

# Assignment 1 (ML for TS) - MVA 2023/2024

Gauthier Multari [gauthier.multari@ens-paris-saclay.fr](mailto:gauthier.multari@ens-paris-saclay.fr)

November 7, 2023

## 1 Introduction

**Objective.** This assignment has three parts: questions about the convolutional dictionary learning, the spectral features and a data study using the DTW.

### Warning and advice.

- Use code from the tutorials as well as from other sources. Do not code yourself well-known procedures (e.g. cross validation or k-means), use an existing implementation.
- The associated notebook contains some hints and several helper functions.
- Be concise. Answers are not expected to be longer than a few sentences (omitting calculations).

### Instructions.

- Fill in your names and emails at the top of the document.
- Hand in your report (one per pair of students) by Tuesday 7<sup>th</sup> November 23:59 PM.
- Rename your report and notebook as follows:  
FirstnameLastname1\_FirstnameLastname2.pdf and  
FirstnameLastname1\_FirstnameLastname2.ipynb.  
For instance, LaurentOudre\_CharlesTruong.pdf.
- Upload your report (PDF file) and notebook (IPYNB file) using this link:  
[docs.google.com/forms/d/e/1FAIpQLSdTwJEyc6QIoYTknjk12kJMtcKlIFvPIWLk5LbyugW0YO7K6Q/viewform?usp=sf\\_link](https://docs.google.com/forms/d/e/1FAIpQLSdTwJEyc6QIoYTknjk12kJMtcKlIFvPIWLk5LbyugW0YO7K6Q/viewform?usp=sf_link).

## 2 Convolution dictionary learning

### Question 1

Consider the following Lasso regression:

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \quad (1)$$

where  $y \in \mathbb{R}^n$  is the response vector,  $X \in \mathbb{R}^{n \times p}$  the design matrix,  $\beta \in \mathbb{R}^p$  the vector of regressors and  $\lambda > 0$  the smoothing parameter.

Show that there exists  $\lambda_{\max}$  such that the minimizer of (1) is  $\mathbf{0}_p$  (a  $p$ -dimensional vector of zeros) for any  $\lambda > \lambda_{\max}$ .

### Answer 1

If  $\lambda$  is big enough, the cost of not being sparse becomes higher than the cost of not fitting the data. The subdifferential of the loss  $L$  is the following:

$$\partial(L(\beta)) = \{X^T(y - X\beta) + \lambda z \mid z \in [-1, 1]\} \quad (2)$$

$\beta = 0$  is a solution to the problem if it belongs to the subdifferential, i.e. :

$$\beta \in \partial(L(0)) \Rightarrow X^T y = \lambda(-z), \quad z \in [-1, 1] \quad (3)$$

This is true if  $-\mathbb{1} \leq \frac{1}{\lambda} X^T y \leq \mathbb{1}$  element-wise, so we have:

$$\lambda \leq \lambda_{max} = \|X^T y\|_\infty \quad (4)$$

### Question 2

For a univariate signal  $\mathbf{x} \in \mathbb{R}^n$  with  $n$  samples, the convolutional dictionary learning task amounts to solving the following optimization problem:

$$\min_{(\mathbf{d}_k)_k, (\mathbf{z}_k)_k \|\mathbf{d}_k\|_2 \leq 1} \left\| \mathbf{x} - \sum_{k=1}^K \mathbf{z}_k * \mathbf{d}_k \right\|_2^2 + \lambda \sum_{k=1}^K \|\mathbf{z}_k\|_1 \quad (5)$$

where  $\mathbf{d}_k \in \mathbb{R}^L$  are the  $K$  dictionary atoms (patterns),  $\mathbf{z}_k \in \mathbb{R}^{N-L+1}$  are activations signals, and  $\lambda > 0$  is the smoothing parameter.

Show that

- for a fixed dictionary, the sparse coding problem is a lasso regression (explicit the response vector and the design matrix);
- for a fixed dictionary, there exists  $\lambda_{max}$  (which depends on the dictionary) such that the sparse codes are only 0 for any  $\lambda > \lambda_{max}$ .

### Answer 2

## 3 Spectral feature

Let  $X_n$  ( $n = 0, \dots, N-1$ ) be a weakly stationary random process with zero mean and autocovariance function  $\gamma(\tau) := \mathbb{E}(X_n X_{n+\tau})$ . Assume the autocovariances are absolutely summable, i.e.  $\sum_{\tau \in \mathbb{Z}} |\gamma(\tau)| < \infty$ , and square summable, i.e.  $\sum_{\tau \in \mathbb{Z}} \gamma^2(\tau) < \infty$ . Denote by  $f_s$  the sampling frequency, meaning that the index  $n$  corresponds to the time instant  $n/f_s$  and for simplicity, let  $N$  be even.

The *power spectrum*  $S$  of the stationary random process  $X$  is defined as the Fourier transform of the autocovariance function:

$$S(f) := \sum_{\tau=-\infty}^{+\infty} \gamma(\tau) e^{-2\pi f \tau / f_s}. \quad (6)$$

The power spectrum describes the distribution of power in the frequency space. Intuitively, large values of  $S(f)$  indicates that the signal contains a sine wave at the frequency  $f$ . There are many estimation procedures to determine this important quantity, which can then be used in a machine learning pipeline. In the following, we discuss about the large sample properties of simple estimation procedures, and the relationship between the power spectrum and the autocorrelation.

(Hint: use the many results on quadratic forms of Gaussian random variables to limit the amount of calculations.)

### Question 3

In this question, let  $X_n$  ( $n = 0, \dots, N - 1$ ) be a Gaussian white noise.

- Calculate the associated autocovariance function and power spectrum. (By analogy with the light, this process is called “white” because of the particular form of its power spectrum.)

### Answer 3

Let  $X_n$  a Gaussian white noise of mean  $\mu = 0$  and standard deviation  $\sigma > 0$ .

Let  $X_i$  the independant samples of the distribution, we then have:

$$\begin{aligned}\gamma(\tau) &= \mathbb{E}[X_n X_{n+\tau}] \\ &= \text{cov}(X_n, X_{n+\tau}) + \mathbb{E}[X_n] \mathbb{E}[X_{n+\tau}]\end{aligned}$$

$$\boxed{\gamma(\tau) = \begin{cases} 0, & \text{if } \tau \neq 0. \\ \sigma^2, & \text{otherwise.} \end{cases}} \quad (7)$$

Power spectrum :

$$\begin{aligned}S(f) &= \sum_{\tau=-\infty}^{+\infty} \gamma(\tau) e^{-2\pi f \tau / f_s} \\ &= \gamma(0) e^0 + \sum_{\tau=-\infty}^{-1} \gamma(\tau) e^{-2\pi f \tau / f_s} + \sum_{\tau=1}^{+\infty} \gamma(\tau) e^{-2\pi f \tau / f_s} \\ &= \gamma(0)\end{aligned}$$

$$\boxed{S(f) = \sigma^2} \quad (8)$$

### Question 4

A natural estimator for the autocorrelation function is the sample autocovariance

$$\hat{\gamma}(\tau) := (1/N) \sum_{n=0}^{N-\tau-1} X_n X_{n+\tau} \quad (9)$$

for  $\tau = 0, 1, \dots, N - 1$  and  $\hat{\gamma}(\tau) := \hat{\gamma}(-\tau)$  for  $\tau = -(N - 1), \dots, -1$ .

- Show that  $\hat{\gamma}(\tau)$  is a biased estimator of  $\gamma(\tau)$  but asymptotically unbiased. What would be a simple way to de-bias this estimator?

#### Answer 4

To check if the estimator is biased, we can compute the average sample autocovariance:

$$\begin{aligned}
 \mathbb{E}[\hat{\gamma}_N(\tau)] &= \mathbb{E}\left[\frac{1}{N} \sum_{n=0}^{N-\tau-1} X_n X_{n+\tau}\right] \\
 &= \frac{1}{N} \sum_{n=0}^{N-\tau-1} \mathbb{E}[X_n X_{n+\tau}] \\
 &= \frac{1}{N} \sum_{n=0}^{N-\tau-1} \gamma(\tau) \\
 &= \frac{N-\tau}{N} \gamma(\tau) \neq \gamma(\tau)
 \end{aligned}$$

$\hat{\gamma}_N(\tau)$  is indeed biased, but we can simply de-bias it by multiplying it by  $\frac{N-\tau}{N}$

#### Question 5

Define the discrete Fourier transform of the random process  $\{X_n\}_n$  by

$$J(f) := (1/\sqrt{N}) \sum_{n=0}^{N-1} X_n e^{-2\pi i f n / f_s} \quad (10)$$

The *periodogram* is the collection of values  $|J(f_0)|^2, |J(f_1)|^2, \dots, |J(f_{N/2})|^2$  where  $f_k = f_s k / N$ . (They can be efficiently computed using the Fast Fourier Transform.)

- Write  $|J(f_k)|^2$  as a function of the sample autocovariances.
- For a frequency  $f$ , define  $f^{(N)}$  the closest Fourier frequency  $f_k$  to  $f$ . Show that  $|J(f^{(N)})|^2$  is an asymptotically unbiased estimator of  $S(f)$  for  $f > 0$ .

#### Answer 5

Write  $|J(f_k)|^2$  as a function of the sample autocovariances :

$$\begin{aligned}
 |J(f_k)|^2 &= \left| \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} X_n e^{-\frac{2\pi i f_k n}{f_s}} \right|^2 \\
 &= \frac{1}{N} \left| \sum_{n=0}^{N-1} X_n e^{-\frac{2\pi i k n}{N}} \right|^2 \\
 &= \frac{1}{N} \left( \sum_{n=0}^{N-1} X_n e^{-\frac{2\pi i k n}{N}} \right) \left( \sum_{m=0}^{N-1} X_m e^{\frac{2\pi i k m}{N}} \right) \\
 &= \frac{1}{N} \left( \sum_{n=0}^{N-1} \sum_{m=0}^{N-1} X_n X_m e^{\frac{2\pi i k (m-n)}{N}} \right) \\
 &= \frac{1}{N} \sum_{\tau=-(N-1)}^{N-1} N \hat{\gamma}(\tau) e^{-\frac{2\pi i k \tau}{N}}
 \end{aligned}$$

Finally, we have:

$$|J(f_k)|^2 = \frac{1}{N} \sum_{\tau=-(N-1)}^{N-1} N \hat{\gamma}(\tau) e^{-\frac{2\pi i k \tau}{N}} \quad (11)$$

For a frequency  $f$ , define  $f^{(N)}$  the closest Fourier frequency  $f_k$  to  $f$ :

$$f^{(N)} = \arg \min_{f_k} |f - f_k| \quad (12)$$

### Question 6

In this question, let  $X_n$  ( $n = 0, \dots, N-1$ ) be a Gaussian white noise with variance  $\sigma^2 = 1$  and set the sampling frequency to  $f_s = 1$  Hz

- For  $N \in \{200, 500, 1000\}$ , compute the *sample autocovariances* ( $\hat{\gamma}(\tau)$  vs  $\tau$ ) for 100 simulations of  $X$ . Plot the average value as well as the average  $\pm$  the standard deviation. What do you observe?
- For  $N \in \{200, 500, 1000\}$ , compute the *periodogram* ( $|J(f_k)|^2$  vs  $f_k$ ) for 100 simulations of  $X$ . Plot the average value as well as the average  $\pm$  the standard deviation. What do you observe?

Add your plots to Figure 1.

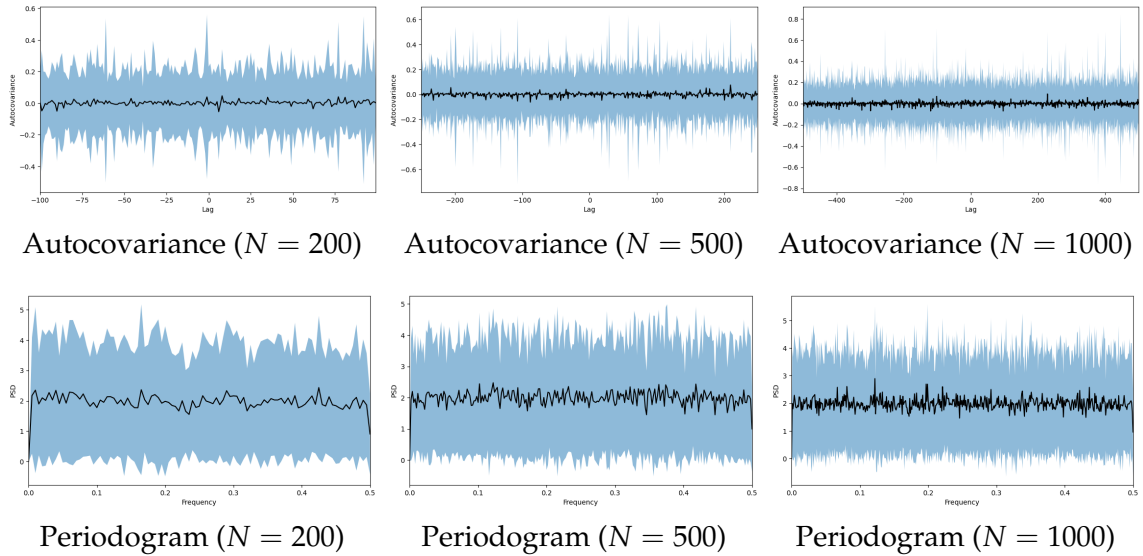


Figure 1: Autocovariances and periodograms of a Gaussian white noise (see Question 6).

### Answer 6

The standard deviation as well as the mean of all signals don't change over time.

## Question 7

We want to show that the estimator  $\hat{\gamma}(\tau)$  is consistent, i.e. it converges in probability when the number  $N$  of samples grows to  $\infty$  to the true value  $\gamma(\tau)$ . In this question, assume that  $X$  is a wide-sense stationary *Gaussian* process.

- Show that for  $\tau > 0$

$$\text{var}(\hat{\gamma}(\tau)) = (1/N) \sum_{n=-(N-\tau-1)}^{n=N-\tau-1} \left(1 - \frac{\tau + |n|}{N}\right) [\gamma^2(n) + \gamma(n-\tau)\gamma(n+\tau)]. \quad (13)$$

(Hint: if  $\{Y_1, Y_2, Y_3, Y_4\}$  are four centered jointly Gaussian variables, then  $\mathbb{E}[Y_1 Y_2 Y_3 Y_4] = \mathbb{E}[Y_1 Y_2] \mathbb{E}[Y_3 Y_4] + \mathbb{E}[Y_1 Y_3] \mathbb{E}[Y_2 Y_4] + \mathbb{E}[Y_1 Y_4] \mathbb{E}[Y_2 Y_3]$ .)

- Conclude that  $\hat{\gamma}(\tau)$  is consistent.

## Answer 7

$$\begin{aligned} \text{Var}(\hat{\gamma}(\tau)) &= \mathbb{E}[\hat{\gamma}(\tau)^2] - \mathbb{E}[\hat{\gamma}(\tau)]^2 \\ &= \frac{1}{N^2} \mathbb{E} \left[ \left( \sum_{n=0}^{N-\tau-1} X_n X_{n+\tau} \right)^2 \right] - \frac{1}{N^2} \left( \sum_{n=0}^{N-\tau-1} \gamma(\tau) \right)^2 \\ &= \frac{1}{N^2} \left[ \sum_{n=0}^{N-\tau-1} \sum_{m=0}^{N-\tau-1} X_n X_{n+\tau} X_m X_{m+\tau} \right] - \frac{1}{N^2} \sum_{n=0}^{N-\tau-1} \sum_{m=0}^{N-\tau-1} \gamma(\tau)^2 \\ &= \frac{1}{N^2} \sum_{n=0}^{N-\tau-1} \sum_{m=0}^{N-\tau-1} (\gamma(\tau)^2 + \gamma(m-n)^2 + \gamma(m+\tau-n)\gamma(m-n-\tau)) - \frac{1}{N^2} \sum_{n=0}^{N-\tau-1} \sum_{m=0}^{N-\tau-1} \gamma(\tau)^2 \\ &= \frac{1}{N^2} \sum_{n=0}^{N-\tau-1} \sum_{m=0}^{N-\tau-1} (\gamma(m)^2 + \gamma(m+\tau)\gamma(m-\tau)) \\ &= \frac{1}{N^2} \sum_{m=-(N-\tau-1)}^{N-\tau-1} (N-\tau-|m|)(\gamma(m)^2 + \gamma(m+\tau)\gamma(m-\tau)) \end{aligned}$$

The sample autocovariance is consistent if  $\lim_{N \rightarrow \infty} \text{Var}(\hat{\gamma}(\tau)) = 0$ :

$$\begin{aligned} \text{Var}(\hat{\gamma}(\tau)) &= \frac{1}{N^2} \sum_{m=-(N-\tau-1)}^{N-\tau-1} (N-\tau-|m|)(\gamma(m)^2 + \gamma(m+\tau)\gamma(m-\tau)) \\ &\leq \frac{1}{N^2} \sum_{m=-(N-\tau-1)}^{N-\tau-1} N(\gamma(m)^2 + \gamma(m+\tau)\gamma(m-\tau)) \\ &= \frac{1}{N^2} \left( \sum_{m=-N}^N \gamma(m)^2 + \left( \sum_{m=-N}^N \gamma(m) \right)^2 \right) \xrightarrow{N \rightarrow +\infty} 0 \end{aligned}$$

$\text{Var}(\hat{\gamma}(\tau)) \xrightarrow{N \rightarrow +\infty} 0$ , the sample autocovariance is consistent.

Contrary to the correlogram, the periodogram is not consistent. It is one of the most well-known estimators that is asymptotically unbiased but not consistent. In the following question, this is proven for a Gaussian white noise but this holds for more general stationary processes.

### Question 8

Assume that  $X$  is a Gaussian white noise (variance  $\sigma^2$ ) and let  $A(f) := \sum_{n=0}^{N-1} X_n \cos(-2\pi f n / f_s)$  and  $B(f) := \sum_{n=0}^{N-1} X_n \sin(-2\pi f n / f_s)$ . Observe that  $J(f) = (1/N)(A(f) + iB(f))$ .

- Derive the mean and variance of  $A(f)$  and  $B(f)$  for  $f = f_0, f_1, \dots, f_{N/2}$  where  $f_k = f_s k / N$ .
- What is the distribution of the periodogram values  $|J(f_0)|^2, |J(f_1)|^2, \dots, |J(f_{N/2})|^2$ .
- What is the variance of the  $|J(f_k)|^2$ ? Conclude that the periodogram is not consistent.
- Explain the erratic behavior of the periodogram in Question 6 by looking at the covariance between the  $|J(f_k)|^2$ .

### Answer 8

*Sinus* and *cosinus* don't change the mean of the gaussian, so:  $\mathbb{E}[A(f)] = \mathbb{E}[B(f)] = 0$

$$\begin{aligned}
 \text{Var}(A(f_k)) &= \mathbb{E}[A(f_k)^2] \\
 &= \mathbb{E} \left[ \sum_{n=0}^{N-1} \sum_{m=0}^{N-1} X_n X_m \cos \frac{2\pi k n}{N} \cos \frac{2\pi k m}{N} \right] \\
 &= \sum_{n=0}^{N-1} \mathbb{E}[X_n^2] \mathbb{E} \left[ \cos \left( \frac{2\pi k n}{N} \right)^2 \right] \\
 &= \sigma^2 \sum_{n=0}^{N-1} \cos^2 \left( \frac{2\pi k n}{N} \right) \\
 &= \frac{\sigma^2}{2} \left( N + \text{Re} \left( \sum_{n=0}^{N-1} e^{\frac{4i\pi k n}{N}} \right) \right) \\
 &= \frac{\sigma^2 N}{2}
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(B(f_k)) &= \mathbb{E}[B(f_k)^2] \\
 &= \sigma^2 \sum_{n=0}^{N-1} \sin^2 \left( \frac{2\pi k n}{N} \right) \\
 &= \sigma^2 \sum_{n=0}^{N-1} 1 - \cos^2 \left( \frac{2\pi k n}{N} \right) \\
 &= \frac{\sigma^2 N}{2}
 \end{aligned}$$

A and B have the same mean and variance.

## Question 9

As seen in the previous question, the problem with the periodogram is the fact that its variance does not decrease with the sample size. A simple procedure to obtain a consistent estimate is to divide the signal in  $K$  sections of equal durations, compute a periodogram on each section and average them. Provided the sections are independent, this has the effect of dividing the variance by  $K$ . This procedure is known as Bartlett's procedure.

- Rerun the experiment of Question 6, but replace the periodogram by Bartlett's estimate (set  $K = 5$ ). What do you observe.

Add your plots to Figure 2.

## Answer 9

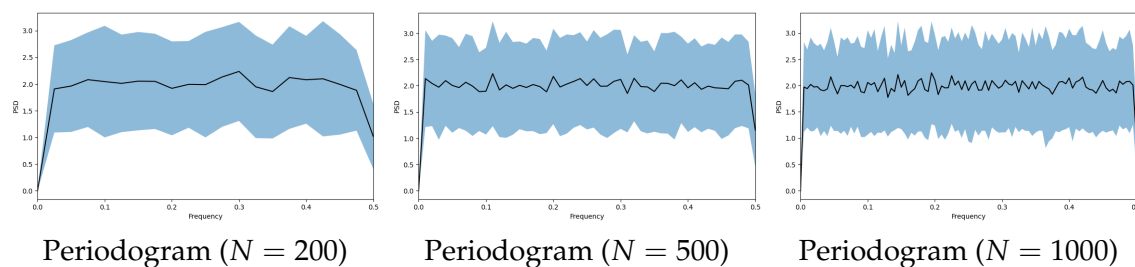


Figure 2: Bartlett's periodograms of a Gaussian white noise (see Question 9).

## 4 Data study

### 4.1 General information

**Context.** The study of human gait is a central problem in medical research with far-reaching consequences in the public health domain. This complex mechanism can be altered by a wide range of pathologies (such as Parkinson's disease, arthritis, stroke,...), often resulting in a significant loss of autonomy and an increased risk of fall. Understanding the influence of such medical disorders on a subject's gait would greatly facilitate early detection and prevention of those possibly harmful situations. To address these issues, clinical and bio-mechanical researchers have worked to objectively quantify gait characteristics.

Among the gait features that have proved their relevance in a medical context, several are linked to the notion of step (step duration, variation in step length, etc.), which can be seen as the core atom of the locomotion process. Many algorithms have therefore been developed to automatically (or semi-automatically) detect gait events (such as heel-strikes, heel-off, etc.) from accelerometer and gyrometer signals.

**Data.** Data are described in the associated notebook.

### 4.2 Step classification with the dynamic time warping (DTW) distance

**Task.** The objective is to classify footsteps then walk signals between healthy and non-healthy.



**Performance metric.** The performance of this binary classification task is measured by the F-score.

### **Question 10**

Combine the DTW and a k-neighbors classifier to classify each step. Find the optimal number of neighbors with 5-fold cross-validation and report the optimal number of neighbors and the associated F-score. Comment briefly.

### **Answer 10**

The optimal number of neighbors is 5 with an associated F-score of 0.782 on the validation set and 0.513 on the test set. The classifier doesn't generalize well on new data. This could be due to a lack of data or that this approach is not fit for this task.

### Question 11

Display on Figure 3 a badly classified step from each class (healthy/non-healthy).

### Answer 11

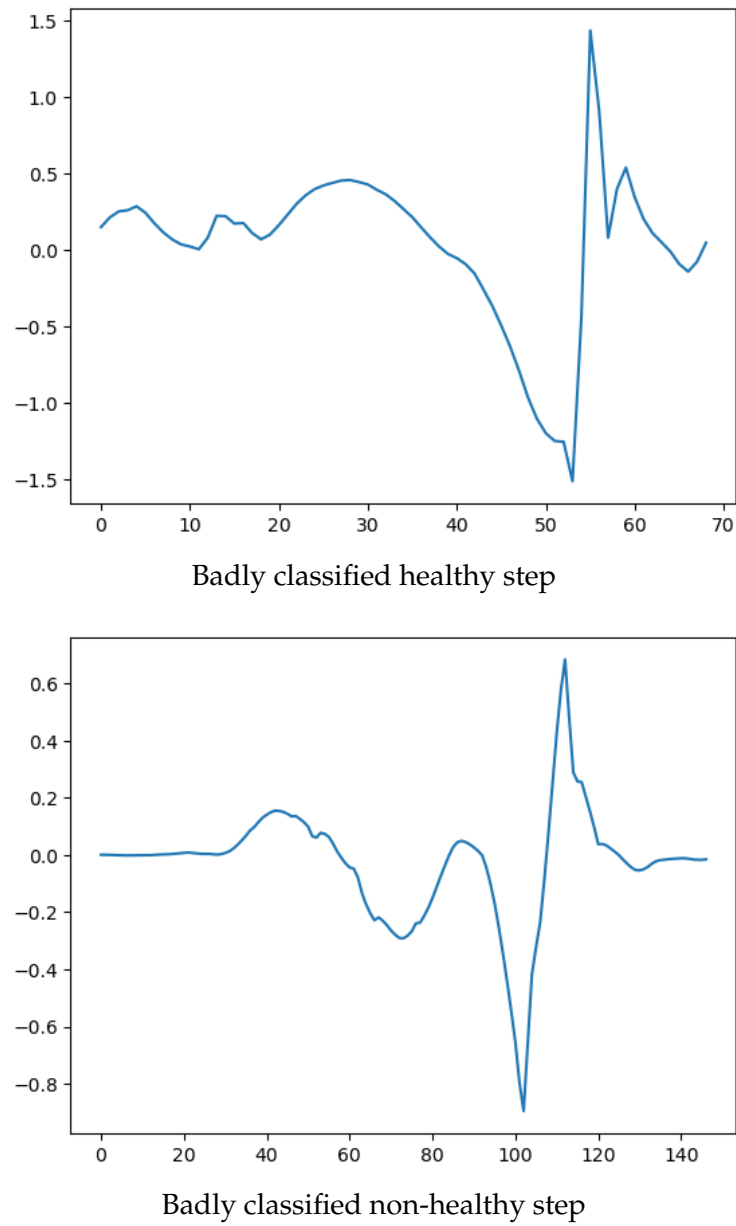


Figure 3: Examples of badly classified steps (see Question 11).