# Feature Engineering - 2

# Feature Optimization

- Normalization

- Interaction between Features

- Binning

- Adaptive Binning

- Thresholding

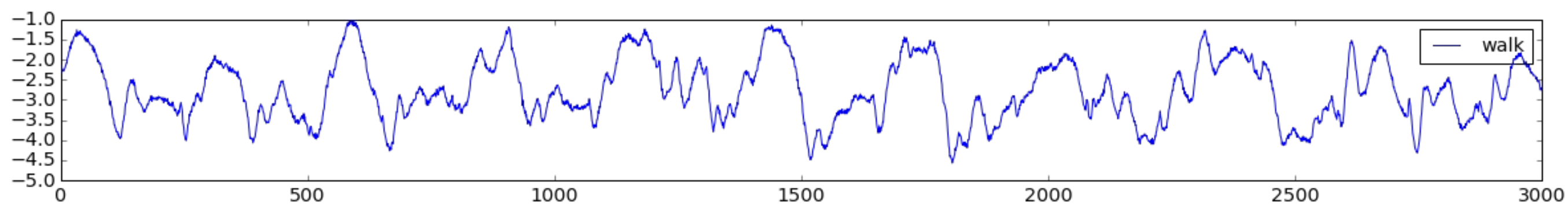- Scaling

- Log Transformation

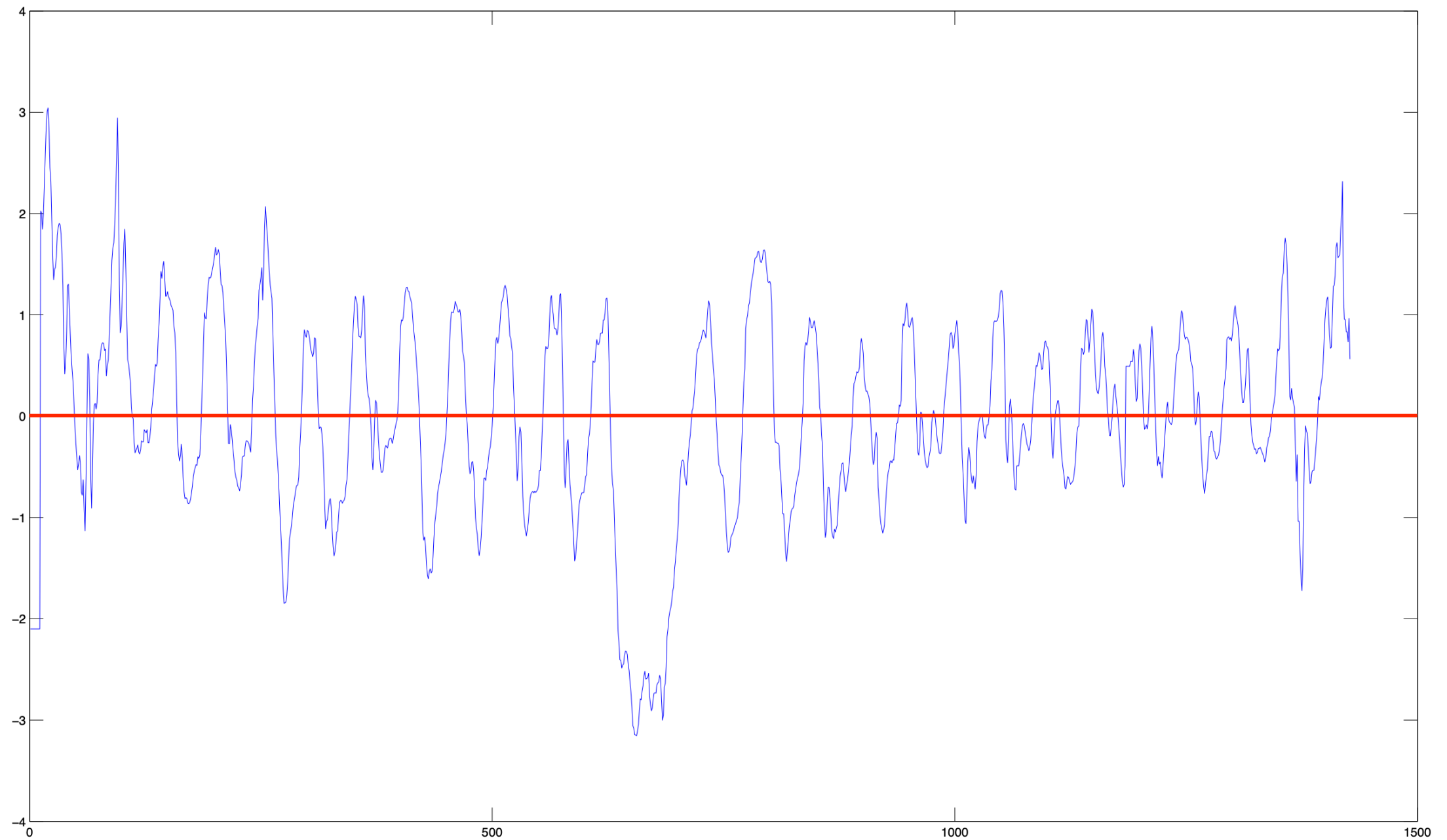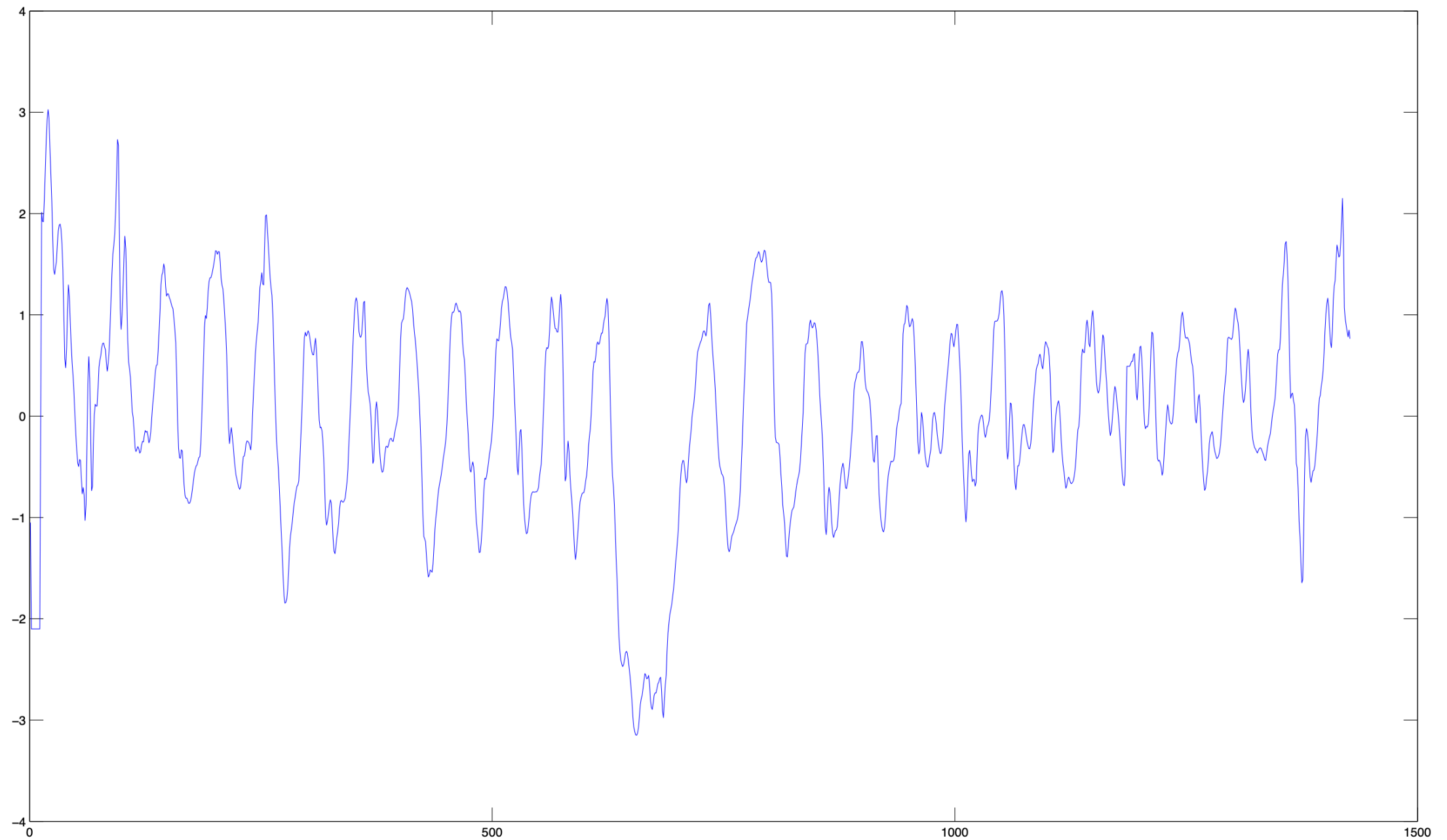# Histograms



*vs.*

# Time-Series Data



- Cleaning or preprocessing the data

- Correct representation

- Level/Magnitude
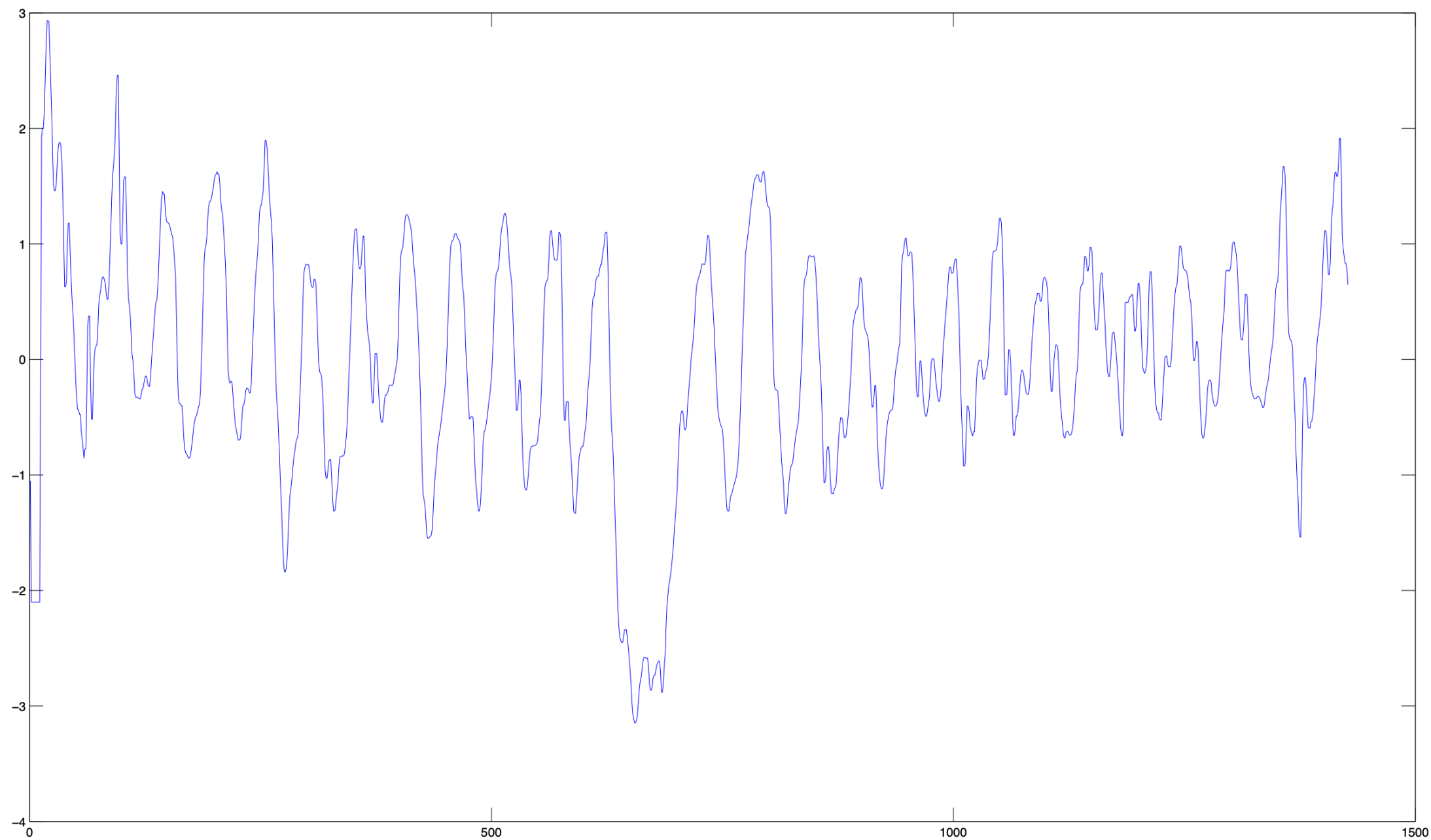
- Repetitions

- Shape of the curve

# Preprocessing
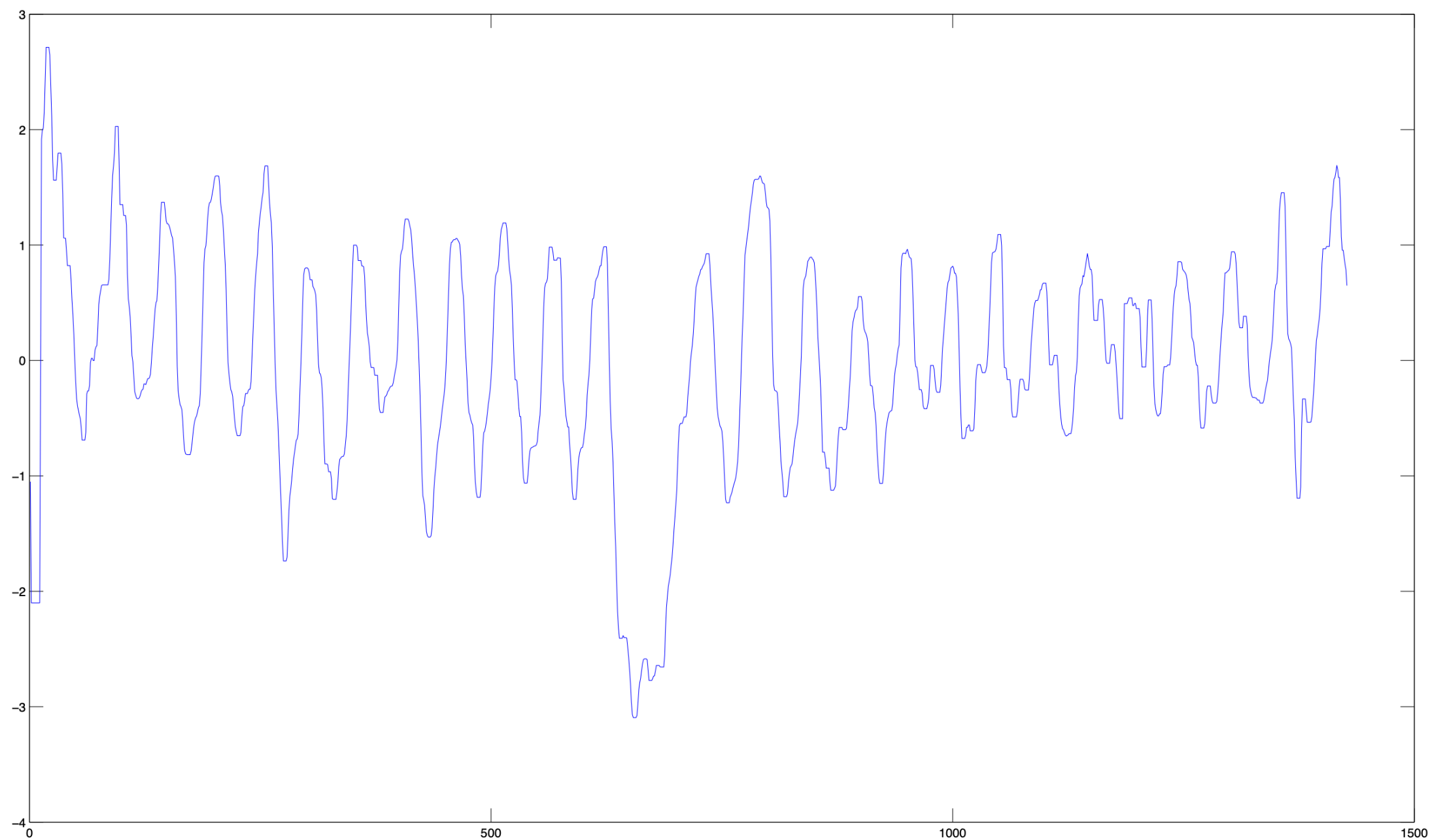


Zero Crossing on clean data

# Preprocessing
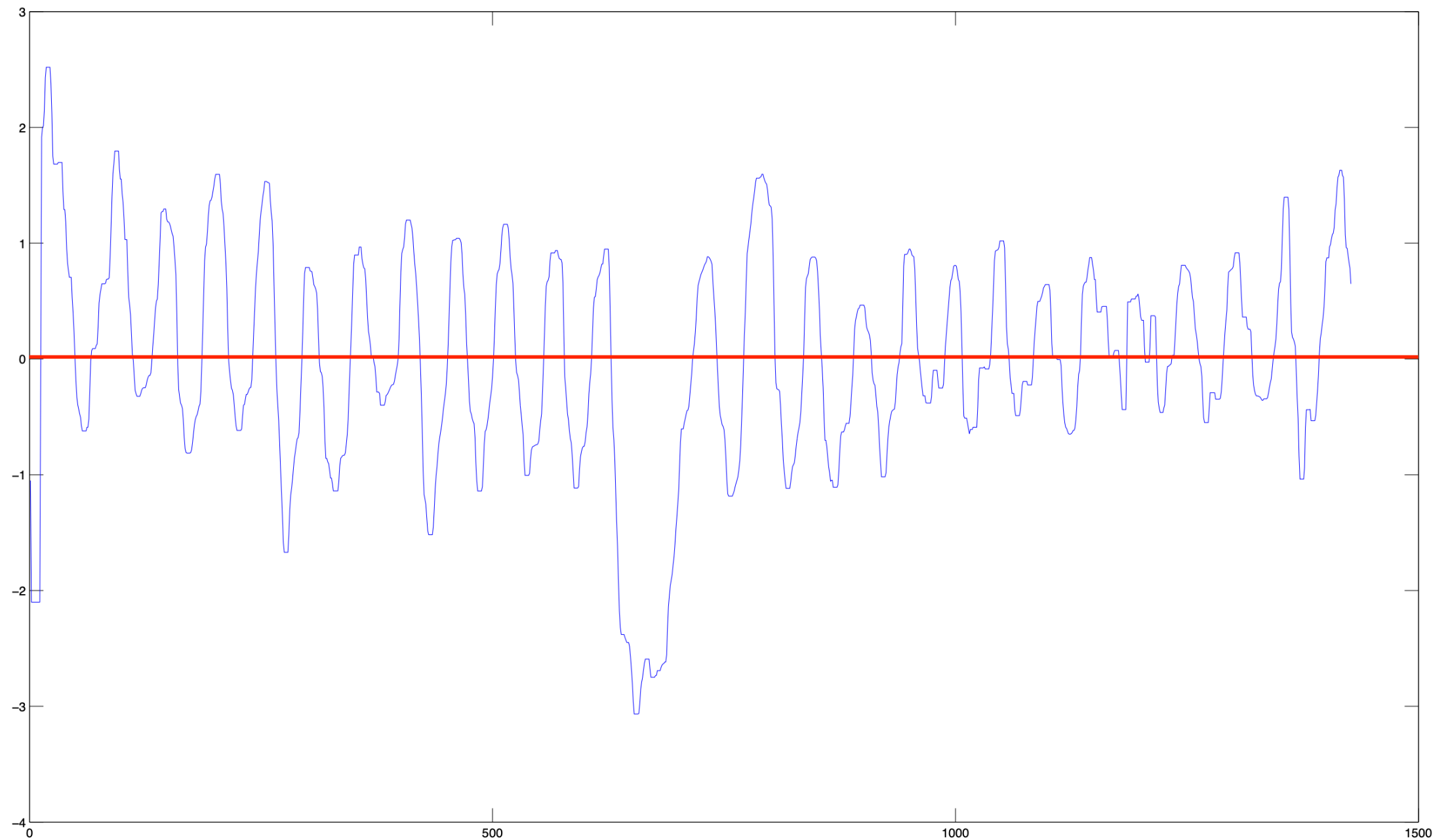


Median Filtering, size = 2

# Preprocessing



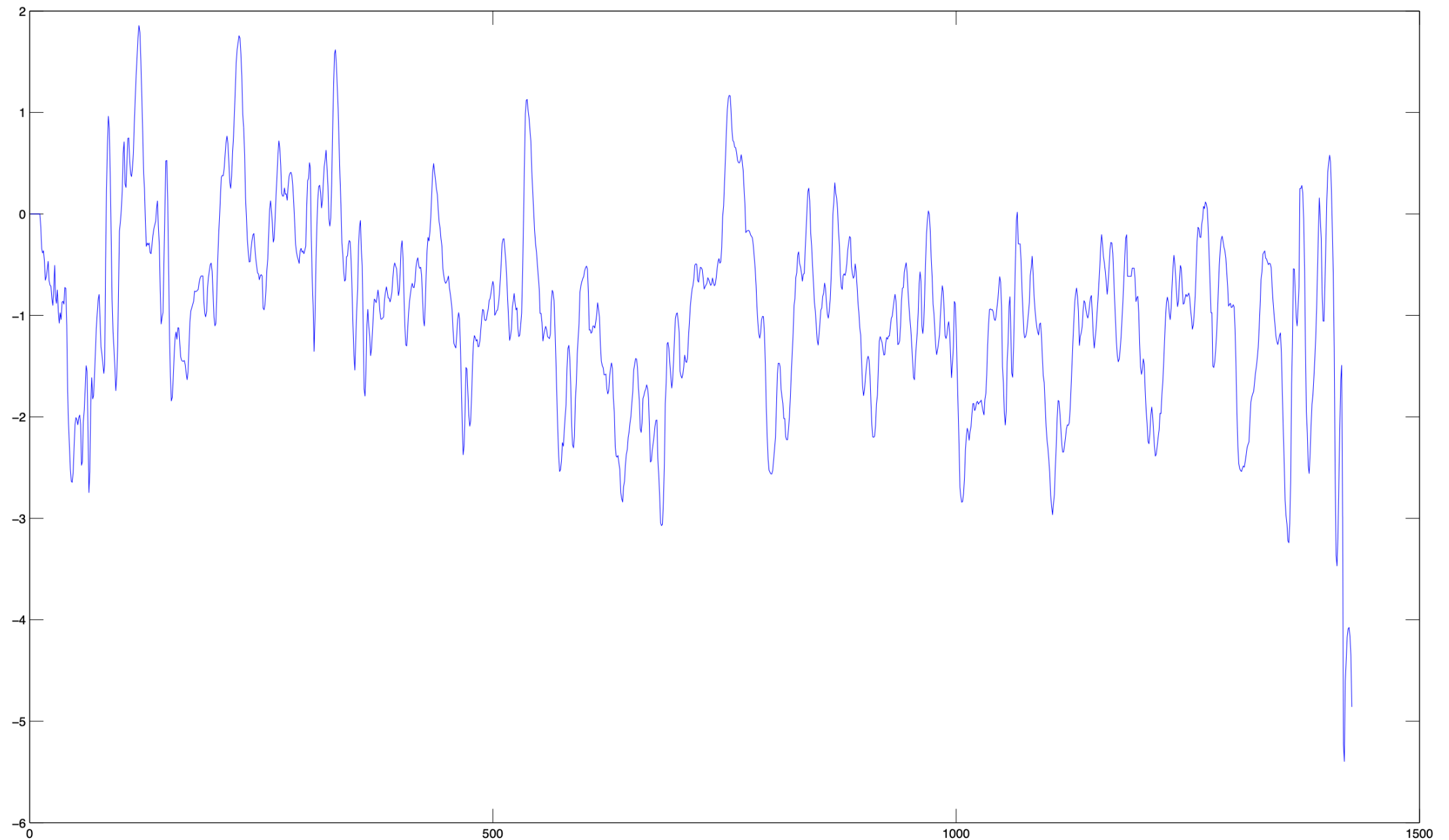Median Filtering, size = 4

# Preprocessing



Median Filtering, size = 6
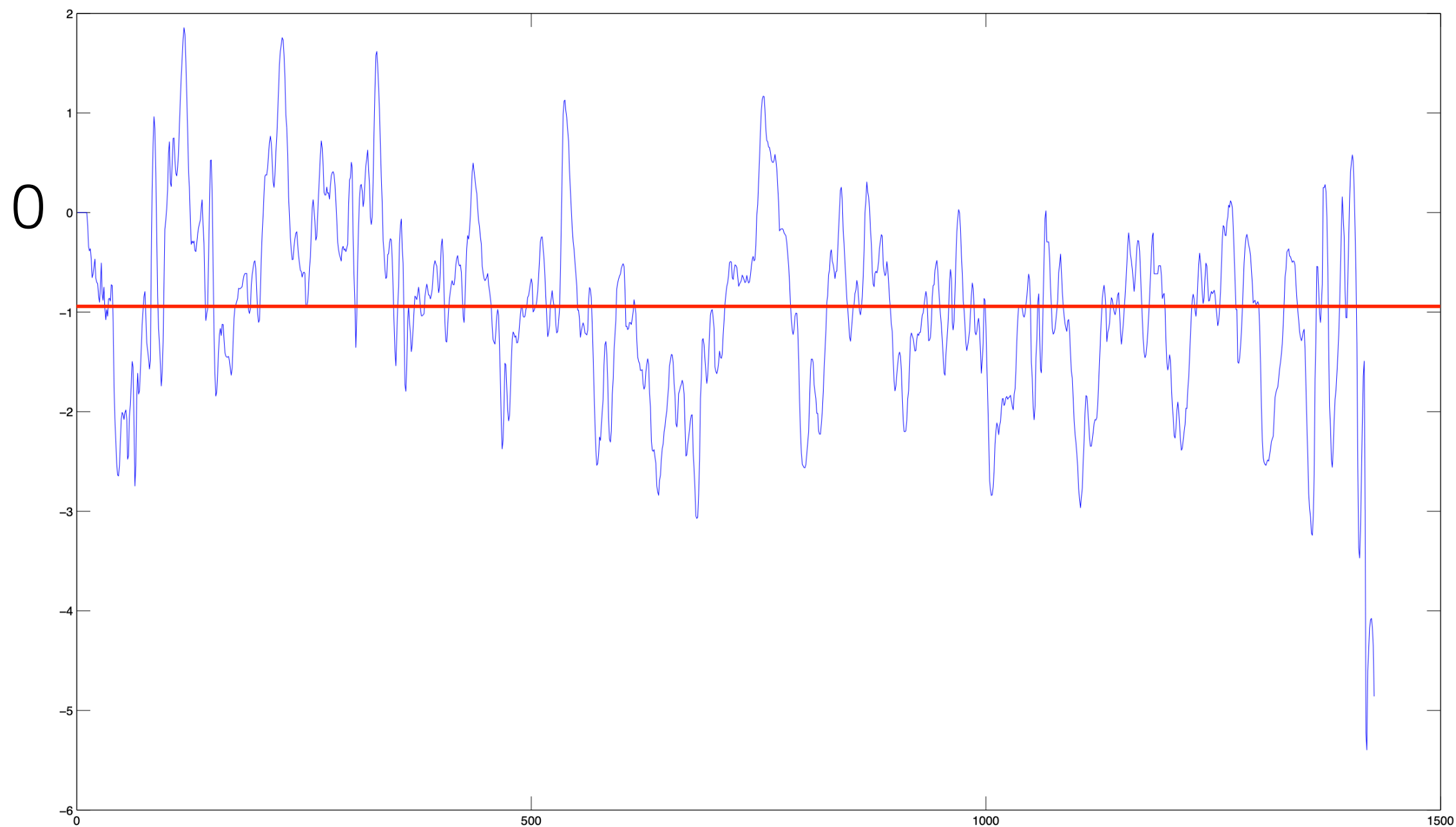
# Preprocessing



Median Filtering, size = 8
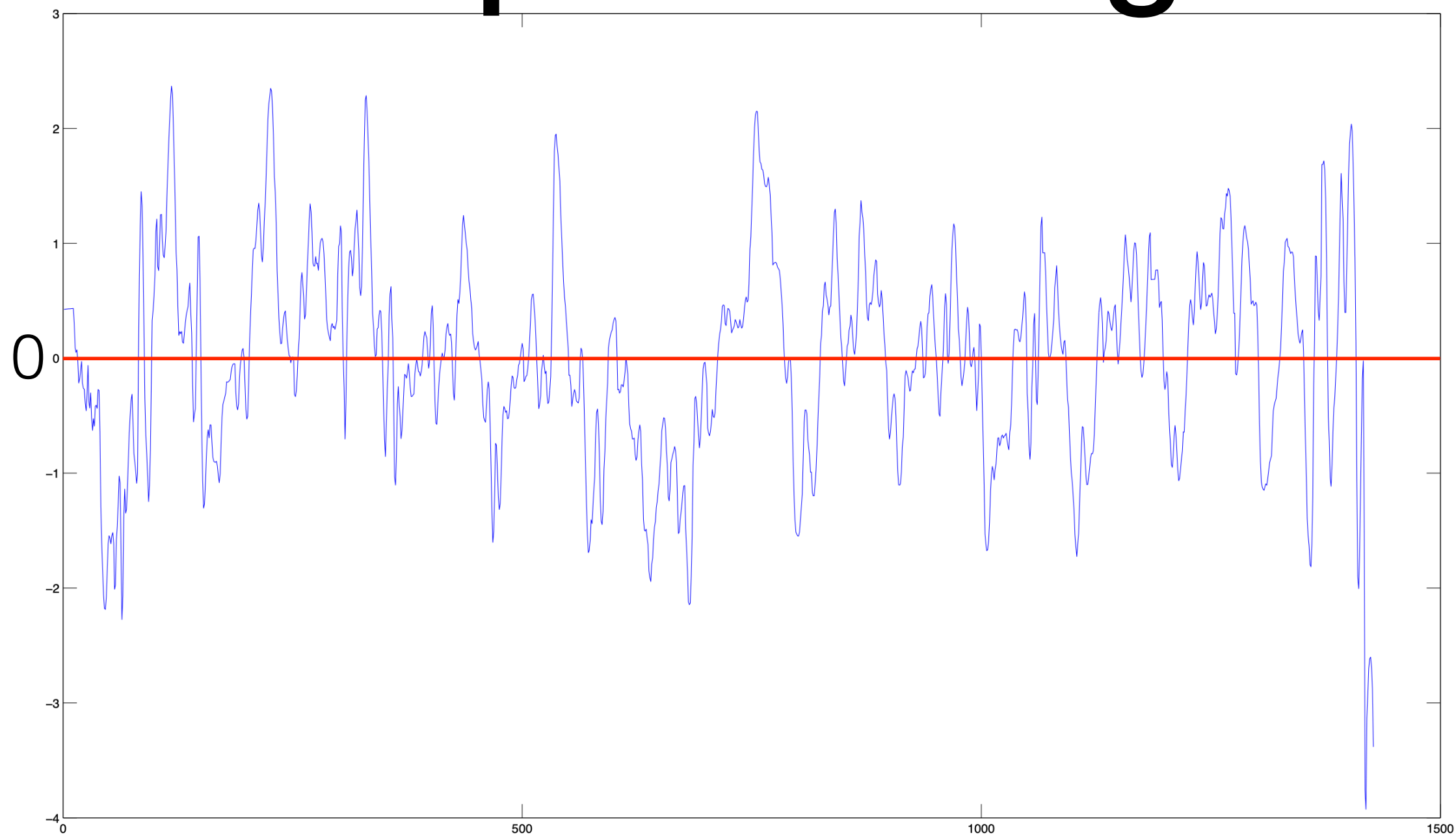
# Preprocessing



It might not always be clean

# Preprocessing


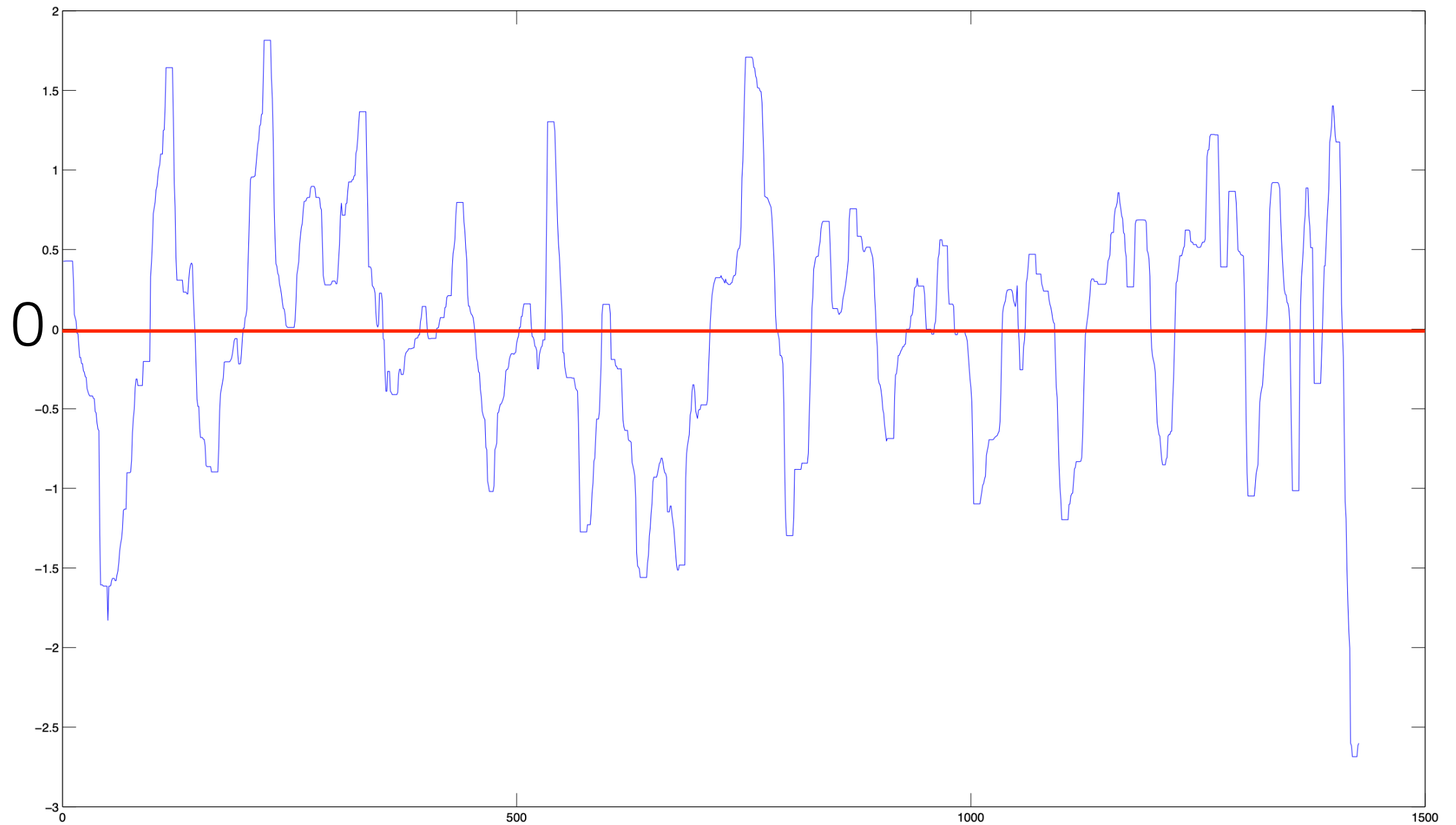
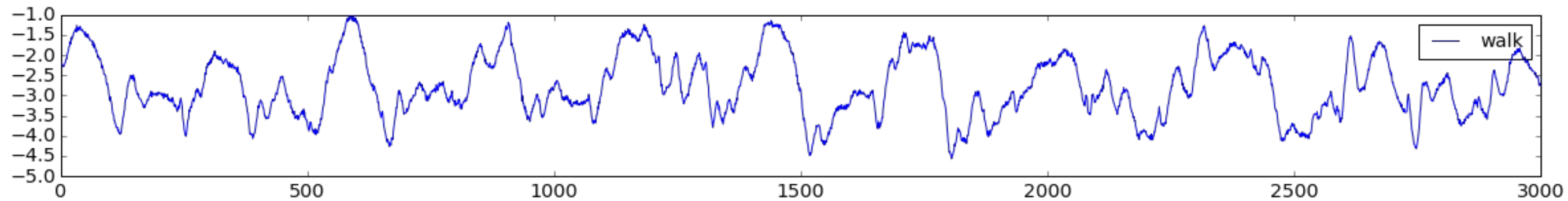Mean is not zero

# Preprocessing



After de-meaning
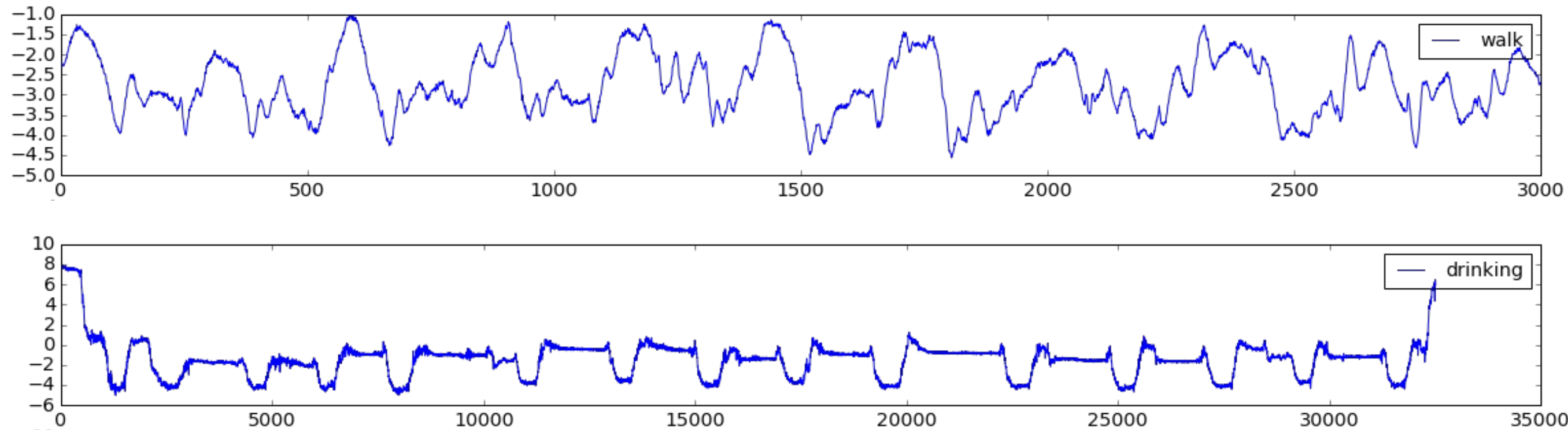
# Preprocessing



After median filtering

# Level/Magnitude



- Max

- Min

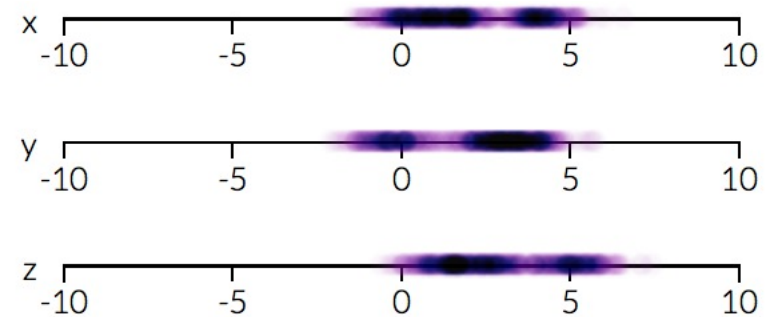- Mean
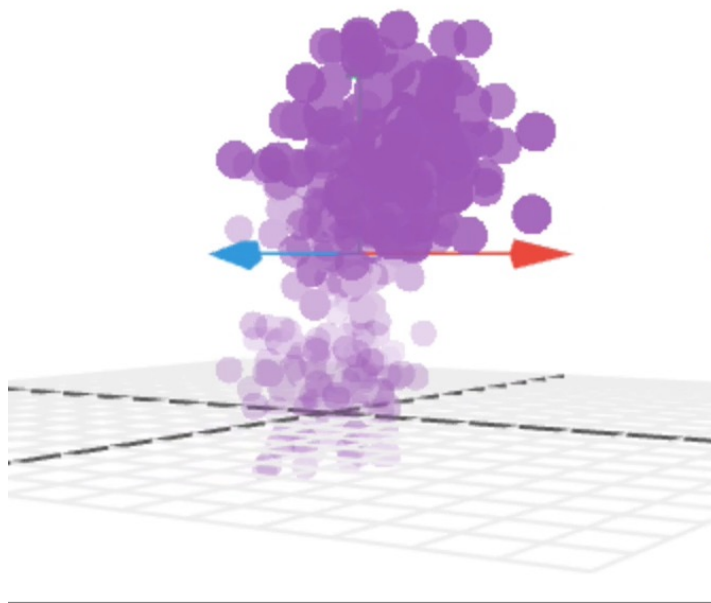
- Median

- Quantiles

# Repetitions



- Zero Crossing

- Frequency Analysis

- Auto-Correlation

# Correct Representation
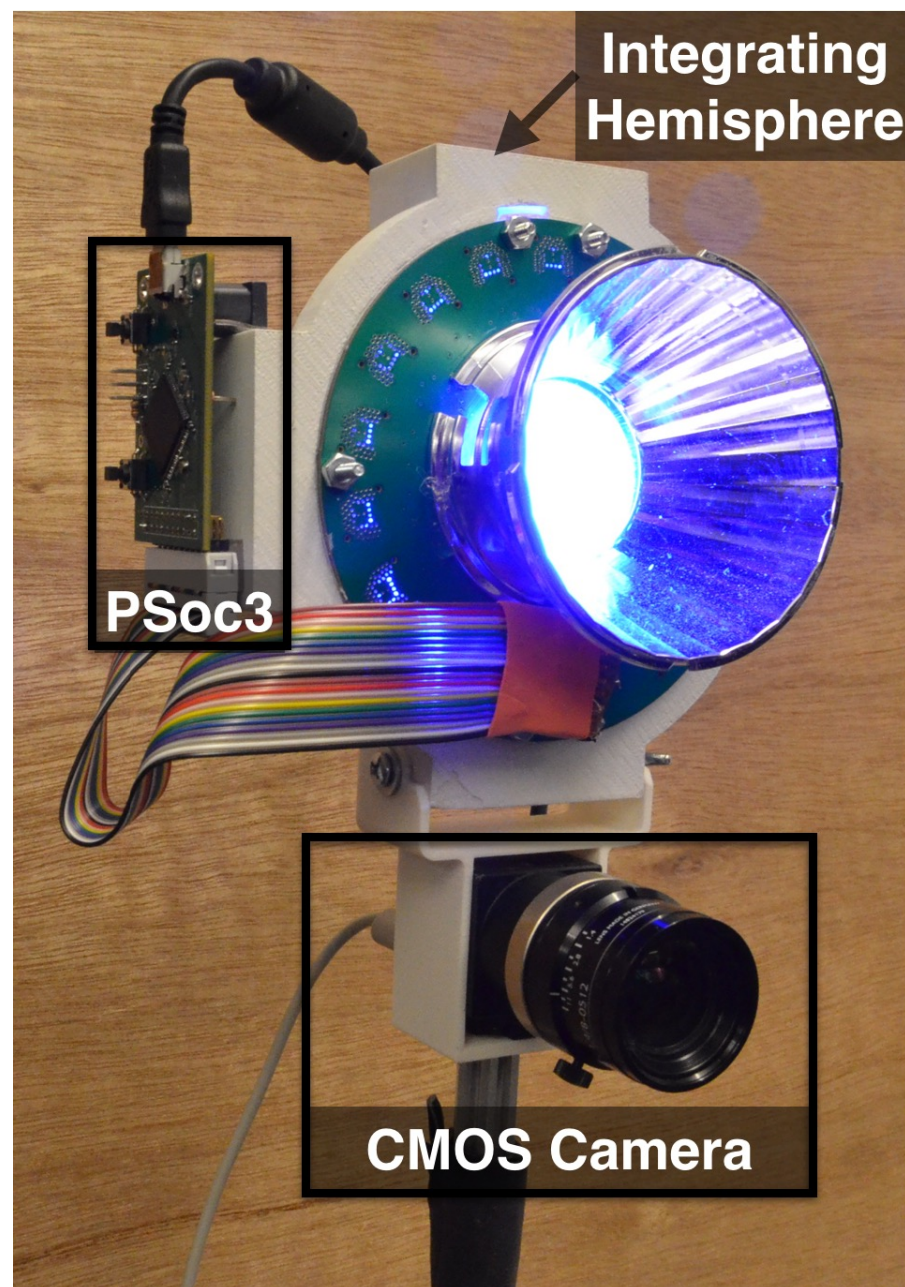*(e.g., Axes)*

- Let us look at some sensor examples

- Domain-dependent

- Magnitude

- Principal Component Analysis
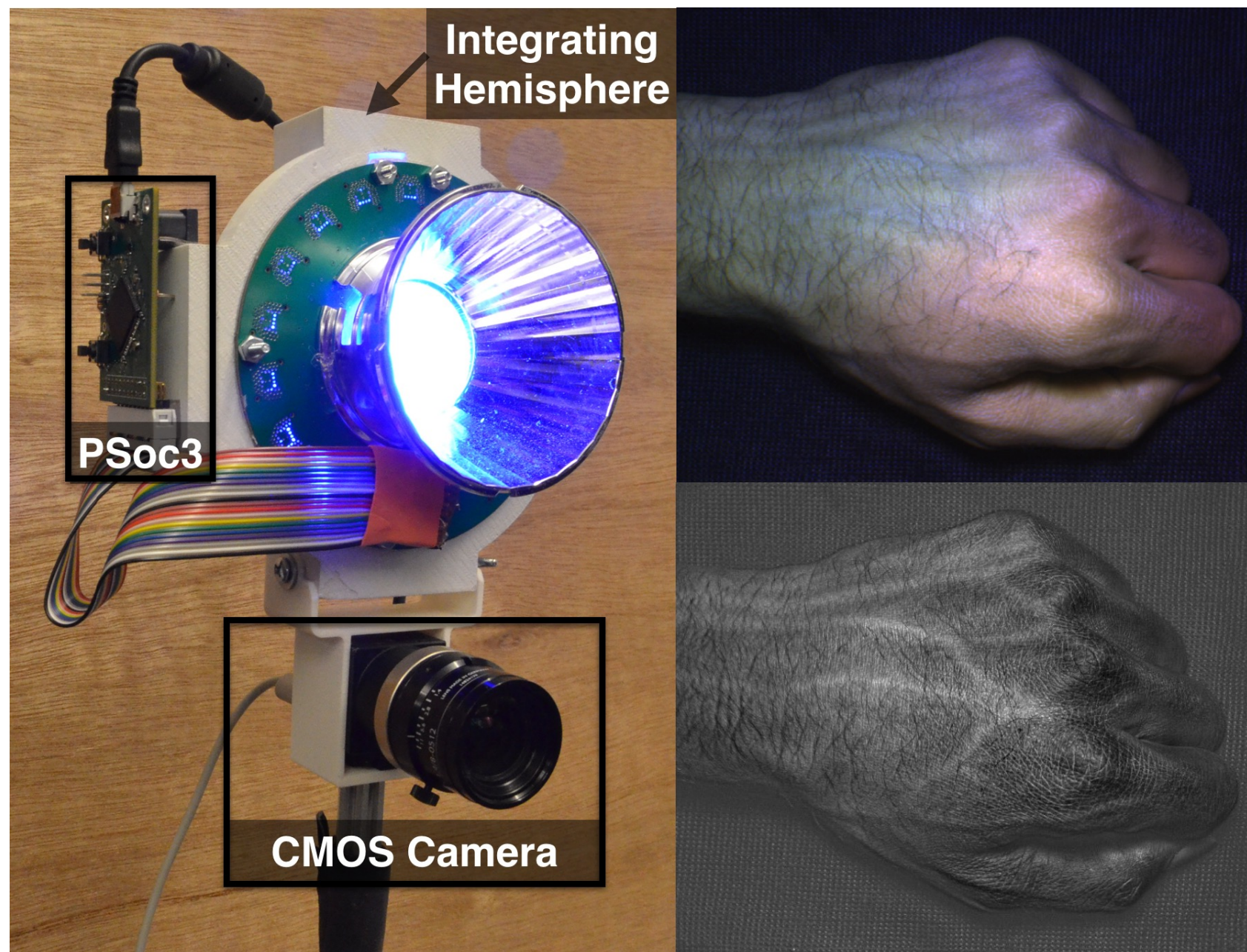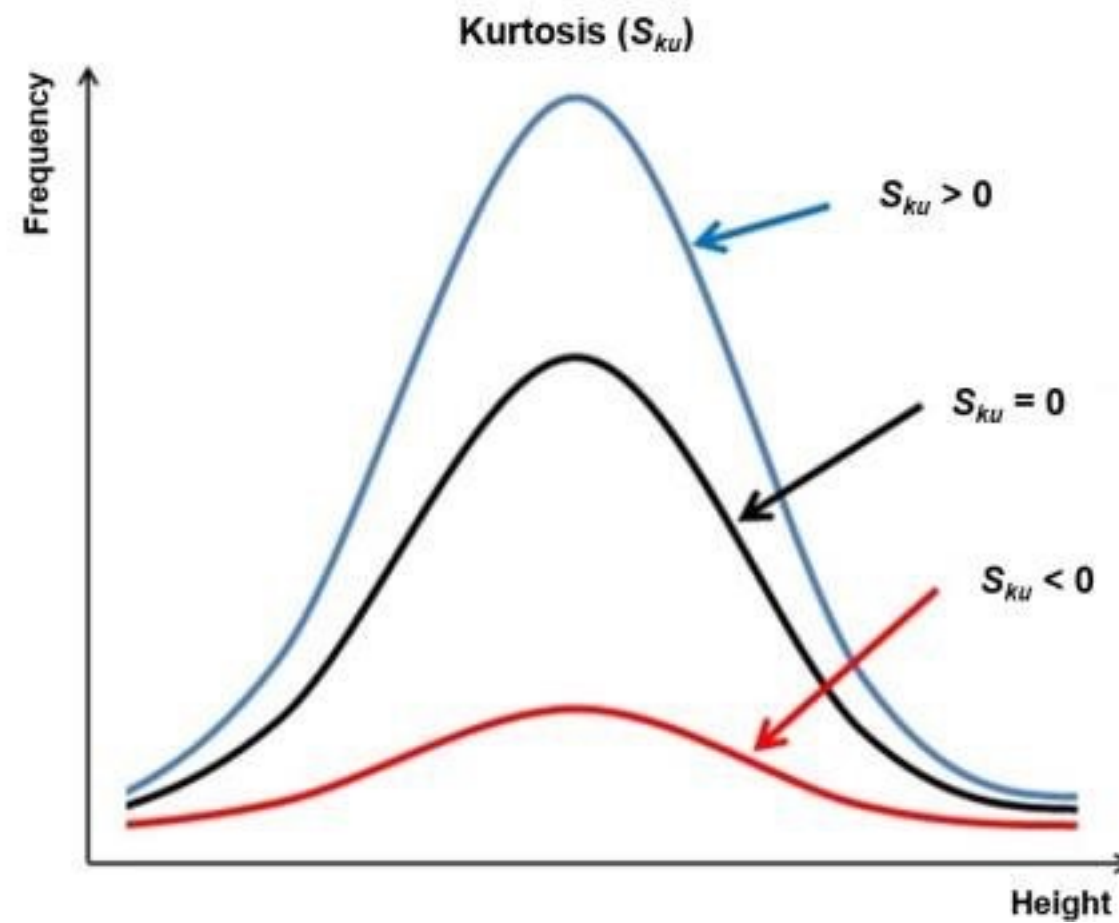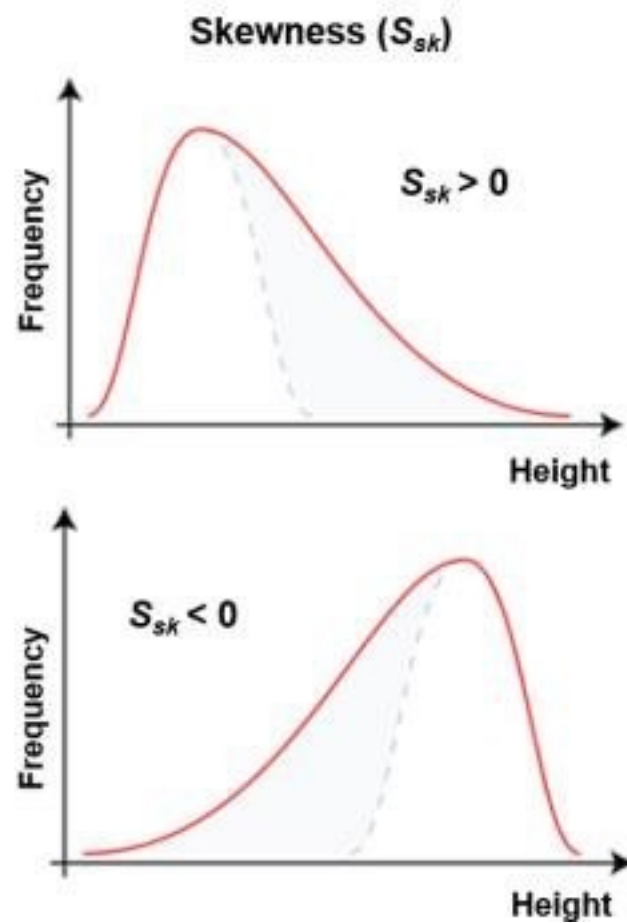
# Principal Component Analysis
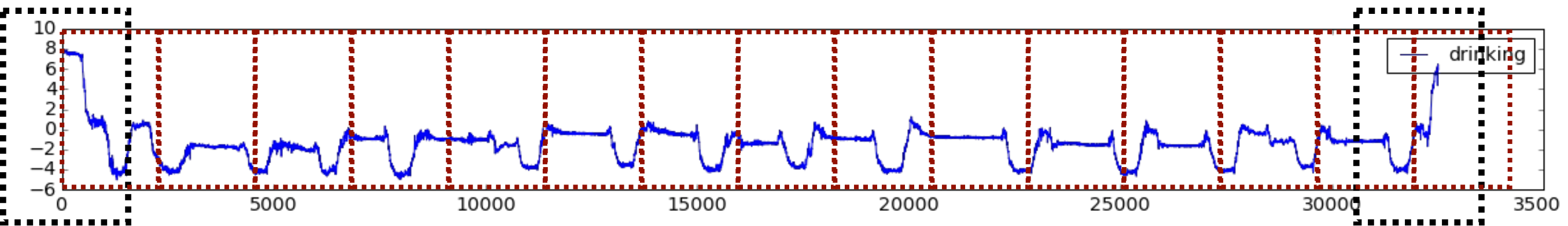
# Principal Component Analysis

# Principal Component Analysis

# Shape of the Curve

# Windowing



`Duration, Number of Peaks, Max, nth Quantile, Skewness,` **Drinking/Walking**

**2 Strategies:**

• Make a decision for each window

• Concatenate information from each window into a feature vector

`Duration, Number of Peaks, Max, nth Quantile, Skewness,` **Drinking/Walking** `for w1`
`Duration, Number of Peaks, Max, nth Quantile, Skewness,` **Drinking/Walking** `for w2`
`Duration, Number of Peaks, Max, nth Quantile, Skewness,` **Drinking/Walking** `for w3`
`Duration, Number of Peaks, Max, nth Quantile, Skewness,` **Drinking/Walking** `for w4`
`- - - - - - - - - - - - - - - - - - - -` **COMBINE** `- - - - - - - - - - - - - - - - - - - - - - -`
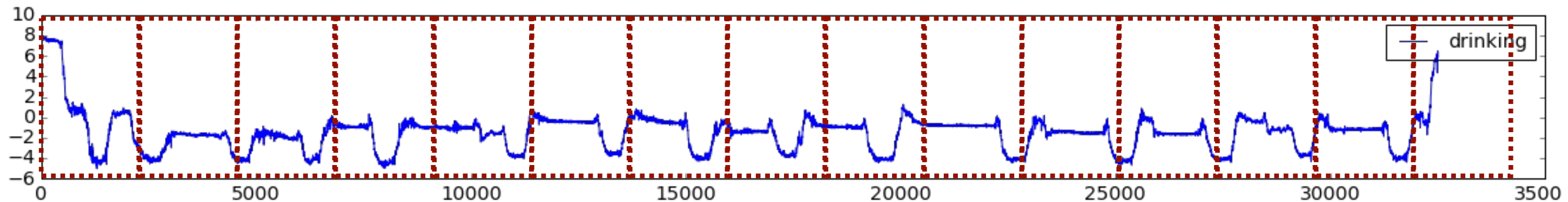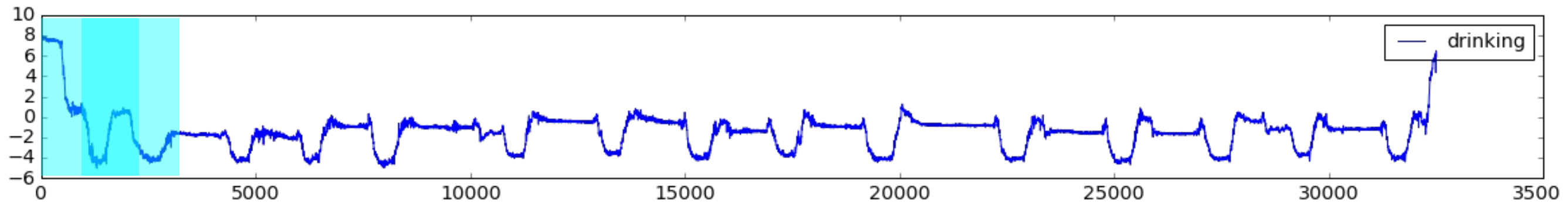
# Windowing



`Duration, Number of Peaks, Max, nth Quantile, Skewness,` **Drinking/Walking**
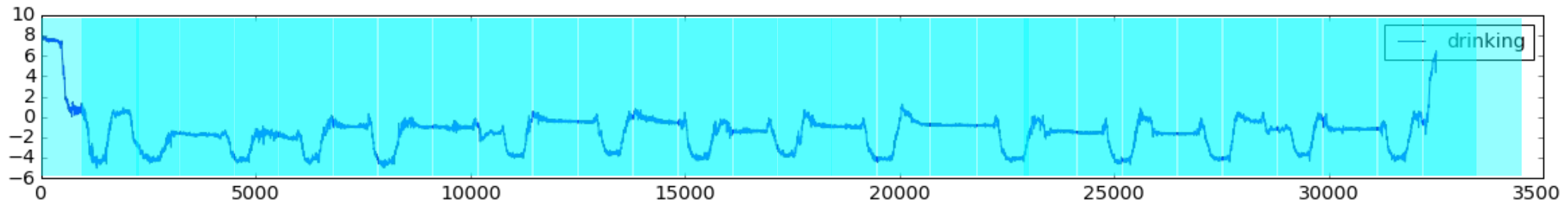
## 2 Strategies:

- Make a decision for each window
- Concatenate information from each window into a feature vector

```
(Duration , Number of Peaks, Max, nth Quantile, Skewness)for w1,
(Duration , Number of Peaks, Max, nth Quantile, Skewness)for w2,
                    and so on for rest of the windows,
```
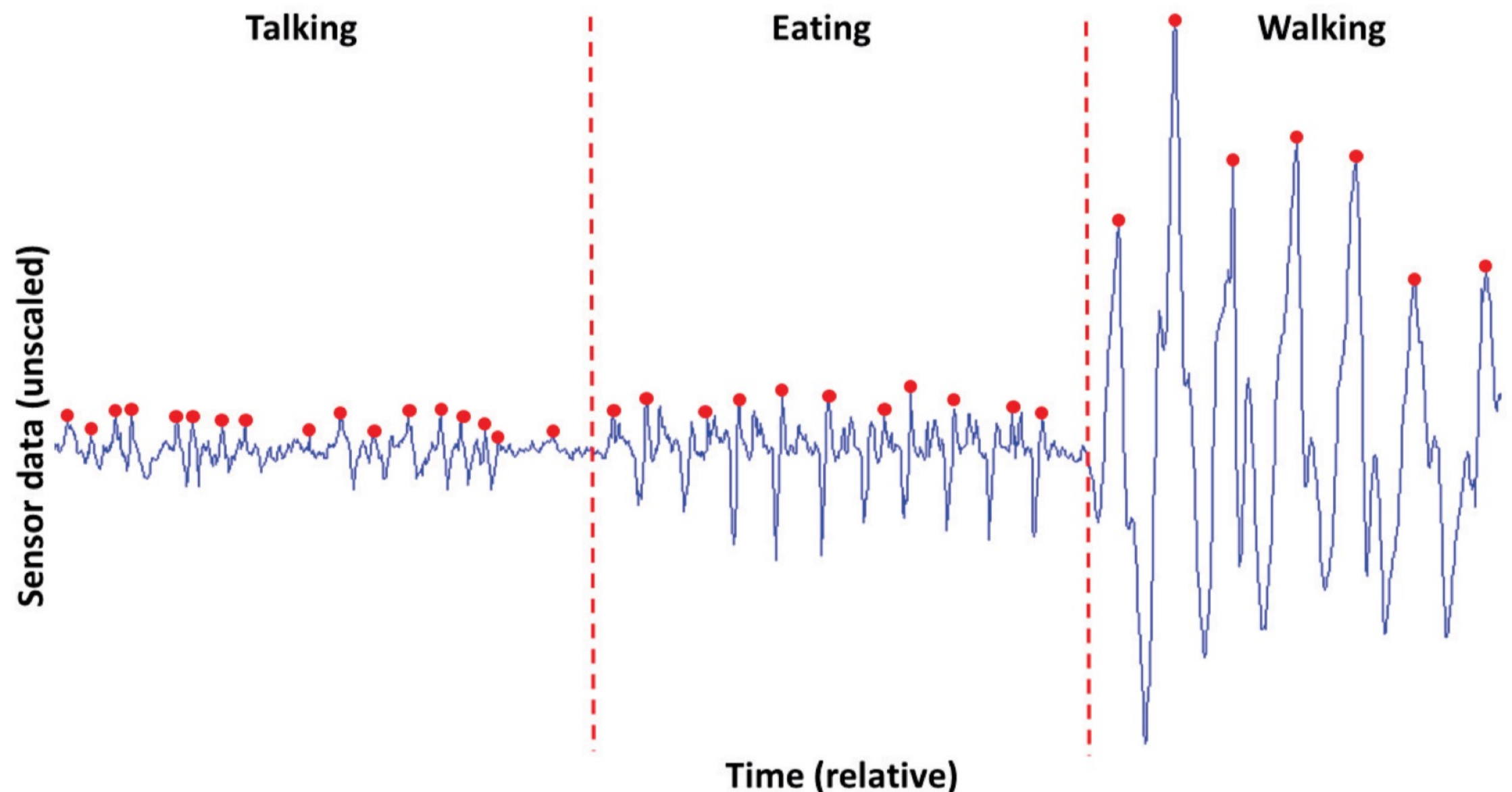single feature vector
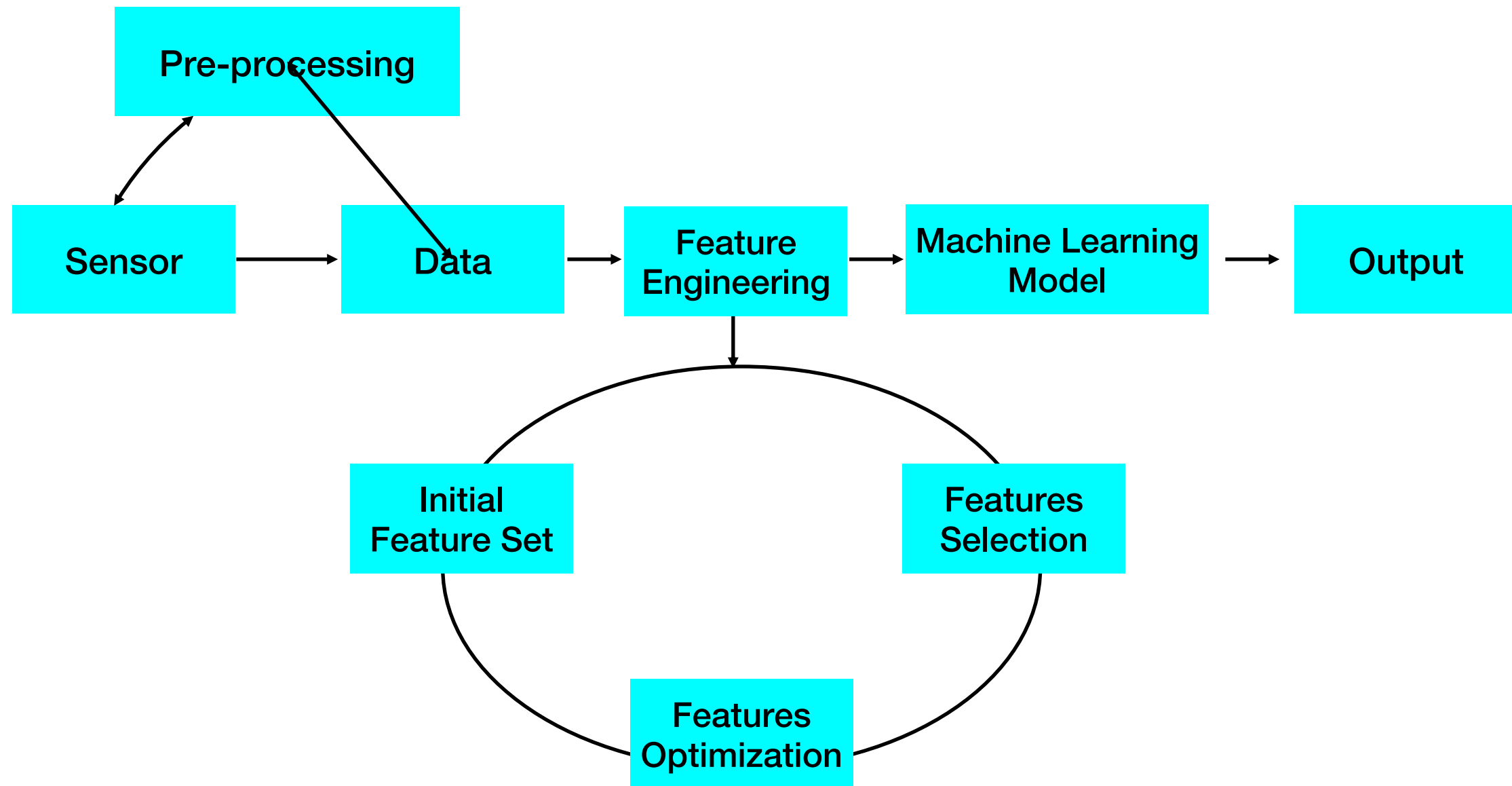
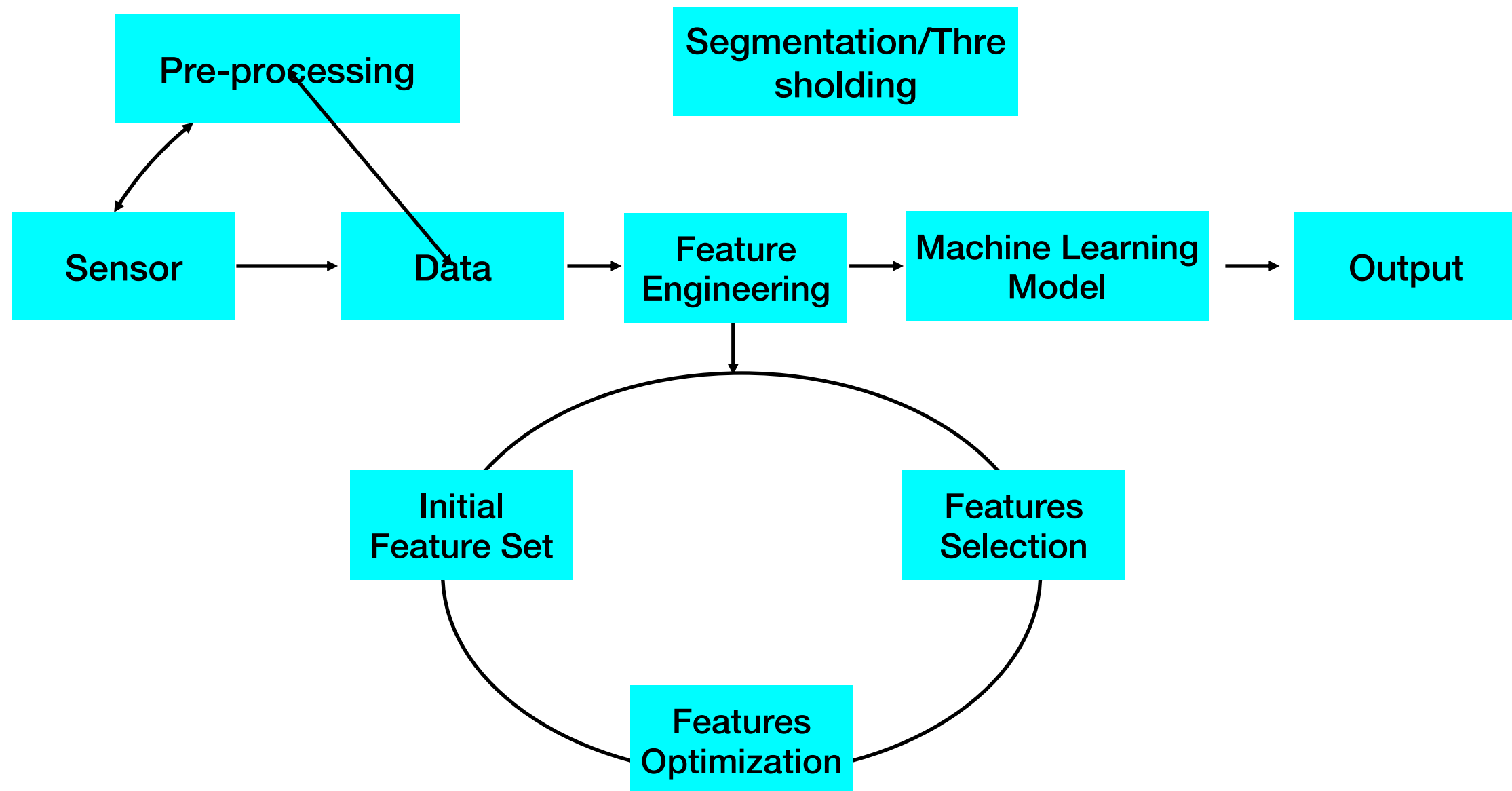**Drinking/Walking**

# Windowing

# Windowing

# Do we keep running all the data through the whole ML pipeline?



*Bedri et al.*

# Segmentation

# Segmentation

# Segmentation