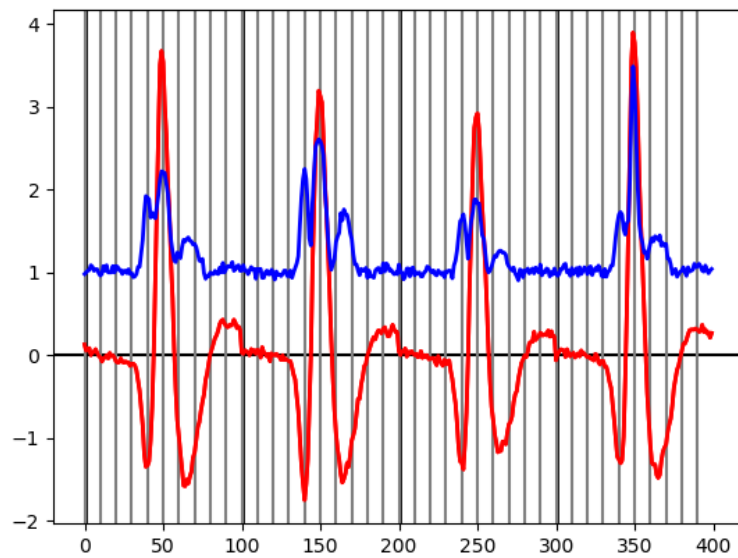
Figure 7: Robust estimates of the central event (black) and of the sample's dispersion around the central event (red) obtained with "long" (100 sampling points) cuts. We see clearly that the dispersion is back to noise level 15 points before the peak and 30 points after the peak.

```
1  evtsE = swp.mk_events(sp0E,np.array(data),14,30)
```

We can visualize the first 200 events with:

```
1  swp.plot_events(evtsE,200)
```



Figure 8:  First 200 events of `evtsE`. Cuts from the four recording sites appear one after the other. The background (white / grey) changes with the site. In red, *robust* estimate of the "central" event obtained by computing the pointwise median. In blue, *robust* estimate of the scale (SD) obtained by computing the pointwise `MAD`.

## 7.2  Noise

Getting an estimate of the noise statistical properties is an essential ingredient to build respectable goodness of fit tests. In our approach "noise events" are essentially anything that is not an "event" is the sense of the previous section. I wrote essentially and not exactly since there is a little twist here which is the minimal distance we are willing to accept between the reference time of a noise event and the reference time of the last preceding and of the first following "event". We could think that keeping a cut length on each side would be enough. That would indeed be the case if *all* events were starting from and returning to zero within a cut. But this is not the case with the cuts parameters we chose previously (that will become clear soon). You might wonder why we chose so short a cut length then. Simply to avoid having to deal with too many superposed events which

14

are the really bothering events for anyone wanting to do proper sorting. To obtain our noise events we are going to use function `mk_noise` which takes the *same* arguments as function `mk_events` plus two numbers:

- `safety_factor` a number by which the cut length is multiplied and which sets the minimal distance between the reference times discussed in the previous paragraph.

- `size` the maximal number of noise events one wants to cut (the actual number obtained might be smaller depending on the data length, the cut length, the safety factor and the number of events).

We cut noise events with a rather large safety factor:

```
1  noiseE = swp.mk_noise(sp0E,np.array(data),14,30,safety_factor=2.5,size=2000)
```

## 7.3 Getting "clean" events

Our spike sorting has two main stages, the first one consist in estimating a **model** and the second one consists in using this model to **classify** the data. Our **model** is going to be built out of reasonably "clean" events. Here by clean we mean events which are not due to a nearly simultaneous firing of two or more neurons; and simultaneity is defined on the time scale of one of our cuts. When the model will be subsequently used to classify data, events are going to decomposed into their (putative) constituent when they are not "clean", that is, **superposition are going to be looked and accounted for**.

In order to eliminate the most obvious superpositions we are going to use a rather brute force approach, looking at the sides of the central peak of our median event and checking if individual events are not too large there, that is do not exhibit extra peaks. We first define a function doing this job:

```
1  def good_evts_fct(samp, thr=3):
2      samp_med = apply(np.median,0,samp)
3      samp_mad = apply(swp.mad,0,samp)
4      above = samp_med > 0
5      samp_r = samp.copy()
6      for i in range(samp.shape[0]): samp_r[i,above] = 0
7      samp_med[above] = 0
8      res = apply(lambda x:
9                  np.all(abs((x-samp_med)/samp_mad) < thr),
10                 1,samp_r)
11     return res
```

We then apply our new function to our sample using a threshold of 8 (set by trial and error):

```
1  goodEvts = good_evts_fct(evtsE,8)
```

Out of 908 events we get 858 "good" ones. As usual, the first 200 good ones can be visualized with:

```
1  swp.plot_events(evtsE[goodEvts,:][:200,:])
```
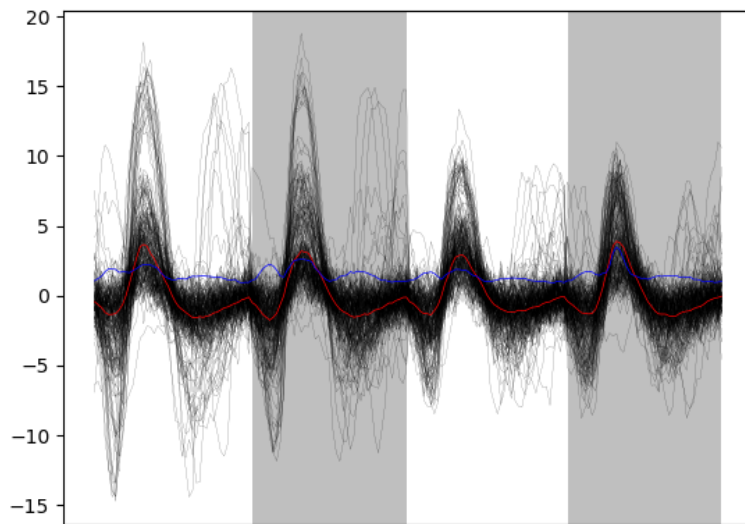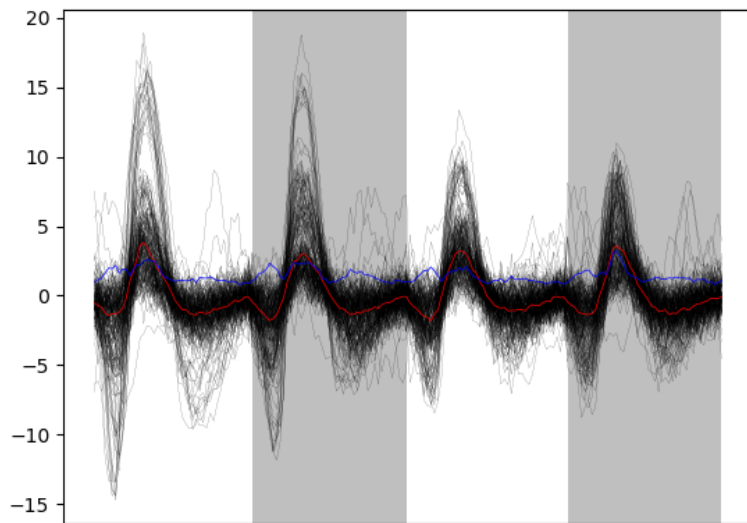
Figure 9: First 200 "good" events of `evtsE`. Cuts from the four recording sites appear one after the other. The background (white / grey) changes with the site. In red, *robust* estimate of the "central" event obtained by computing the pointwise median. In blue, *robust* estimate of the scale (SD) obtained by computing the pointwise `MAD`.

# 8 Dimension reduction

## 8.1 Principal Component Analysis (PCA)

Our events are living right now in an 180 dimensional space (our cuts are 45 sampling points long and we are working with 4 recording sites simultaneously). It turns out that it hard for most humans to perceive structures in such spaces. It also hard, not to say impossible with a realistic sample size, to estimate probability densities (which is what model based clustering algorithms are actually doing) in such spaces, unless one is ready to make strong assumptions about these densities. It is therefore usually a good practice to try to reduce the dimension of the sample space used to represent the data. We are going to that with principal component analysis (PCA), using it on our "good" events.

```
1  from numpy.linalg import svd
2  varcovmat = np.cov(evtsE[goodEvts,:].T)
3  u, s, v = svd(varcovmat)
```

With this "back to the roots" approach, `u` should be an orthonormal matrix whose column are made of the `principal components` (and `v` should be the transpose of `u` since our matrix `varcovmat` is symmetric and real by construction). `s` is a vector containing the amount of sample variance explained by each principal component.

## 8.2 Exploring `PCA` results

PCA is a rather abstract procedure to most of its users, at least when they start using it. But one way to grasp what it does is to plot the `mean event` plus or minus, say five times, each principal components like:

```
1  evt_idx = range(180)
2  evtsE_good_mean = np.mean(evtsE[goodEvts,:],0)
3  for i in range(4):
4      plt.subplot(2,2,i+1)
5      plt.plot(evt_idx,evtsE_good_mean, 'black',evt_idx,
6               evtsE_good_mean + 5 * u[:,i],
7               'red',evt_idx,evtsE_good_mean - 5 * u[:,i], 'blue')
8      plt.title('PC' + str(i) + ': ' + str(round(s[i]/sum(s)*100)) +'%')
```

We can see on Fig. **??** that the first 3 PCs correspond to pure amplitude variations. An event with a large projection (`score`) on the first PC is smaller than the average event on recording sites 1, 2 and 3, but not on 4. An event with a large projection on PC 1 is larger than average on site 1, smaller than average on site 2 and 3 and identical to the average on site 4. An event with a large projection on PC 2 is larger than the average on site 4 only. PC 3 is the first principal component corresponding to a change in *shape* as opposed to *amplitude*. A large projection on PC 3 means that the event as a shallower first valley and a deeper second valley than the average event on all recording sites.

We now look at the next 4 principal components:

```
1  for i in range(4,8):
```

Figure 10: PCA of `evtsE` (for "good" events) exploration (PC 1 to 4). Each of the 4 graphs shows the mean waveform (black), the mean waveform + 5 x PC (red), the mean - 5 x PC (blue) for each of the first 4 PCs. The fraction of the total variance "explained" by the component appears in the title of each graph.

```
2        plt.subplot(2,2,i−3)
3        plt.plot(evt_idx,evtsE_good_mean, 'black',
4                 evt_idx,evtsE_good_mean + 5 * u[:,i], 'red',
5                 evt_idx,evtsE_good_mean − 5 * u[:,i], 'blue')
6        plt.title('PC' + str(i) + ': ' + str(round(s[i]/sum(s)*100)) +'%')
```
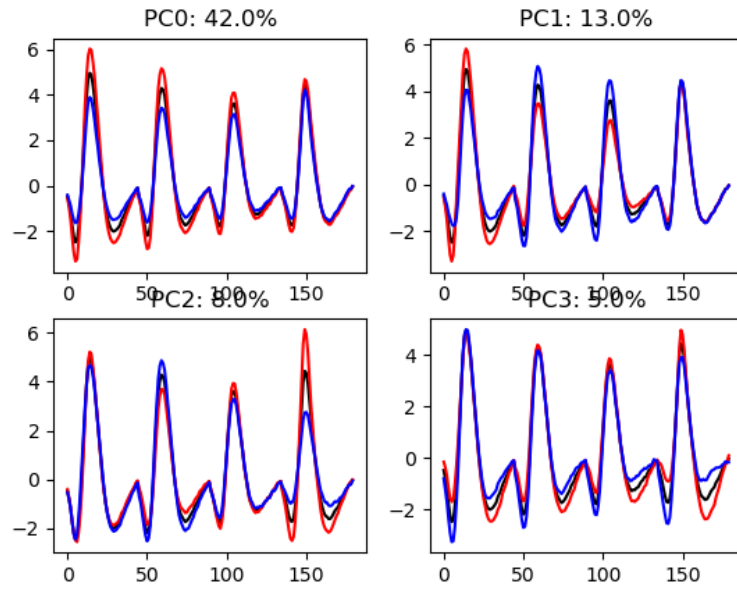


Figure 11: PCA of `evtsE` (for "good" events) exploration (PC 4 to 7). Each of the 4 graphs shows the mean waveform (black), the mean waveform + 5 x PC (red), the mean - 5 x PC (blue). The fraction of the total variance "explained" by the component appears in between parenthesis in the title of each graph.

An event with a large projection on PC 4 (Fig. **??**) tends to be "slower" than the average event. An event with a large projection on PC 5 exhibits a slower kinetics of its second valley than the average event. PC 4 and 5 correspond to effects shared among recording sites. PC 6 correspond also to a "change of shape" effect on all sites except the first. Events with a large projection on PC 7 rise slightly faster and decay slightly slower than the average event on all recording site. Notice also that PC 7 has a "noisier" aspect than the other suggesting that we are reaching the limit of the "events extra variability" compared to the variability present in the background noise.

## 8.3 Static representation of the projected data

We can build a `scatter plot matrix` showing the projections of our "good" events sample onto the plane defined by pairs of the few first PCs:

```
1  evtsE_good_P0_to_P3 = np.dot(evtsE[goodEvts,:],u[:,0:4])
2  from pandas.tools.plotting import scatter_matrix
3  import pandas as pd
4  df = pd.DataFrame(evtsE_good_P0_to_P3)
5  scatter_matrix(df,alpha=0.2,s=4,c='k',figsize=(6,6),
6                 diagonal='kde',marker=".")
```



Figure 12: Scatter plot matrix of the projections of the good events in `evtsE` onto the planes defined by the first 4 PCs. The diagonal shows a smooth (Gaussian kernel based) density estimate of the projection of the sample on the corresponding PC. Using the first 8 PCs does not make finner structure visible.

## 8.4 Dynamic visualization of the data with `GGobi`

The best way to discern structures in "high dimensional" data is to dynamically visualize them. To this end, the tool of choice is GGobi, an open source software available on `Linux`, `Windows` and `MacOS`. We start by exporting our data in `csv` format to our disk:

```
1  import csv
2  f = open('evtsE.csv','w')
3  w = csv.writer(f)
4  w.writerows(np.dot(evtsE[goodEvts,:],u[:,:8]))
5  f.close()
```

The following terse procedure should allow the reader to get going with `GGobi`:

- Launch `GGobi`

- In menu: `File -> Open`, select `evtsE.csv`.

- Since the glyphs are rather large, start by changing them for smaller ones:
    - Go to menu: `Interaction -> Brush`.
    - On the Brush panel which appeared check the `Persistent` box.
    - Click on `Choose color & glyph...`.
    - On the chooser which pops out, click on the small dot on the upper left of the left panel.
    - Go back to the window with the data points.
    - Right click on the lower right corner of the rectangle which appeared on the figure after you selected `Brush`.
    - Dragg the rectangle corner in order to cover the whole set of points.
    - Go back to the `Interaction` menu and select the first row to go back where you were at the start.

- Select menu: `View -> Rotation`.

- Adjust the speed of the rotation in order to see things properly.

We easily discern 10 rather well separated clusters. Meaning that an automatic clustering with 10 clusters on the first 3 principal components should do the job.

## 9 Clustering with K-Means

Since our dynamic visualization shows 10 well separated clusters in 3 dimension, a simple k-means should do the job. We are using here the KMeans class of scikit-learn:

```
1  from sklearn.cluster import KMeans
2  km10 = KMeans(n_clusters=10, init='k-means++', n_init=100, max_iter=100)
3  km10.fit(np.dot(evtsE[goodEvts,:],u[:,0:3]))
4  c10 = km10.fit_predict(np.dot(evtsE[goodEvts,:],u[:,0:3]))
```

In order to facilitate comparison when models with different numbers of clusters or when different models are used, clusters are sorted by "size". The size is defined here as the sum of the absolute value of the median of the cluster (an L1 norm):

```
1  cluster_median = list ([( i ,
2                           np.apply_along_axis(np.median ,0 ,
3                                                evtsE[goodEvts ,:][c10 == i ,:]))
4                                               for i in range(10)
5                                               if sum(c10 == i) > 0])
6  cluster_size = list ([np.sum(np.abs(x[1])) for x in cluster_median])
7  new_order = list(reversed(np.argsort(cluster_size)))
8  new_order_reverse = sorted(range(len(new_order)), key=new_order.__getitem__)
9  c10b = [new_order_reverse[i] for i in c10]
```

## 9.1 Cluster specific plots

Looking at the first 5 clusters we get Fig. **??** with:

```
1  plt.subplot(511)
2  swp.plot_events(evtsE[goodEvts ,:][np.array(c10b) == 0 ,:])
3  plt.ylim([−15,20])
4  plt.subplot(512)
5  swp.plot_events(evtsE[goodEvts ,:][np.array(c10b) == 1 ,:])
6  plt.ylim([−15,20])
7  plt.subplot(513)
8  swp.plot_events(evtsE[goodEvts ,:][np.array(c10b) == 2 ,:])
9  plt.ylim([−15,20])
10 plt.subplot(514)
11 swp.plot_events(evtsE[goodEvts ,:][np.array(c10b) == 3 ,:])
12 plt.ylim([−15,20])
13 plt.subplot(515)
14 swp.plot_events(evtsE[goodEvts ,:][np.array(c10b) == 4 ,:])
15 plt.ylim([−15,20])
```

Looking at the last 5 clusters we get Fig. **??** with:

```
1  plt.subplot(511)
2  swp.plot_events(evtsE[goodEvts ,:][np.array(c10b) == 5 ,:])
3  plt.ylim([−10,10])
4  plt.subplot(512)
5  swp.plot_events(evtsE[goodEvts ,:][np.array(c10b) == 6 ,:])
6  plt.ylim([−10,10])
7  plt.subplot(513)
8  swp.plot_events(evtsE[goodEvts ,:][np.array(c10b) == 7 ,:])
9  plt.ylim([−10,10])
10 plt.subplot(514)
11 swp.plot_events(evtsE[goodEvts ,:][np.array(c10b) == 8 ,:])
12 plt.ylim([−10,10])
13 plt.subplot(515)
14 swp.plot_events(evtsE[goodEvts ,:][np.array(c10b) == 9 ,:])
15 plt.ylim([−10,10])
```

Figure 13: First 5 clusters. Cluster 0 at the top, cluster 4 at the bottom. Red, cluster specific central / median event. Blue, cluster specific `MAD`.

## 9.2 Results inspection with `GGobi`

We start by checking our clustering quality with `GGobi`. To this end we export the data and the labels of each event:

```
1  f = open('evtsEsorted.csv','w')
2  w = csv.writer(f)
3  w.writerows(np.concatenate((np.dot(evtsE[goodEvts,:],u[:,:8]),
4                              np.array([c10b]).T),
5                              axis=1))
6  f.close()
```

An again succinct description of how to do the dynamical visual check is:

- Load the new data into GGobi like before.

- In menu: `Display -> New Scatterplot Display`, select `evtsEsorted.csv`.

- Change the glyphs like before.

- In menu: `Tools -> Color Schemes`, select a scheme with 10 colors, like `Spectral`, `Spectral 10`.

- In menu: `Tools -> Automatic Brushing`, select `evtsEsorted.csv` tab and, within this tab, select variable `c10b`. Then click on `Apply`.
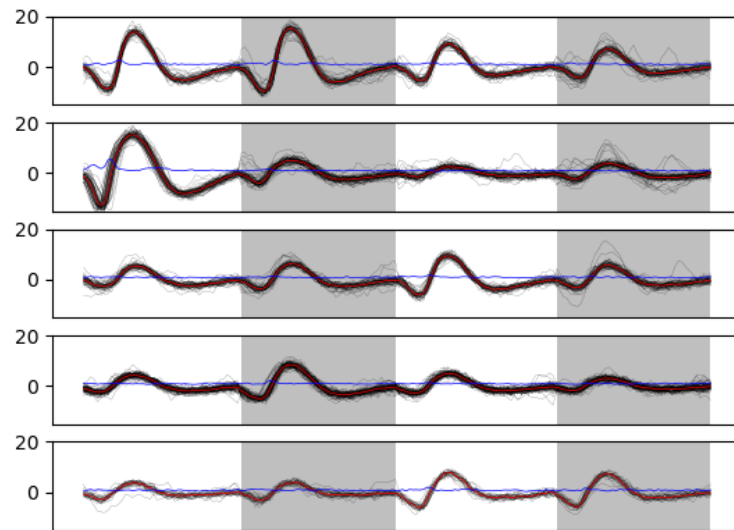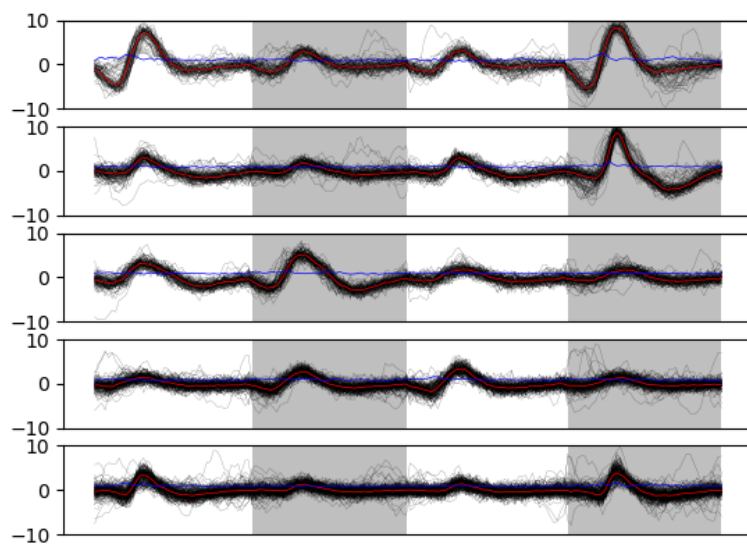
Figure 14: Last 5 clusters. Cluster 5 at the top, cluster 9 at the bottom. Red, cluster specific central / median event. Blue, cluster specific MAD. Notice the change in ordinate scale compared to the previous figure.

- Select `View -> Rotation` like before and see your result.

## 10 Spike "peeling": a "Brute force" superposition resolution

We are going to resolve (the most "obvious") superpositions by a "recursive peeling method":

1. Events are detected and cut from the raw data *or from an already peeled version of the data*.

2. The closest center (in term of Euclidean distance) to the event is found.

3. If the residual sum of squares (`RSS`), that is: (actual data - best center)$^2$, is smaller than the squared norm of a cut, the best center is subtracted from the data on which detection was performed—jitter is again compensated for at this stage.

4. Go back to step 1 or stop.

To apply this procedure, we need, for each cluster, estimates of its center and of its first two derivatives. Function `mk_center_dictionary` does the job for us. We must moreover build our clusters' centers such that they can be used for subtraction, *this implies that we should make them long enough, on both side of the peak, to see them go back to baseline*. Formal parameters `before` and `after` bellow should therefore be set to larger values than the ones used for clustering:

```
1  centers = { "Cluster " + str(i) :
2              swp.mk_center_dictionary(sp0E[goodEvts][np.array(c10b)==i],
3                                        np.array(data))
4              for i in range(10)}
```

### 10.1 First peeling

Function `classify_and_align_evt` is used next. For each detected event, it matches the closest template, correcting for the jitter, if the closest template is close enough:

```
1  swp.classify_and_align_evt(sp0[0],np.array(data),centers)
```

```
['Cluster 7', 281, -0.14107833394834743]
```

We can use the function on every detected event. A trick here is to store the matrix version of the data in order to avoid the conversion of the list of vectors (making the data of the different channels) into a matrix for each detected event:

```
1  data0 = np.array(data)
2  round0 = [swp.classify_and_align_evt(sp0[i],data0,centers)
3            for i in range(len(sp0))]
```

We can check how many events got unclassified on a total of 1795 :

```
1  len([x[1] for x in round0 if x[0] == '?'])
```

22

Using function `predict_data`, we create an ideal data trace given events' positions, events' origins and a clusters' catalog:

```
1  pred0 = swp.predict_data(round0,centers)
```

We then subtract the prediction (`pred0`) from the data (`data0`) to get the "peeled" data (`data1`):

```
1  data1 = data0 - pred0
```

We can compare the original data with the result of the "first peeling" to get Fig. **??**:

```
1   plt.plot(tt, data0[0,], color='black')
2   plt.plot(tt, data1[0,], color='red',lw=0.3)
3   plt.plot(tt, data0[1,]-15, color='black')
4   plt.plot(tt, data1[1,]-15, color='red',lw=0.3)
5   plt.plot(tt, data0[2,]-25, color='black')
6   plt.plot(tt, data1[2,]-25, color='red',lw=0.3)
7   plt.plot(tt, data0[3,]-40, color='black')
8   plt.plot(tt, data1[3,]-40, color='red',lw=0.3)
9   plt.xlabel('Time␣(s)')
10  plt.xlim([0.9,1])
```

## 10.2 Second peeling

We then take `data1` as our former `data0` and we repeat the procedure. We do it with slight modifications: detection is done on a single recording site and a shorter filter length is used before detecting the events. Doing detection on a single site (here site 0) allows us to correct some drawbacks of our crude spike detection method. When we used it the first time we summed the filtered and rectified versions of the data before looking at peaks. This summation can lead to badly defined spike times when two neurons that are large on different recording sites, say site 0 and site 1 fire at nearly the same time. The summed event can then have a peak in between the two true peaks and our jitter correction cannot resolve that. We are therefore going to perform detection on the different sites. The jitter estimation and the subtraction are always going to be done on the 4 recording sites:

```
1  data_filtered = np.apply_along_axis(lambda x:
2                                     fftconvolve(x,np.array([1,1,1])/3.,
3                                                'same'),
4                                     1,data1)
5  data_filtered = (data_filtered.transpose() /
6                    np.apply_along_axis(swp.mad,1,
7                                        data_filtered)).transpose()
```

26

Figure 15:   100 ms of the locust data set.  Black, original data; red, after first peeling.

```
8    data_filtered[data_filtered < 4] = 0
9    sp1 = swp.peak(data_filtered[0,:])
```

We classify the events and obtain the new prediction and the new "data":

```
1    round1 = [swp.classify_and_align_evt(sp1[i],data1,centers)
2              for i in range(len(sp1))]
3    pred1 = swp.predict_data(round1,centers)
4    data2 = data1 − pred1
```

We can check how many events got unclassified on a total of 244:

```
1    len([x[1] for x in round1 if x[0] == '?'])
```

We can compare the first peeling with the second one (Fig. **??**):

```
1    plt.plot(tt, data1[0,], color='black')
2    plt.plot(tt, data2[0,], color='red',lw=0.3)
3    plt.plot(tt, data1[1,]−15, color='black')
4    plt.plot(tt, data2[1,]−15, color='red',lw=0.3)
5    plt.plot(tt, data1[2,]−25, color='black')
6    plt.plot(tt, data2[2,]−25, color='red',lw=0.3)
7    plt.plot(tt, data1[3,]−40, color='black')
```

```
8  plt.plot(tt, data2[3,]−40, color='red',lw=0.3)
9  plt.xlabel('Time␣(s)')
10 plt.xlim([0.9,1])
```



Figure 16: 100 ms of the locust data set. Black, first peeling; red, second peeling.

## 10.3 Third peeling

We take `data2` as our former `data1` and we repeat the procedure detecting on channel 1:

```
1  data_filtered = apply(lambda x:
2                        fftconvolve(x,np.array([1,1,1])/3.,'same'),
3                        1,data2)
4  data_filtered = (data_filtered.transpose() / \
5                  apply(swp.mad,1,data_filtered)).transpose()
6  data_filtered[data_filtered < 4] = 0
7  sp2 = swp.peak(data_filtered[1,:])
8  len(sp2)
```
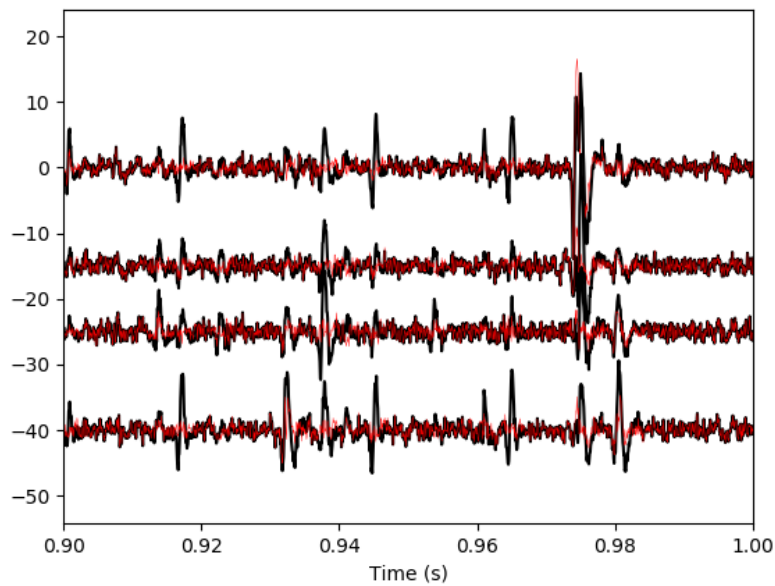
129

The classification follows with the prediction and the number of unclassified events:

```
1  round2 = [swp.classify_and_align_evt(sp2[i],data2,centers) for i in range(len(sp2))]
2  pred2 = swp.predict_data(round2,centers)
```

28

```
3  data3 = data2 − pred2
4  len ([ x [1]  for  x  in  round2  if  x [0] == '?'])
```

We can compare the second peeling with the third one (Fig. **??**):

```
1   plt.plot(tt, data2[0,], color='black')
2   plt.plot(tt, data3[0,], color='red',lw=0.3)
3   plt.plot(tt, data2[1,]−15, color='black')
4   plt.plot(tt, data3[1,]−15, color='red',lw=0.3)
5   plt.plot(tt, data2[2,]−25, color='black')
6   plt.plot(tt, data3[2,]−25, color='red',lw=0.3)
7   plt.plot(tt, data2[3,]−40, color='black')
8   plt.plot(tt, data3[3,]−40, color='red',lw=0.3)
9   plt.xlabel('Time␣(s)')
10  plt.xlim([0.9,1])
```



Figure 17:  100 ms of the locust data set.  Black, second peeling; red, third peeling. *In this portion of data we see events but none belonging to our centers catalog.*

## 10.4  Fourth peeling

We take `data3` as our former `data2` and we repeat the procedure detecting on channel 2:

29

```
1  data_filtered = apply(lambda x:
2                        fftconvolve(x,np.array([1,1,1])/3.,'same'),
3                        1,data3)
4  data_filtered = (data_filtered.transpose() / \
5                   apply(swp.mad,1,data_filtered)).transpose()
6  data_filtered[data_filtered < 4] = 0
7  sp3 = swp.peak(data_filtered[2,:])
8  len(sp3)
```

99

The classification follows with the prediction and the number of unclassified events:

```
1  round3 = [swp.classify_and_align_evt(sp3[i],data3,centers) for i in range(len(sp3))]
2  pred3 = swp.predict_data(round3,centers)
3  data4 = data3 - pred3
4  len([x[1] for x in round3 if x[0] == '?'])
```

16

We can compare the third peeling with the fourth one (Fig. **??**) looking at a different part of the data than on the previous figures:

```
1  plt.plot(tt, data3[0,], color='black')
2  plt.plot(tt, data4[0,], color='red',lw=0.3)
3  plt.plot(tt, data3[1,]-15, color='black')
4  plt.plot(tt, data4[1,]-15, color='red',lw=0.3)
5  plt.plot(tt, data3[2,]-25, color='black')
6  plt.plot(tt, data4[2,]-25, color='red',lw=0.3)
7  plt.plot(tt, data3[3,]-40, color='black')
8  plt.plot(tt, data4[3,]-40, color='red',lw=0.3)
9  plt.xlabel('Time␣(s)')
10 plt.xlim([3.9,4])
```
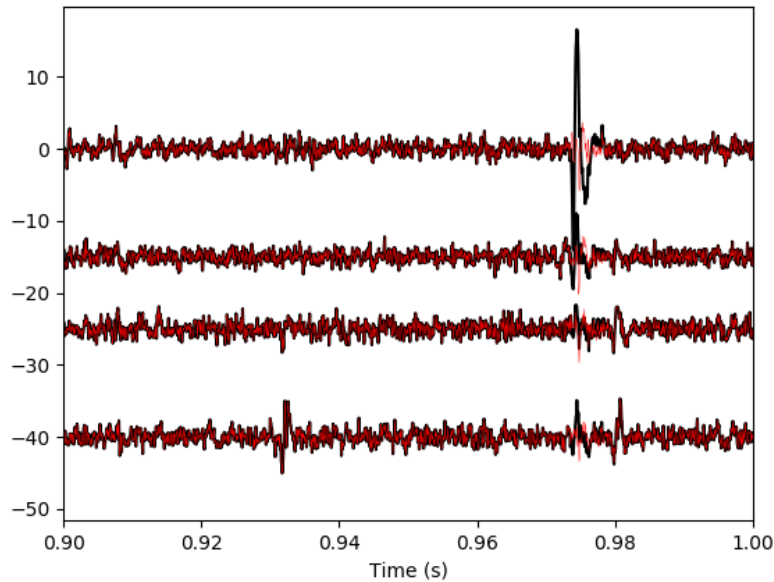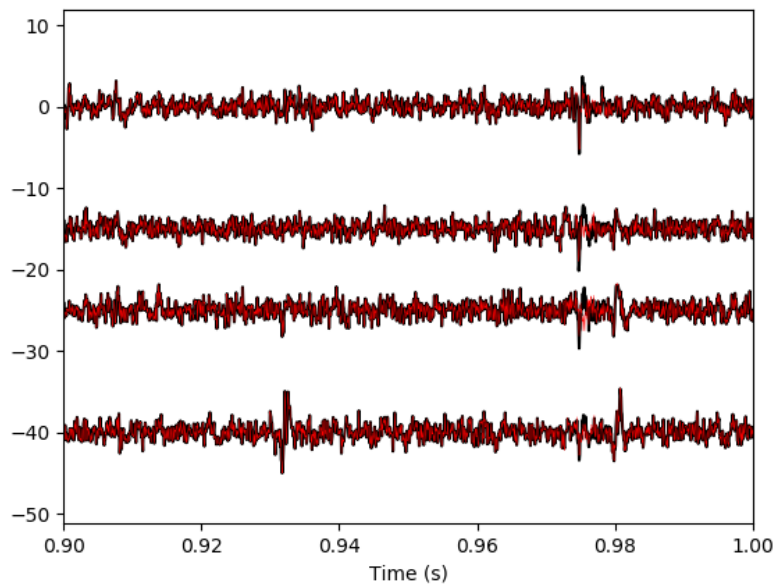
## 10.5 Fifth peeling

We take `data4` as our former `data3` and we repeat the procedure detecting on channel 3:

```
1  data_filtered = apply(lambda x:
2                        fftconvolve(x,np.array([1,1,1])/3.,'same'),
3                        1,data4)
4  data_filtered = (data_filtered.transpose() / \
5                   apply(swp.mad,1,data_filtered)).transpose()
6  data_filtered[data_filtered < 4] = 0
7  sp4 = swp.peak(data_filtered[3,:])
8  len(sp4)
```
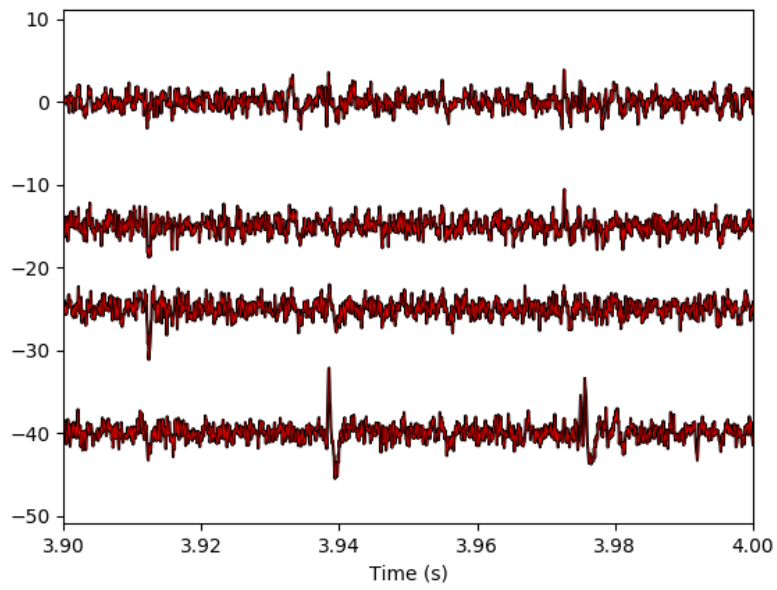
Figure 18: 100 ms of the locust data set (different time frame than on the previous plot). Black, third peeling; red, fourth peeling. *On this portion of the trace, nothing was detected on site 2 (the third one, remember that* `Python` *starts numbering at 0).*

The classification follows with the prediction and the number of unclassified events:

```
1  round4 = [swp.classify_and_align_evt(sp4[i],data4,centers) for i in range(len(sp4))]
2  pred4 = swp.predict_data(round4,centers)
3  data5 = data4 - pred4
4  len([x[1] for x in round4 if x[0] == '?'])
```

We can compare the third peeling with the fourth one (Fig. **??**):

```
1   plt.plot(tt, data4[0,], color='black')
2   plt.plot(tt, data5[0,], color='red',lw=0.3)
3   plt.plot(tt, data4[1,]-15, color='black')
4   plt.plot(tt, data5[1,]-15, color='red',lw=0.3)
5   plt.plot(tt, data4[2,]-25, color='black')
6   plt.plot(tt, data5[2,]-25, color='red',lw=0.3)
7   plt.plot(tt, data4[3,]-40, color='black')
8   plt.plot(tt, data5[3,]-40, color='red',lw=0.3)
9   plt.xlabel('Time␣(s)')
10  plt.xlim([3.9,4])
```

## 10.6 General comparison

We can compare the raw data with the fifth peeling on the first second (Fig. **??**):

```
1   plt.plot(tt, data0[0,], color='black')
2   plt.plot(tt, data5[0,], color='red',lw=0.3)
3   plt.plot(tt, data0[1,]-15, color='black')
4   plt.plot(tt, data5[1,]-15, color='red',lw=0.3)
5   plt.plot(tt, data0[2,]-25, color='black')
6   plt.plot(tt, data5[2,]-25, color='red',lw=0.3)
7   plt.plot(tt, data0[3,]-40, color='black')
8   plt.plot(tt, data5[3,]-40, color='red',lw=0.3)
9   plt.xlabel('Time␣(s)')
10  plt.xlim([0,1])
```

We can also look at the remaining unclassified events; they don't look like any of our templates (Fig. **??**):

```
1  bad_ones = [x[1] for x in round4 if x[0] == '?']
2  r4BE = swp.mk_events(bad_ones, data4)
3  swp.plot_events(r4BE)
```

Figure 19: 100 ms of the locust data set. Black, fourth peeling; red, fifth peeling. Two events got detected on channel 3 and subtracted.



Figure 20: The first second of the locust data set. Black, raw data; red, fifth peeling.

Figure 21: The 53 remaining bad events after the fifth peeling.

## 11 Getting the spike trains

Once we have decided to stop the peeling iterations we can extract our spike trains with (notice the syntax difference between `Python 3` and `Python 2`):

```
1  round_all = round0.copy() # Python 3
2  # round_all = round0[:] # Python 2
3  round_all.extend(round1)
4  round_all.extend(round2)
5  round_all.extend(round3)
6  round_all.extend(round4)
7  spike_trains = { n : np.sort([x[1] + x[2] for x in round_all
8                               if x[0] == n]) for n in list(centers)}
```

The number of spikes attributed to each neuron is:

```
1  [(n, len(spike_trains[n])) for n in list(centers)]
```

```
[('Cluster 0', 92),
 ('Cluster 1', 173),
 ('Cluster 2', 101),
 ('Cluster 3', 173),
 ('Cluster 4', 63),
 ('Cluster 5', 149),
 ('Cluster 6', 238),
 ('Cluster 7', 233),
```

```
('Cluster 8', 456),
('Cluster 9', 588)]
```

# 12 Individual function definitions

Short function are presented in 'one piece'. The longer ones are presented with their `docstring` first followed by the body of the function. To get the actual function you should replace the «`docstring`» appearing in the function definition by the actual `doctring`. This is just a direct application of the literate programming paradigm. More complicated functions are split into more parts with their own descriptions.

## 12.1 `plot_data_list`

We define a function, `plot_data_list`, making our raw data like displaying command lighter, starting with the `docstring`:

```
1  """Plots data when individual recording channels make up elements
2  of a list.
3
4  Parameters
5  ----------
6  data_list: a list of numpy arrays of dimension 1 that should all
7           be of the same length (not checked).
8  time_axes: an array with as many elements as the components of
9           data_list. The time values of the abscissa.
10  linewidth: the width of the lines drawing the curves.
11  color: the color of the curves.
12
13  Returns
14  -------
15  Nothing is returned, the function is used for its side effect: a
16  plot is generated.
17  """
```

Then the definition of the function per se:

```
1  def plot_data_list(data_list,
2                     time_axes,
3                     linewidth=0.2,
4                     color='black'):
5      <<plot_data_list-doctring>>
6      nb_chan = len(data_list)
7      data_min = [np.min(x) for x in data_list]
8      data_max = [np.max(x) for x in data_list]
9      display_offset = list(np.cumsum(np.array([0] +
10                                               [data_max[i]-
11                                                data_min[i-1]
12                                               for i in
13                                               range(1,nb_chan)]))))
```

35

```
14      for  i  in range(nb_chan):
15          plt.plot(time_axes,data_list[i]−display_offset[i],
16                  linewidth=linewidth,color=color)
17      plt.yticks([])
18      plt.xlabel("Time␣(s)")
```

## 12.2 peak

We define function `peak` which detects local maxima using an estimate of the derivative of the signal. Only putative maxima that are farther apart than `minimal_dist` sampling points are kept. The function returns a vector of indices. Its `docstring` is:

```
1   """Find peaks on one dimensional arrays.
2
3   Parameters
4   ——————————
5   x: a one dimensional array on which scipy.signal.fftconvolve can
6       be called.
7   minimal_dist: the minimal distance between two successive peaks.
8   not_zero: the smallest value above which the absolute value of
9   the derivative is considered not null.
10
11  Returns
12  ———————
13  An array of (peak) indices is returned.
14  """
```

And the function per se:

```
1   def peak(x, minimal_dist=15, not_zero=1e−3):
2       <<peak−docstring>>
3       ## Get the first derivative
4       dx = scipy.signal.fftconvolve(x,np.array([1,0,−1])/2.,'same')
5       dx[np.abs(dx) < not_zero] = 0
6       dx = np.diff(np.sign(dx))
7       pos = np.arange(len(dx))[dx < 0]
8       return pos[:−1][np.diff(pos) > minimal_dist]
```

## 12.3 cut_sgl_evt

Function `mk_events` (defined next) that we will use directly will call `cut_sgl_evt`. As its name says cuts a single event (an return a vector with the cuts on the different recording sites glued one after the other). Its `docstring` is:

```
1   """Cuts an 'event' at 'evt_pos' on 'data'.
2
3   Parameters
4   ——————————
5   evt_pos: an integer, the index (location) of the (peak of) the
6           event.
```

36

```
7   data: a matrix whose rows contains the recording channels.
8   before: an integer, how many points should be within the cut
9           before the reference index / time given by evt_pos.
10  after: an integer, how many points should be within the cut
11          after the reference index / time given by evt_pos.
12
13  Returns
14  -------
15  A vector with the cuts on the different recording sites glued
16  one after the other.
17  """
```

And the function per se:

```
1   def cut_sgl_evt(evt_pos,data,before=14, after=30):
2       <<cut_sgl_evt-docstring>>
3       ns = data.shape[0] ## Number of recording sites
4       dl = data.shape[1] ## Number of sampling points
5       cl = before+after+1 ## The length of the cut
6       cs = cl*ns ## The 'size' of a cut
7       cut = np.zeros((ns,cl))
8       idx = np.arange(-before,after+1)
9       keep = idx + evt_pos
10      within = np.bitwise_and(0 <= keep, keep < dl)
11      kw = keep[within]
12      cut[:,within] = data[:,kw].copy()
13      return cut.reshape(cs)
```

## 12.4 `mk_events`

Function `mk_events` takes a vector of indices as its first argument and returns a matrix with has many rows as events. Its `docstring is`

```
1   """Make events matrix out of data and events positions.
2
3   Parameters
4   ----------
5   positions: a vector containing the indices of the events.
6   data: a matrix whose rows contains the recording channels.
7   before: an integer, how many points should be within the cut
8           before the reference index / time given by evt_pos.
9   after: an integer, how many points should be within the cut
10          after the reference index / time given by evt_pos.
11
12  Returns
13  -------
14  A matrix with as many rows as events and whose rows are the cuts
15  on the different recording sites glued one after the other.
16  """
```

And the function per se:

```
1  def mk_events(positions, data, before=14, after=30):
2      <<mk_events−docstring >>
3      res = np.zeros((len(positions),(before+after+1)∗data.shape[0]))
4      for i,p in enumerate(positions):
5          res[i,:] = cut_sgl_evt(p,data,before,after)
6      return res
```

## 12.5 `plot_events`

In order to facilitate events display, we define an event specific plotting function starting with its `docstring`:

```
1  """Plot events.
2
3  Parameters
4  −−−−−−−−−−
5  evts_matrix: a matrix of events. Rows are events. Cuts from
6              different recording sites are glued one after the
7              other on each row.
8  n_plot: an integer, the number of events to plot (if 'None',
9          default, all are shown).
10 n_channels: an integer, the number of recording channels.
11 events_color: the color used to display events.
12 events_lw: the line width used to display events.
13 show_median: should the median event be displayed?
14 median_color: color used to display the median event.
15 median_lw: line width used to display the median event.
16 show_mad: should the MAD be displayed?
17 mad_color: color used to display the MAD.
18 mad_lw: line width used to display the MAD.
19
20 Returns
21 −−−−−−−
22 Noting, the function is used for its side effect.
23 """
```

And the function per se:

```
1  def plot_events(evts_matrix,
2                  n_plot=None,
3                  n_channels=4,
4                  events_color='black',
5                  events_lw=0.1,
6                  show_median=True,
7                  median_color='red',
8                  median_lw=0.5,
9                  show_mad=True,
10                 mad_color='blue',
11                 mad_lw=0.5):
12     <<plot_events−docstring >>
13     if n_plot is None:
```

```
14            n_plot = evts_matrix.shape[0]
15
16        cut_length = evts_matrix.shape[1] // n_channels
17
18        for i in range(n_plot):
19            plt.plot(evts_matrix[i,:], color=events_color, lw=events_lw)
20        if show_median:
21            MEDIAN = np.apply_along_axis(np.median,0,evts_matrix)
22            plt.plot(MEDIAN, color=median_color, lw=median_lw)
23
24        if show_mad:
25            MAD = np.apply_along_axis(mad,0,evts_matrix)
26            plt.plot(MAD, color=mad_color, lw=mad_lw)
27
28        left_boundary = np.arange(cut_length,
29                                  evts_matrix.shape[1],
30                                  cut_length*2)
31        for l in left_boundary:
32            plt.axvspan(l,l+cut_length-1,
33                        facecolor='grey',alpha=0.5,edgecolor='none')
34        plt.xticks([])
35        return
```

## 12.6 plot_data_list_and_detection

We define a function, plot_data_list_and_detection, making our data and detection displaying command lighter. Its docstring:

```
1  """Plots data together with detected events.
2
3  Parameters
4  ----------
5  data_list: a list of numpy arrays of dimension 1 that should all
6             be of the same length (not checked).
7  time_axes: an array with as many elements as the components of
8             data_list. The time values of the abscissa.
9  evts_pos: a vector containing the indices of the detected
10            events.
11  linewidth: the width of the lines drawing the curves.
12  color: the color of the curves.
13
14  Returns
15  -------
16  Nothing is returned, the function is used for its side effect: a
17  plot is generated.
18  """
```

And the function:

```
1  def plot_data_list_and_detection(data_list,
2                                   time_axes,
3                                   evts_pos,
```

```
4                                    linewidth =0.2 ,
5                                    color=' black ' ):
6       <<plot_data_list_and_detection-docstring >>
7       nb_chan = len ( data_list )
8       data_min = [ np . min( x ) for x in data_list ]
9       data_max = [ np . max( x ) for x in data_list ]
10      display_offset = list ( np . cumsum ( np . array ( [ 0 ] +
11                                        [ data_max [ i ]-
12                                          data_min [ i -1] for i in
13                                        range ( 1 , nb_chan ) ] ) ) )
14      for i in range ( nb_chan ):
15          plt . plot ( time_axes , data_list [ i ]-display_offset [ i ] ,
16                     linewidth=linewidth , color=color )
17          plt . plot ( time_axes [ evts_pos ] ,
18                     data_list [ i ][ evts_pos]-display_offset [ i ] , ' ro ' )
19      plt . yticks ( [ ] )
20      plt . xlabel ( "Time␣( s )" )
```

## 12.7 `mk_noise`

Getting an estimate of the noise statistical properties is an essential ingredient to build respectable goodness of fit tests. In our approach "noise events" are essentially anything that is not an "event". I wrote essentially and not exactly since there is a little twist here which is the minimal distance we are willing to accept between the reference time of a noise event and the reference time of the last preceding and of the first following "event". We could think that keeping a cut length on each side would be enough. That would indeed be the case if *all* events were starting from and returning to zero within a cut. But this is not the case with the cuts parameters we chose previously (that will become clear soon). You might wonder why we chose so short a cut length then. Simply to avoid having to deal with too many superposed events which are the really bothering events for anyone wanting to do proper sorting. To obtain our noise events we are going to use function `mk_noise` which takes the *same* arguments as function `mk_events` plus two numbers:

- `safety_factor` a number by which the cut length is multiplied and which sets the minimal distance between the reference times discussed in the previous paragraph.

- `size` the maximal number of noise events one wants to cut (the actual number obtained might be smaller depending on the data length, the cut length, the safety factor and the number of events).

We define now function `mk_noise` starting with its `docstring`:

```
1    """Constructs a noise sample .
2
3    Parameters
4    ----------
5    positions : a vector containing the indices of the events .
6    data : a matrix whose rows contains the recording channels .
7    before : an integer , how many points should be within the cut
```

```
 8            before the reference index / time given by evt_pos.
 9  after: an integer, how many points should be within the cut
10            after the reference index / time given by evt_pos.
11  safety_factor: a number by which the cut length is multiplied
12                  and which sets the minimal distance between the
13                  reference times discussed in the previous
14                  paragraph.
15  size: the maximal number of noise events one wants to cut (the
16         actual number obtained might be smaller depending on the
17         data length, the cut length, the safety factor and the
18         number of events).
19
20  Returns
21  -------
22  A matrix with as many rows as noise events and whose rows are
23  the cuts on the different recording sites glued one after the
24  other.
25  """
```

And the function:

```
 1  def mk_noise(positions, data, before=14, after=30, safety_factor=2, size=2000):
 2      <<mk_noise-docstring>>
 3      sl = before+after+1 ## cut length
 4      ns = data.shape[0] ## number of recording sites
 5      i1 = np.diff(positions) ## inter-event intervals
 6      minimal_length = round(sl*safety_factor)
 7      ## Get next the number of noise sweeps that can be
 8      ## cut between each detected event with a safety factor
 9      nb_i = (i1-minimal_length)//sl
10      ## Get the number of noise sweeps that are going to be cut
11      nb_possible = min(size,sum(nb_i[nb_i>0]))
12      res = np.zeros((nb_possible,sl*data.shape[0]))
13      ## Create next a list containing the indices of the inter event
14      ## intervals that are long enough
15      idx_l = [i for i in range(len(i1)) if nb_i[i] > 0]
16      ## Make next an index running over the inter event intervals
17      ## from which at least one noise cut can be made
18      interval_idx = 0
19      ## noise_positions = np.zeros(nb_possible,dtype=numpy.int)
20      n_idx = 0
21      while n_idx < nb_possible:
22          within_idx = 0 ## an index of the noise cut with a long enough
23                         ## interval
24          i_pos = positions[idx_l[interval_idx]] + minimal_length
25          ## Variable defined next contains the number of noise cuts
26          ## that can be made from the "currently" considered long-enough
27          ## inter event interval
28          n_at_interval_idx = nb_i[idx_l[interval_idx]]
29          while within_idx < n_at_interval_idx and n_idx < nb_possible:
30              res[n_idx,:]= cut_sgl_evt(int(i_pos),data,before,after)
31              ## noise_positions[n_idx] = i_pos
32              n_idx += 1
```

41

```
33              i_pos += sl
34              within_idx += 1
35           interval_idx += 1
36      ## return (res, noise_positions)
37      return res
```

## 12.8 `mad`

We define the `mad` function in one piece since it is very short:

```
1  def mad(x):
2      """Returns the Median Absolute Deviation of its argument.
3      """
4      return np.median(np.absolute(x - np.median(x)))*1.4826
```

## 12.9 `mk_aligned_events`

### 12.9.1 The jitter: A worked out example

Function `mk_aligned_events` is somehow the "heavy part" of this document. Its job is to align events on their templates while taking care of two jitter sources: the sampling and the noise one. Rather than getting into a theoretical discussion, we illustrate the problem with one of the events detected on our data set. Cluster 1 is the cluster exhibiting the largest sampling jitter effects, since it has the largest time derivative, in absolute value, of its median event . This is clearly seen when we superpose the 50th event from this cluster with the median event (remember that we start numbering at 0). So we get first our estimate for center or template of cluster 1:

```
1  c1_median = apply(np.median,0,evtsE[goodEvts,:][np.array(c10b)==1,:])
```

And we do the plot (Fig. **??**):

```
1  plt.plot(c1_median,color='red')
2  plt.plot(evtsE[goodEvts,:][np.array(c10b)==1,:][50,:],color='black')
```

A Taylor expansion shows that if we write $g(t)$ the observed 50th event, the sampling jitter and $f(t)$ the actual waveform of the event then:

$$g(t) = f(t+) + (t) \approx f(t) + f'(t) +^2 /2\, f''(t) + (t)\,; \qquad (1)$$

where is a Gaussian process and where $f'$ and $f''$ stand for the first and second time derivatives of $f$. Therefore, if we can get estimates of $f'$ and $f''$ we should be able to estimate by linear regression (if we neglect the $^2$ term as well as the potentially non null correlation in ) or by non linear regression (if we keep the latter). We start by getting the derivatives estimates:

```
1  dataD = apply(lambda x: fftconvolve(x,np.array([1,0,-1])/2.,'same'),
2                1, data)
```

Figure 22: The median event of cluster 1 (red) together with event 50 of the same cluster (black).

```
3  evtsED = swp.mk_events(sp0E,dataD,14,30)
4  dataDD = apply(lambda x: fftconvolve(x,np.array([1,0,-1])/2.,'same'),
5                 1, dataD)
6  evtsEDD = swp.mk_events(sp0E,dataDD,14,30)
7  c1D_median = apply(np.median,0,
8                      evtsED[goodEvts,:][np.array(c10b)==1,:])
9  c1DD_median = apply(np.median,0,
10                     evtsEDD[goodEvts,:][np.array(c10b)==1,:])
```

We then get something like Fig. **??**:

```
1  plt.plot(evtsE[goodEvts,:][np.array(c10b)==1,:][50,:]-\
2           c1_median,color='red',lw=2)
3  plt.plot(1.5*c1D_median,color='blue',lw=2)
4  plt.plot(1.5*c1D_median+1.5**2/2*c1DD_median,color='black',lw=2)
```

If we neglect the $^2$ term we quickly arrive at:

$$\hat{} = \frac{\mathbf{f}' \cdot (\mathbf{g} - \mathbf{f})}{\|\mathbf{f}'\|^2} ; \tag{2}$$

where the 'vectorial' notation like $\mathbf{a} \cdot \mathbf{b}$ stands here for:

$$\sum_{i=0}^{179} a_i b_i .$$

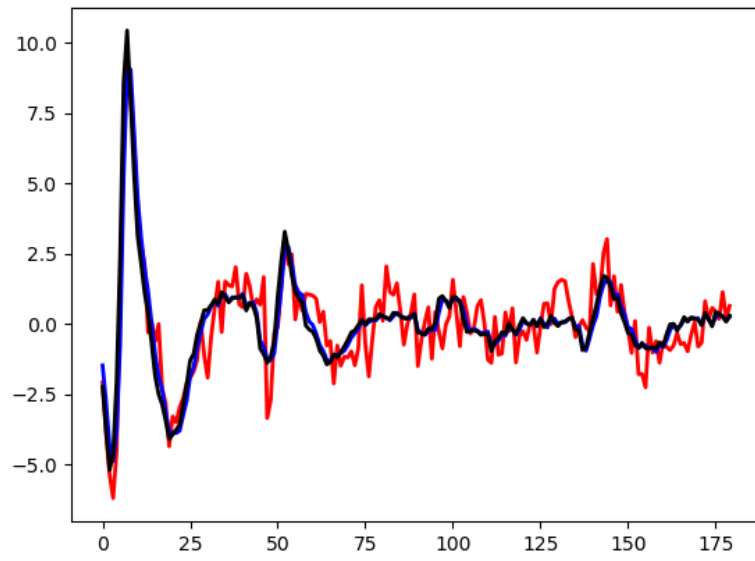Figure 23: The median event of cluster 1 subtracted from event 50 of the same cluster (red); 1.5 times the first derivative of the median event (blue)—corresponding to =1.5—; 1.5 times the first derivative + 1.5ˆ2/2 times the second (black)—corresponding again to =1.5—.

For the 50th event of the cluster we get:

```
1  delta_hat = np.dot(c1D_median,
2                     evtsE[goodEvts,:][np.array(c10b)==1,:][50,:]-\
3                     c1_median)/np.dot(c1D_median,c1D_median)
4  delta_hat
```

1.4917182304326997

We can use this estimated value of `delta_hat` as an initial guess for a procedure refining the estimate using also the $^2$ term. The obvious quantity we should try to minimize is the residual sum of square, RSS defined by:

$$\mathrm{RSS}() = \|\mathbf{g} - \mathbf{f} - \mathbf{f}' - {}^2/2\,\mathbf{f}''\|^2 \; .$$

We can define a function returning the RSS for a given value of  as well as an event `evt` a cluster center (median event of the cluster) `center` and its first two derivatives, `centerD` and `centerDD`:

```
1  def rss_fct(delta, evt, center, centerD, centerDD):
2      return np.sum((evt - center - delta*centerD - delta**2/2*centerDD)**2)
```

To create quickly a graph of the RSS as a function of  for the specific case we are dealing with now (51st element of cluster 1) we create a vectorized or *universal* function version of the `rss_for_alignment` we just defined:

```
1  urss_fct = np.frompyfunc(lambda x:
2                           rss_fct(x,
3                                   evtsE[goodEvts,:]\
4                                   [np.array(c10b)==1,:][50,:],
5                                   c1_median,c1D_median,c1DD_median),1,1)
```

We then get the Fig. **??** with:

```
1  plt.subplot(1,2,1)
2  dd = np.arange(-5,5,0.05)
3  plt.plot(dd,urss_fct(dd),color='black',lw=2)
4  plt.subplot(1,2,2)
5  dd_fine = np.linspace(delta_hat-0.5,delta_hat+0.5,501)
6  plt.plot(dd_fine,urss_fct(dd_fine),color='black',lw=2)
7  plt.axvline(x=delta_hat,color='red')
```

The left panel of the above figure shows that our initial guess for ˆ is not bad but still approximately 0.2 units away from the actual minimum. The classical way to refine our  estimate—in 'nice situations' where the function we are trying to minimize is locally convex—is to use the <span style="color:magenta">Newton-Raphson algorithm</span> which consists in approximating locally the 'target function' (here our RSS function) by a parabola having locally the same first and second derivatives, before jumping to the minimum of this approximating parabola. If we develop our previous expression of RSS() we get:

$$\mathrm{RSS}() = \|\mathbf{h}\|^2 - 2\,\mathbf{h}\cdot\mathbf{f}' + {}^2\left(\|\mathbf{f}'\|^2 - \mathbf{h}\cdot\mathbf{f}''\right) + {}^3\,\mathbf{f}'\cdot\mathbf{f}'' + \frac{{}^4}{4}\|\mathbf{f}''\|^2 \; ;$$
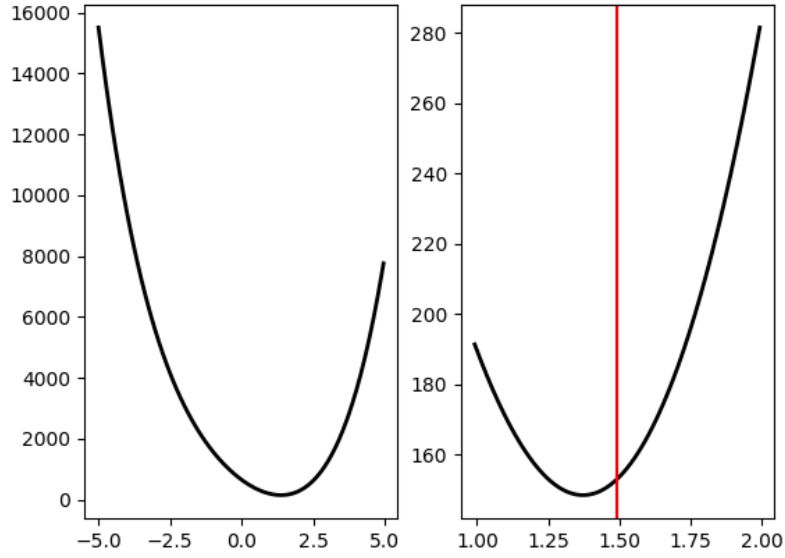
45

Figure 24: The RSS as a function of for event 50 of cluster 1. Left, $\in [-5, 5]$; right, $\in [\hat{} - 0.5, \hat{} + 0.5]$ and the red vertical line shows $\hat{}$.

where $\mathbf{h}$ stands for $\mathbf{g} - \mathbf{f}$. By differentiation with respect to we get:

$$\mathrm{RSS}'() = -2\,\mathbf{h} \cdot \mathbf{f}' + 2\,\left(\|\mathbf{f}'\|^2 - \mathbf{h} \cdot \mathbf{f}''\right) + 3\,^2\,\mathbf{f}' \cdot \mathbf{f}'' + ^3\,\|\mathbf{f}''\|^2 \,.$$

And a second differentiation leads to:

$$\mathrm{RSS}''() = 2\,\left(\|\mathbf{f}'\|^2 - \mathbf{h} \cdot \mathbf{f}''\right) + 6\,\mathbf{f}' \cdot \mathbf{f}'' + 3\,^2\|\mathbf{f}''\|^2 \,.$$

The equation of the approximating parabola at $^{(k)}$ is then:

$$\mathrm{RSS}(^{(k)}+) \approx \mathrm{RSS}(^{(k)}) + \,\mathrm{RSS}'(^{(k)}) + \frac{^2}{2}\,\mathrm{RSS}''(^{(k)}) \,,$$

and its minimum—if $\mathrm{RSS}''() > 0$—is located at:

$$^{(k+1)} = ^{(k)} - \frac{\mathrm{RSS}'(^{(k)})}{\mathrm{RSS}''(^{(k)})} \,.$$

Defining functions returning the required derivatives:

```
def rssD_fct(delta, evt, center, centerD, centerDD):
    h = evt - center
    return -2*np.dot(h, centerD) + \
        2*delta*(np.dot(centerD, centerD) - np.dot(h, centerDD)) + \
        3*delta**2*np.dot(centerD, centerDD) + \
```

```
6               delta**3*np.dot(centerDD,centerDD)
7
8    def rssDD_fct(delta,evt,center,centerD,centerDD):
9        h = evt − center
10       return 2*(np.dot(centerD,centerD) − np.dot(h,centerDD)) + \
11           6*delta*np.dot(centerD,centerDD) + \
12           3*delta**2*np.dot(centerDD,centerDD)
```

we can get a graphical representation (Fig. **??**) of a single step of the Newton-Raphson algorithm:

```
1    rss_at_delta0 = rss_fct(delta_hat,
2                            evtsE[goodEvts,:][np.array(c10b)==1,:][50,:],
3                            c1_median,c1D_median,c1DD_median)
4    rssD_at_delta0 = rssD_fct(delta_hat,
5                            evtsE[goodEvts,:][np.array(c10b)==1,:][50,:],
6                            c1_median,c1D_median,c1DD_median)
7    rssDD_at_delta0 = rssDD_fct(delta_hat,
8                            evtsE[goodEvts,:][np.array(c10b)==1,:]\
9                            [50,:],c1_median,c1D_median,c1DD_median)
10   delta_1 = delta_hat − rssD_at_delta0/rssDD_at_delta0
```

```
1    plt.plot(dd_fine,urss_fct(dd_fine),color='black',lw=2)
2    plt.axvline(x=delta_hat,color='red')
3    plt.plot(dd_fine,
4             rss_at_delta0 + (dd_fine−delta_hat)*rssD_at_delta0 + \
5             (dd_fine−delta_hat)**2/2*rssDD_at_delta0,color='blue',lw=2)
6    plt.axvline(x=delta_1,color='grey')
```

Subtracting the second order in approximation of f(t+) from the observed 50th event of cluster 1 we get Fig. **??**:

```
1    plt.plot(evtsE[goodEvts,:][np.array(c10b)==1,:][50,:]−\
2             c1_median−delta_1*c1D_median−delta_1**2/2*c1DD_median,
3             color='red',lw=2)
4    plt.plot(evtsE[goodEvts,:][np.array(c10b)==1,:][50,:],
5             color='black',lw=2)
6    plt.plot(c1_median+delta_1*c1D_median+delta_1**2/2*c1DD_median,
7             color='blue',lw=1)
```

### 12.9.2 Function definition

We start with the chunk importing the required functions from the different modules («mk_aligned_events-import-functions»):

```
1    from scipy.signal import fftconvolve
2    from numpy import apply_along_axis as apply
3    from scipy.spatial.distance import squareform
```

We then get the first and second derivatives of the data:

Figure 25: The RSS as a function of  for event 50 of cluster 1 (black), the red vertical line
shows ˆ. In blue, the approximating parabola at ˆ. The grey vertical line shows
the minimum of the approximating parabola.

Figure 26: Event 50 of cluster 1 (black), second order approximation of f(t+) (blue) and residual (red) for —obtained by a succession of a linear regression (order 1) and a single Newton-Raphson step—equal to: 1.3748048144324905.

```
1  dataD = apply(lambda x: fftconvolve(x,np.array([1,0,-1])/2., 'same'),
2              1, data)
3  dataDD = apply(lambda x: fftconvolve(x,np.array([1,0,-1])/2., 'same'),
4              1, dataD)
```

Events are cut from the different data 'versions', derivatives of order 0, 1 and 2 («mk_aligned_events-get-eve

```
1  evts = mk_events(positions, data, before, after)
2  evtsD = mk_events(positions, dataD, before, after)
3  evtsDD = mk_events(positions, dataDD, before, after)
```

A center or template is obtained by taking the pointwise median of the events we just got on the three versions of the data («mk_aligned_events-get-centers»):

```
1  center = apply(np.median,0,evts)
2  centerD = apply(np.median,0,evtsD)
3  centerD_norm2 = np.dot(centerD,centerD)
4  centerDD = apply(np.median,0,evtsDD)
5  centerDD_norm2 = np.dot(centerDD,centerDD)
6  centerD_dot_centerDD = np.dot(centerD,centerDD)
```
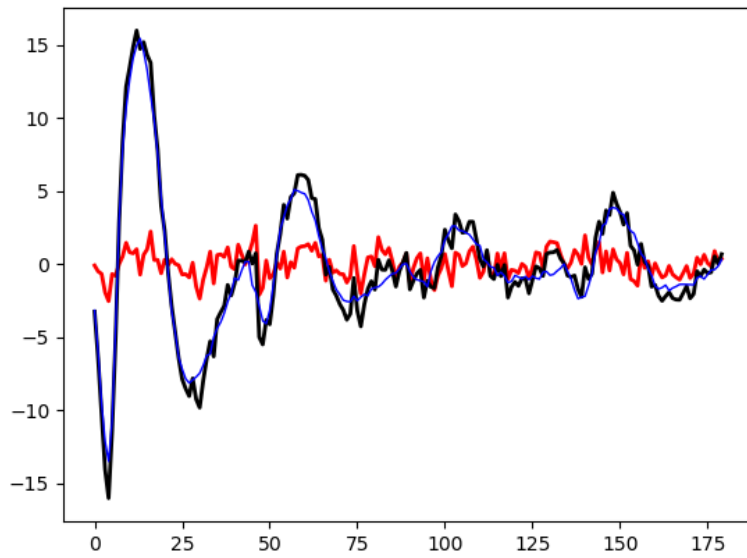
Given an event, make a first order jitter estimation and compute the norm of the initial residual, h_order0_norm2, and of its first order jitter corrected version, h_order1_norm2 («mk_aligned_events-do-job-on-single-event-order1»):

```
1  h = evt - center
2  h_order0_norm2 = sum(h**2)
3  h_dot_centerD = np.dot(h,centerD)
4  jitter0 = h_dot_centerD/centerD_norm2
5  h_order1_norm2 = sum((h-jitter0*centerD)**2)
```

If the residual's norm decrease upon first order jitter correction, try a second order one. At the end compare the norm of the second order jitter corrected residual (h_order2_norm2) with the one of the first order (h_order1_norm2). If the former is larger or equal than the latter, set the estimated jitter to its first order value («mk_aligned_events-do-job-on-single-event-order

```
1  h_dot_centerDD = np.dot(h,centerDD)
2  first = -2*h_dot_centerD + \
3     2*jitter0*(centerD_norm2 - h_dot_centerDD) + \
4     3*jitter0**2*centerD_dot_centerDD + \
5     jitter0**3*centerDD_norm2
6  second = 2*(centerD_norm2 - h_dot_centerDD) + \
7     6*jitter0*centerD_dot_centerDD + \
8     3*jitter0**2*centerDD_norm2
9  jitter1 = jitter0 - first/second
10 h_order2_norm2 = sum((h-jitter1*centerD - \
11                    jitter1**2/2*centerDD)**2)
12 if h_order1_norm2 <= h_order2_norm2:
13     jitter1 = jitter0
```

And now the function's docstring («mk_aligned_events-docstring»):

```
1   """Align events on the central event using first or second order
2   Taylor expansion.
3
4   Parameters
5   ——————————
6   positions: a vector of indices with the positions of the
7              detected events.
8   data: a matrix whose rows contains the recording channels.
9   before: an integer, how many points should be within the cut
10         before the reference index / time given by positions.
11  after: an integer, how many points should be within the cut
12         after the reference index / time given by positions.
13
14  Returns
15  ———————
16  A tuple whose elements are:
17    A matrix with as many rows as events and whose rows are the
18    cuts on the different recording sites glued one after the
19    other. These events have been jitter corrected using the
20    second order Taylor expansion.
21    A vector of events positions where "actual" positions have
22    been rounded to the nearest index.
23    A vector of jitter values.
24
25  Details
26  ———————
27  (1) The data first and second derivatives are estimated first.
28  (2) Events are cut next on each of the three versions of the data.
29  (3) The global median event for each of the three versions are
30  obtained.
31  (4) Each event is then aligned on the median using a first order
32  Taylor expansion.
33  (5) If this alignment decreases the squared norm of the event
34  (6) an improvement is looked for using a second order expansion.
35  If this second order expansion still decreases the squared norm
36  and if the estimated jitter is larger than 1, the whole procedure
37  is repeated after cutting a new the event based on a better peak
38  position (7).
39  """
```

To end up with the function itself:

```
1   def mk_aligned_events(positions, data, before=14, after=30):
2       <<mk_aligned_events-docstring>>
3       <<mk_aligned_events-import-functions>>
4       n_evts = len(positions)
5       new_positions = positions.copy()
6       jitters = np.zeros(n_evts)
7       # Details (1)
8       <<mk_aligned_events-dataD-and-dataDD>>
9       # Details (2)
10      <<mk_aligned_events-get-events>>
```

```
11        # Details (3)
12        <<mk_aligned_events-get-centers>>
13        # Details (4)
14        for evt_idx in range(n_evts):
15            # Details (5)
16            evt = evts[evt_idx,:]
17            evt_pos = positions[evt_idx]
18            <<mk_aligned_events-do-job-on-single-event-order1>>
19            if h_order0_norm2 > h_order1_norm2:
20                # Details (6)
21                <<mk_aligned_events-do-job-on-single-event-order2>>
22            else:
23                jitter1 = 0
24            if abs(round(jitter1)) > 0:
25                # Details (7)
26                evt_pos -= int(round(jitter1))
27                evt = cut_sgl_evt(evt_pos,data=data,
28                                  before=before, after=after)
29                <<mk_aligned_events-do-job-on-single-event-order1>>
30                if h_order0_norm2 > h_order1_norm2:
31                    <<mk_aligned_events-do-job-on-single-event-order2>>
32                else:
33                    jitter1 = 0
34            if sum(evt**2) > sum((h-jitter1*centerD-
35                                  jitter1**2/2*centerDD)**2):
36                evts[evt_idx,:] = evt-jitter1*centerD- \
37                    jitter1**2/2*centerDD
38            new_positions[evt_idx] = evt_pos
39            jitters[evt_idx] = jitter1
40        return (evts, new_positions,jitters)
```

## 12.10 `mk_center_dictionary`

We define function `mk_center_dictionary` starting with its `docstring`:

```
1  """ Computes clusters 'centers' or templates and associated data.
2
3  Clusters' centers should be built such that they can be used for
4  subtraction, this implies that we should make them long enough, on
5  both side of the peak, to see them go back to baseline. Formal
6  parameters before and after bellow should therefore be set to
7  larger values than the ones used for clustering.
8
9  Parameters
10 ----------
11 positions : a vector of spike times, that should all come from the
12             same cluster and correspond to reasonably 'clean'
13             events.
14 data : a data matrix.
15 before : the number of sampling point to keep before the peak.
16 after : the number of sampling point to keep after the peak.
17
```

```
18    Returns
19    -------
20  A dictionary with the following components:
21     center: the estimate of the center (obtained from the median).
22     centerD: the estimate of the center's derivative (obtained from
23              the median of events cut on the derivative of data).
24     centerDD: the estimate of the center's second derivative
25              (obtained from the median of events cut on the second
26              derivative of data).
27     centerD_norm2: the squared norm of the center's derivative.
28     centerDD_norm2: the squared norm of the center's second
29                     derivative.
30     centerD_dot_centerDD: the scalar product of the center's first
31                           and second derivatives.
32     center_idx: an array of indices generated by
33                 np.arange(-before, after+1).
34     """
```

The function starts by evaluating the first two derivatives of the data («get-derivatives»):

```
1  from scipy.signal import fftconvolve
2  from numpy import apply_along_axis as apply
3  dataD = apply(lambda x:
4                fftconvolve(x,np.array([1,0,-1])/2.,'same'),
5                1, data)
6  dataDD = apply(lambda x:
7                 fftconvolve(x,np.array([1,0,-1])/2.,'same'),
8                 1, dataD)
```

The function is defined next:

```
1  def mk_center_dictionary(positions, data, before=49, after=80):
2      <<mk_center_dictionary-docstring>>
3      <<mk_center_dictionary-get-derivatives>>
4      evts = mk_events(positions, data, before, after)
5      evtsD = mk_events(positions, dataD, before, after)
6      evtsDD = mk_events(positions, dataDD, before, after)
7      evts_median = apply(np.median,0,evts)
8      evtsD_median = apply(np.median,0,evtsD)
9      evtsDD_median = apply(np.median,0,evtsDD)
10     return {"center" : evts_median,
11             "centerD" : evtsD_median,
12             "centerDD" : evtsDD_median,
13             "centerD_norm2" : np.dot(evtsD_median,evtsD_median),
14             "centerDD_norm2" : np.dot(evtsDD_median,evtsDD_median),
15             "centerD_dot_centerDD" : np.dot(evtsD_median,
16                                             evtsDD_median),
17             "center_idx" : np.arange(-before,after+1)}
```

## 12.11 classify_and_align_evt

We now define with the following docstring («classify_and_align_evt-docstring»):

```
1  """Compares a single event to a dictionary of centers and returns
2  the name of the closest center if it is close enough or '?', the
3  corrected peak position and the remaining jitter.
4
5  Parameters
6  ——————————
7  evt_pos : a sampling point at which an event was detected.
8  data : a data matrix.
9  centers : a centers' dictionary returned by mk_center_dictionary.
10 before : the number of sampling point to consider before the peak.
11 after : the number of sampling point to consider after the peak.
12
13 Returns
14 ———————
15 A list with the following components:
16    The name of the closest center if it was close enough or '?'.
17    The nearest sampling point to the events peak.
18    The jitter: difference between the estimated actual peak
19    position and the nearest sampling point.
20 """
```

The first chunk of the function takes a dictionary of centers, `centers`, generated by `mk_center_dictionary`, defines two variables, `cluster_names` and `n_sites`, and builds a matrix of centers, `centersM`:

```
1  cluster_names = np.sort(list(centers))
2  n_sites = data.shape[0]
3  centersM = np.array([centers[c_name]["center"]\
4                      [np.tile((-before <= centers[c_name]\
5                              ["center_idx"]).\
6                              __and__(centers[c_name]["center_idx"] \
7                                      <= after), n_sites)]
8                                      for c_name in cluster_names])
```

Extract the event, `evt`, to classify and subtract each center from it, `delta`, to find the closest one, `cluster_idx`, using the Euclidean squared norm («cluster_idx»):

```
1  evt = cut_sgl_evt(evt_pos, data=data, before=before, after=after)
2  delta = -(centersM - evt)
3  cluster_idx = np.argmin(np.sum(delta**2, axis=1))
```

Get the name of the selected cluster, `good_cluster_name`, and its 'time indices', `good_cluster_idx`. Then, extract the first two derivatives of the center, `centerD` and `centerDD`, their squared norms, `centerD_norm2` and `centerDD_norm2`, and their dot product, `centerD_dot_centerDD` («get-centers»):

```
1  good_cluster_name = cluster_names[cluster_idx]
2  good_cluster_idx = np.tile((-before <= centers[good_cluster_name]\
3                              ["center_idx"]).\
4                              __and__(centers[good_cluster_name]\
5                                      ["center_idx"] <= after),
6                                      n_sites)
```

```
7   centerD = centers [ good_cluster_name ] [ "centerD" ] [ good_cluster_idx ]
8   centerD_norm2 = np.dot ( centerD , centerD )
9   centerDD = centers [ good_cluster_name ] [ "centerDD" ] [ good_cluster_idx ]
10  centerDD_norm2 = np.dot ( centerDD , centerDD )
11  centerD_dot_centerDD = np.dot ( centerD , centerDD )
```

Do a first order jitter correction where `h` contains the difference between the event and the center. Obtain the estimated jitter, `jitter0` and the squared norm of the first order corrected residual, `h_order1_norm2` («jitter-order-1»):

```
1   h_order0_norm2 = sum ( h**2 )
2   h_dot_centerD = np.dot ( h , centerD )
3   jitter0 = h_dot_centerD / centerD_norm2
4   h_order1_norm2 = sum ( ( h−jitter0*centerD )**2 )
```

Do a second order jitter correction. Obtain the estimated jitter, `jitter1` and the squared norm of the second order corrected residual, `h_order2_norm2` («jitter-order-2»):

```
1   h_dot_centerDD = np.dot ( h , centerDD )
2   first = −2*h_dot_centerD + \
3     2*jitter0*( centerD_norm2 − h_dot_centerDD ) + \
4     3*jitter0**2*centerD_dot_centerDD + \
5     jitter0**3*centerDD_norm2
6   second = 2*( centerD_norm2 − h_dot_centerDD ) + \
7     6*jitter0*centerD_dot_centerDD + \
8     3*jitter0**2*centerDD_norm2
9   jitter1 = jitter0 − first/second
10  h_order2_norm2 = sum ( ( h−jitter1*centerD−jitter1**2/2*centerDD )**2 )
```

Now define the function:

```
1   def classify_and_align_evt ( evt_pos , data , centers ,
2                                 before =14 , after =30):
3       <<classify_and_align_evt−docstring >>
4       <<classify_and_align_evt−centersM >>
5       <<classify_and_align_evt−cluster_idx >>
6       <<classify_and_align_evt−get−centers >>
7       h = delta [ cluster_idx ,:]
8       <<classify_and_align_evt−jitter−order−1>>
9       if h_order0_norm2 > h_order1_norm2 :
10          <<classify_and_align_evt−jitter−order−2>>
11          if h_order1_norm2 <= h_order2_norm2 :
12              jitter1 = jitter0
13      else :
14          jitter1 = 0
15      if abs ( round ( jitter1 )) > 0:
16          evt_pos −= int ( round ( jitter1 ))
17          evt = cut_sgl_evt ( evt_pos , data=data ,
18                              before=before , after=after )
19          h = evt − centers [ good_cluster_name ] [ "center" ]\
20            [ good_cluster_idx ]
21          <<classify_and_align_evt−jitter−order−1>>
22          if h_order0_norm2 > h_order1_norm2 :
```

```
23              <<classify_and_align_evt−jitter−order−2>>
24              if h_order1_norm2 <= h_order2_norm2:
25                  jitter1 = jitter0
26          else:
27              jitter1 = 0
28      if sum(evt**2) > sum((h−jitter1*centerD−jitter1**2/2*centerDD)**2):
29          return [cluster_names[cluster_idx], evt_pos, jitter1]
30      else:
31          return ['?',evt_pos, jitter1]
```

## 12.12 `predict_data`

We define function `predict_data` that creates an ideal data trace given events' positions, events' origins and a clusters' catalog. We start with the `docstring`:

```
1  """Predicts ideal data given a list of centers' names, positions,
2  jitters and a dictionary of centers.
3
4  Parameters
5  ——————————
6  class_pos_jitter_list : a list of lists returned by
7                          classify_and_align_evt.
8  centers_dictionary : a centers' dictionary returned by
9                       mk_center_dictionary.
10 nb_channels : the number of recording channels.
11 data_length : the number of sampling points.
12
13 Returns
14 ———————
15 A matrix of ideal (noise free) data with nb_channels rows and
16 data_length columns.
17 """
```

And the function:

```
1  def predict_data(class_pos_jitter_list,
2                   centers_dictionary,
3                   nb_channels=4,
4                   data_length=300000):
5      <<predict_data−docstring>>
6      ## Create next a matrix that will contain the results
7      res = np.zeros((nb_channels,data_length))
8      ## Go through every list element
9      for class_pos_jitter in class_pos_jitter_list:
10         cluster_name = class_pos_jitter[0]
11         if cluster_name != '?':
12             center = centers_dictionary[cluster_name]["center"]
13             centerD = centers_dictionary[cluster_name]["centerD"]
14             centerDD = centers_dictionary[cluster_name]["centerDD"]
15             jitter = class_pos_jitter[2]
16             pred = center + jitter*centerD + jitter**2/2*centerDD
17             pred = pred.reshape((nb_channels,len(center)//nb_channels))
```

```
18              idx = centers_dictionary[cluster_name]["center_idx"] + \
19                class_pos_jitter[1]
20              ## Make sure that the event is not too close to the
21              ## boundaries
22              within = np.bitwise_and(0 <= idx, idx < data_length)
23              kw = idx[within]
24              res[:,kw] += pred[:,within]
25      return res
```