

# An Interpretable Distance Measure for Multivariate Non-Stationary Physiological Signals

Sylvain W. Combettes, Charles Truong, and Laurent Oudre

{sylvain.combettes, charles.truong, laurent.oudre}@ens-paris-saclay.fr  
Université Paris-Saclay, Université Paris Cité, ENS Paris-Saclay, CNRS, SSA, INSERM, Centre Borelli  
Gif-sur-Yvette, France

**Abstract**—We introduce  $d_{symb}$ , a novel distance measure for comparing multivariate non-stationary physiological signals. Unlike most distance measures on multivariate signals such as variants of Dynamic Time Warping (DTW),  $d_{symb}$  can take into account their non-stationarity thanks to a symbolization step. This step is based on a change-point detection procedure, that splits a non-stationary signal into several stationary segments, followed by quantization using  $K$ -means clustering. The proposed distance measure leverages the general edit distance that is applied to the symbolic sequences. The performance of  $d_{symb}$  compared to two commonly used DTW variants is illustrated by applying it to physiological signals recorded during walking protocols. In particular,  $d_{symb}$  is shown to be interpretable: its symbolization detects the segments that correspond to salient behaviors. An open source GitHub repository is made available to reproduce all the experiments in Python.

**Index Terms**—distance measure, multivariate signals, non-stationarity, symbolic representation, change-point detection, interpretability

## I. INTRODUCTION

In the past ten years, sensors have become incredibly prevalent in various fields such as meteorology, finance, healthcare, monitoring, and epidemiology. In healthcare, sensors can be worn by subjects and are capable of measuring different variables. For example, foot-worn Inertial Measurement Units (IMUs) can provide the accelerations and angular velocities in the 3D space [21]. The amount of physiological signals generated by these sensors is massive, and studying them is of paramount importance. For example, the study of human locomotion can lead to early detection and prevention of the risk of fall of pathological subjects. In order to perform inter-individual comparisons or longitudinal follow-up using these physiological signals, it is crucial to define an appropriate distance measure between them. However, when recorded over a long period of time or during complex protocols, the signals are often non-stationary, i.e. their statistical properties change over time. Think, for example, of a connected watch worn for an entire day, during which the subject performs several different activities. A crude comparison of waveforms obtained over two consecutive days (in the time or time-frequency domain) is likely to produce irrelevant results, because, intuitively, the comparison should be made at the level of *actions*, i.e. stationary phases.

Our challenge in this article is to define an interpretable distance measure between multivariate and non-stationary signals. Defining such a distance is difficult because we have

to take into account both interactions between dimensions and abrupt changes in signals caused by the non-stationarity. Several distance measures for multivariate signals exist in the literature (see Section II for an overview). For instance, numerous variants of the popular Dynamic Time Warping (DTW) can be applied to multivariate time series [1]. Those variants can also compare signals of different lengths and align temporal misalignments. However, DTW distances are designed to find alignments directly on the raw waveforms: between samples (if both signals are processed in the time domain) or frames (if they are processed in the time-frequency domain). Therefore, the concept of *stationary phases* is not explicitly taken into account in the distance computation. As will be seen in Section V, this property makes them less suitable and less interpretable in the context of non-stationary signals.

In this article, we propose to handle the non-stationarity by using an adaptive symbolization process. First, we apply change-point detection to the multivariate non-stationary signal to divide the signal into several stationary segments. Next, a clustering procedure assigns a symbol to each stationary segment. These symbols are directly interpretable, as each symbol represents a specific type of behavior within the signal. Once the signal is transformed into a symbolic sequence, we construct a distance measure called  $d_{symb}$  which is inspired from bioinformatics and specifically crafted for these symbolic sequences.

We apply  $d_{symb}$  on time-frequency representations of physiological signals obtained in the context of gait analysis. We show that the symbolic sequences resulting from this data allow for a native and interpretable analysis. Furthermore, compared to the multivariate variants of DTW,  $d_{symb}$  makes more sense regarding the data and metadata, and provides more intuitive results when comparing different exercises or subjects.

The remainder of the paper is organized as follows. Section II provides an overview of distance measures on real-valued multivariate signals and of distance measures on symbolic sequences obtained by symbolic representations. Section III introduces the novel  $d_{symb}$  distance measure, associated with its novel multivariate symbolic representation. Section IV presents the data set of physiological signals (human locomotion) on which  $d_{symb}$  is applied. Section V contains an experimental evaluation of  $d_{symb}$  compared to

variants of DTW. Section VI provides concluding remarks.

## II. BACKGROUND

In this section, we review popular distance measures on real-valued multivariate signals and on symbolic sequences. We then propose a short overview of the symbolization techniques that allows to transform a real-valued multivariate signals into a symbolic sequence.

### A. Distance measures on multivariate signals

We focus here on distances that can cope with signals of different lengths. Such distances are called elastic distance measures. Dynamic Time Warping (DTW) [2], which is arguably the most popular of this group, has been used in numerous data mining tasks [1]. One important feature of DTW is its robustness to time warping, that is, a contraction or dilatation of the time axis. However, it is only defined for univariate signals. Recently, several strategies have been designed to extend DTW to multivariate signals. Two popular approaches are often used in practice: the independent and dependent strategies [3]. In the independent strategy, the univariate DTW is applied to each dimension separately, and the resulting distances on each dimension are summed. The dependent strategy considers the multivariate series as a single series in which each timestamp is associated to a single multidimensional point. The DTW scheme is then applied using Euclidean distances between the multidimensional points of the two series. Other variants of multivariate DTW exist. One such variant is Derivative DTW [4] which applies DTW, not directly on the raw signals, but on their first derivative, in order to prevent unnatural warpings when there is variability in the signals. Another variant, known as Weighted DTW [5], aims at avoiding large warpings by penalizing them using a non-linear multiplicative weight. Both Derivative DTW and Weighted DTW can be combined into a variant called Weighted Derivative DTW.

Distances based on DTW have been used successfully to compare small extracts from multivariate data, but tend to become less suitable for long non-stationary signals. This is because DTW distances are based on a realignment procedure that operates at the level of individual samples and does not take into account the concept of *stationary phase*.

### B. Distance measures on symbolic sequences

A popular distance measure on strings is the edit distance, which is the minimal cost of a sequence of operations that transform a string into another. The edit distance is also known as the Levenshtein distance [18]. The allowed simple operations are the following:

- Insertion of a character in a string.
- Deletion of a character in a string.
- Substitution of characters in both strings.

Each operation is associated with a cost, that also varies upon the transformed characters. The total cost is the sum of the costs of the simple operations. The edit distance can handle symbolic sequences of varying lengths, thanks to the insertions and deletions operations. Many other distances have been

defined for symbolic sequences such as the Longest common subsequence (LCSS) [19], [20]. We refer an interested reader to [17] and [1] for an extensive review.

### C. Symbolization

The basic principle of the distance measure we propose between multivariate signals is to build on existing distance measures for symbolic sequences through the use of symbolization.

Symbolization transforms a real-valued signal  $Q$  of arbitrary length  $n$  into a discrete-valued signals  $\hat{Q}$  of smaller length  $w \leq n$ , called a symbolic sequence. A symbolic representation is often characterized by the number  $A$  of possible values for the symbols, also known as the alphabet size. A common symbolic representation for univariate signals is *Symbolic Aggregate approximation (SAX)* [6], [8] that has successfully being used in several data-mining task such as classification [8], [11], clustering [8] or indexing [14]. Other symbolization techniques for univariate signals include *Symbolic Fourier Approximation (SFA)* [9] or *Adaptive Brownian Bridge-based Aggregation (ABBA)* [10].

Symbolization methods specifically designed for multivariate signals are rarer in the literature. *Trend-based and Valued-based Approximation (TVA)* [12] allows multivariate signals classification by discretizing the means and the trends on uniform segments: U for upwards, D for downwards, or S for straight trend, then combining them on each dimension. SAX-ARM [13] uses SAX to mine association rules efficiently among the deviant events of multivariate time series. However, for both these methods, the symbolization does not natively handle multivariation, but consists in multiple univariate symbolizations that are then handled by a data mining algorithm.

## III. THE $d_{symp}$ DISTANCE MEASURE

This section introduces  $d_{symp}$ , a novel distance measure on multivariate signals of possibly different lengths. This distance measure is designed to handle non-stationarity and to be interpretable. The proposed distance is computed in several steps:

- 1) The multivariate signal is partitioned into stationary segments using a change-point detection procedure,
- 2) Each stationary segment is assigned a symbol through  $K$ -means clustering,
- 3) The final distance  $d_{symp}$  is computed as the general edit distance between the symbolized version of the signals.

Before describing in detail each step, we show in Figure 1 an example of the proposed symbolization on a spectrogram from the Gait data set that will be described in Section IV. Its symbolic representation is displayed below the spectrogram. The different regimes in the spectrogram seem to be detected.

Let  $Q = (q_1, \dots, q_n)$  and  $C = (c_1, \dots, c_m)$  be two real-valued multivariate signals of dimension  $d$ , of lengths  $n$  and  $m$  respectively. We assume that each dimension of  $Q$  and  $C$  have been normalized to zero mean and unit variance.

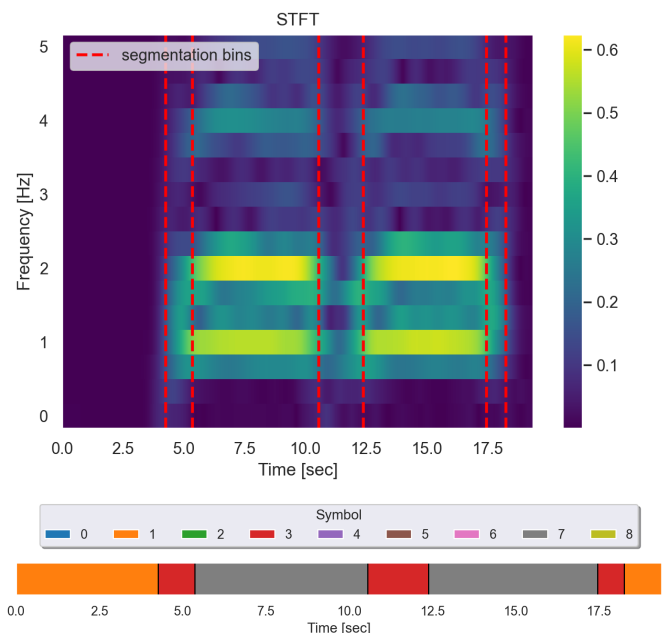


Fig. 1. Example of symbolization of a spectrogram from the Gait data set (see Section IV). The top plot is the spectrogram with the obtained segmentation bins. The bottom plot, called a *color bar*, is a visualization of the resulting symbolic sequence 1373731. The segment limits are indicated with vertical dashed lines.

### A. Step 1: Adaptive Segmentation

Adaptive signal segmentation consists in applying a change-point detection algorithm on the signal at hand, say signal  $Q$  with  $n$  samples. In a nutshell, change-point detection finds the  $w^*$  unknown instants  $t_1^* < t_2^* < \dots < t_{w^*+1}^*$  where some characteristics (here, the mean) of  $Q$  change abruptly. A recent review of such methods is given in [15]. In the context of our symbolization, the number of changes  $w^*$  is unknown and must be estimated too.

The change-point estimates  $\hat{t}_1, \dots, \hat{t}_{\hat{w}+1}$  ( $\hat{w}$  is the number of detected changes) are the minimizers of a discrete optimization problem:

$$\begin{aligned}
 & (\hat{w}, \hat{t}_1, \dots, \hat{t}_{\hat{w}+1}) \\
 & := \arg \min_{(w, t_1, \dots, t_{w+1})} \sum_{k=0}^{w+1} \sum_{t=t_k}^{t_{k+1}-1} \|q_t - \bar{q}_{t_k:t_{k+1}}\|^2 + \lambda(w+1) \quad (1)
 \end{aligned}$$

where  $\bar{q}_{t_k:t_{k+1}}$  is the empirical mean of  $\{q_{t_k}, \dots, q_{t_{k+1}-1}\}$  and  $\lambda > 0$  is a penalization parameter. (By convention,  $t_0 := 0$  and  $t_{w+1} := n$ .) The penalized formulation (1) seeks a compromise between the reconstruction error given by the sum of quadratic errors and the complexity given by the number of change-points. Problem (1) is solved using the Pruned Exact Linear Time (PELT) algorithm [16], which is shown to have  $\mathcal{O}(n)$  complexity under the assumption that the segment lengths are randomly drawn from a uniform distribution.

Intuitively, the  $\lambda$  parameter penalizes the introduction of a new change-point: when  $\lambda$  is small, many change-points are detected. Once the user chooses a penalty  $\lambda$ , the segmentation

procedure returns the segment bins and the estimated number of segments. For calibration purposes, we use the standard scaling  $\lambda = \ln(n)$  [15].

Examples of segmentation of two spectrograms are given Figure 1. The proposed segmentation returns symbolic sequences of different lengths and splits the non-stationary multivariate signal into several segments that look stationary.

### B. Step 2: Symbolization

Once the segment boundaries have been determined for all multivariate signals in our data set, the mean per segment (in dimension  $d$ ) is computed for each multivariate signal. The means per segment are centered and scaled to unit variance. Thanks to Step 1, the segments correspond to mean-shifts, thus it is reasonable to represent each segment by its mean value. Then, these means per segment, from all segments of all multivariate signals in our data set, are clustered using the  $K$ -means algorithm where the number of clusters is the desired number of symbols  $A$ . Finally, each segment is attributed a symbol: the label of its assigned cluster. Due to the symbolization through clustering, there is no reason for the obtained symbols to be equiprobable. Having non-equiprobable symbols can be useful in several tasks such as anomaly detection or outlier removal. Indeed, the segments that are attributed to rarely obtained symbols can be considered as containing anomalies.

### C. Step 3: Compute the Distance Measure $d_{symp}$

The proposed  $d_{symp}$  distance measure leverages the general edit distance described in Section II-B. The operation costs of the edit distance are defined so that they take into account the dissimilarity between individual symbols:

- The substitution cost  $\text{sub}(a, b)$  for individual symbols  $a$  and  $b$  is the Euclidean distance between the cluster center  $G_a$  of symbol  $a$  and the cluster center  $G_b$  of symbol  $b$ 

$$\text{sub}(a, b) = \|G_a - G_b\|_2. \quad (2)$$
- For all characters, the insertion and deletion costs are fixed to  $\text{sub}_{\max}$ , where  $\text{sub}_{\max}$  is the maximum value of the modified substitute costs in Formula (2).

Given the costs,  $d_{symp}$  should do more substitutions than insertions or deletions.

The input of  $d_{symp}$  is a replicated version of the symbolic sequences. Let  $\tilde{Q}$  and  $\tilde{C}$  be the symbolic representations of  $Q$  and  $C$  respectively, of word lengths  $w_Q$  and  $w_C$  respectively. The segments obtained from the adaptive segmentation step are of varying lengths. We incorporate the segment length information into the symbolic sequences by replicating each symbol proportionally to its segment length. Let  $l_{Q,1}, \dots, l_{Q,w_Q}$  and  $l_{C,1}, \dots, l_{C,w_C}$  be the segment lengths resulting from the segmentation of signals  $Q$  and  $C$  respectively. Each segment length is divided by the minimum of all segments lengths of all involved symbolic sequences (here  $\tilde{Q}$  and  $\tilde{C}$ ) to obtain the normalized segment lengths  $\hat{l}_{Q,1}, \dots, \hat{l}_{Q,w_Q}$  and  $\hat{l}_{C,1}, \dots, \hat{l}_{C,w_C}$ . Then, the symbolic sequences are replicated by the normalized lengths. In

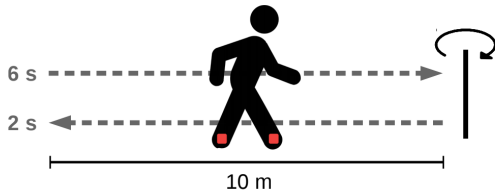


Fig. 2. Schematic protocol recorded by the sensors that are located with red squares. Source: [21].

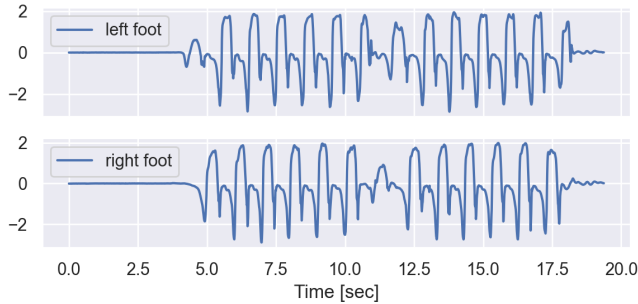


Fig. 3. Scaled univariate gait signals of the left and right foot activities corresponding to the same recording from the human locomotion data set [21], with 1,938 samples each.

the example of the symbolic representation of  $Q$ , the symbol of the first segment  $\ell_{Q,1}$  times, then the symbol of the second segment  $\hat{\ell}_{Q,2}$  times, etc. Finally,  $d_{symb}(Q, C)$  is equal to the general edit distance between these replicated symbolic sequences.

#### IV. PHYSIOLOGICAL SIGNALS: HUMAN LOCOMOTION

In this article, we apply the  $d_{symb}$  distance measure to an open-access data set of real-world physiological signals recorded in the context of human locomotion analysis [21].

##### A. Data and transformations

In this article, walking data consists of angular velocity recorded on the left and right feet using a pair of MTw XSens sensors (sampling frequency: 100 Hz). The gait protocol for all subjects is depicted in Figure 2: standing still for 6 sec, walking 10 meters at the speed they felt comfortable with, turning around, walking back to the initial position, and standing for 2 sec. Note that since the subjects perform the protocol at different speeds, each recording is of different duration. Examples of raw signals, recorded for a healthy subject, are displayed on Figure 3.

Since locomotion is an activity that has a strong periodic component, it is common in the literature to process such signals in the time-frequency domain. For each univariate gait signal, we compute its *Short Time Fourier Transform (STFT)*, with a window length equal to 300 samples (3 seconds) and overlap length of 299 samples (providing one frame each 0.01 seconds). Only the 0–5 Hz frequency band, where phenomena of interest are contained, is kept. The norms of the STFT coefficients are computed and concatenated, providing  $d = 16$

frequency bins per frame (14 per signal). The output data will be seen as a  $d$ -dimensional multivariate signal. The spectrogram associated to the left foot of Figure 3 is displayed on the top of Figure 1. Comparing Figure 1 to Figure 3, we can observe that the static phase at the beginning of the recording and the change of periodicity around the middle appear clearly on the spectrogram. The spectrogram can therefore be seen as a non-stationary multivariate signal, on which we can apply the  $d_{symb}$  distance measure.

##### B. Subjects

The data set is composed of 221 recordings: 192 from healthy subjects, 21 from patients with neurological pathologies (such as cerebellar disorders) and 8 from patients with orthopedic pathologies (such as knee injuries). Note that each recording is associated with two signals: one for each foot.

#### V. EXPERIMENTAL RESULTS

In this section,  $d_{symb}$  is applied to the open-access human locomotion data set described in Section IV, which is composed of 442 spectrograms. The number of symbols  $A$  is set to 9, as in previous publications [10]. A Python implementation of  $d_{symb}$ , along with codes to reproduce the figures and scores in this paper, can be found in a GitHub repository<sup>1</sup>.

##### A. Interpretation of the symbolization

Symbolizations for an extract of 60 of the 442 spectrograms are shown in Figure 4. In this visual, each symbol is associated with a different color, allowing symbolisations to be displayed as color bars.

Three comments can be done from visual inspection. First, the general structure of the symbolic sequences are coherent with the protocol defined in Section IV. We observe a alternation of several stationary segments: the standing still segment (which is always associated to symbol 1), one segment corresponding to the initiation step, one walking segment, one U-turn segment, one walking segment, one segment corresponding to the termination step, and one final standing still segment. The change-point detection procedure allows to precisely detect the boundaries of these segments and to capture the non-stationary structure of the signals.

Second, we notice that each symbol is associated with a specific type of behavior. This can easily be highlighted by plotting the centroid corresponding to each symbol, which is a vector of dimension 16 that captures the average spectrogram frame for all its segments. These centroids can be interpreted as Power Spectral Densities (PSDs) and are displayed on Figure 5. Three types of behaviors can be found:

- Flat behavior (symbol 1) which is likely to correspond to the static phase.
- Harmonic behavior (symbols 0, 2, 4, and 7) where the spectral content is mostly carried by one fundamental frequency and its first harmonic, and which corresponds to regular walking with a periodic structure.

<sup>1</sup><https://github.com/sylvaincom/d-symb>

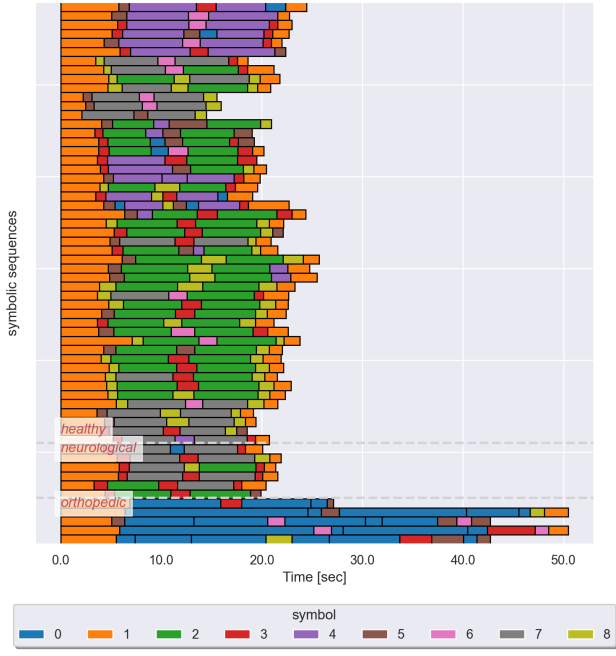


Fig. 4. Color bars for several signals of different pathology groups that are separated by white dashed horizontal lines. Each row is the color bar corresponding to a symbolic sequence. Only a subset of 60 out of the 442 signals is shared for conciseness of visualization.

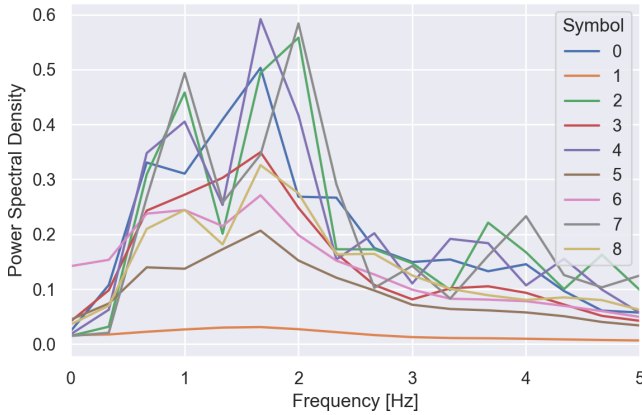


Fig. 5. Power spectral density for each symbol.

- Low-pass behavior (symbols 3, 5, 6, and 8) that is linked to initiation/termination steps and U-turn.

These observations are also confirmed when running a hierarchical clustering algorithm on the distance matrix between centroids and visualising the associated dendrogram (see Figure 6).

Third, one can observe that each symbol is not only characteristic of a human locomotion regime, but it can also be the "signature" of a pathology group or a laterality. Indeed, as displayed in Figure 7, some symbols are specific to the left or right foot, or to a particular pathology group. For example, symbol 6 mostly corresponds to the right foot and symbol 7

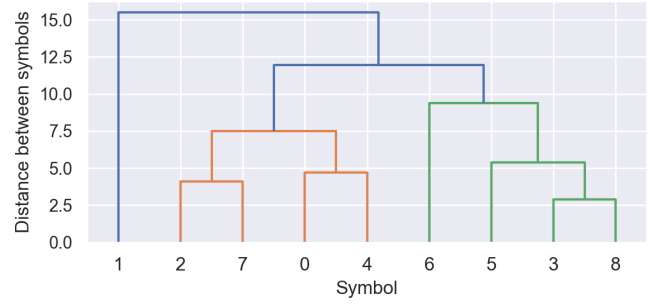


Fig. 6. Dendrogram: distance between the individual symbols and how they are grouped according to hierarchical clustering.

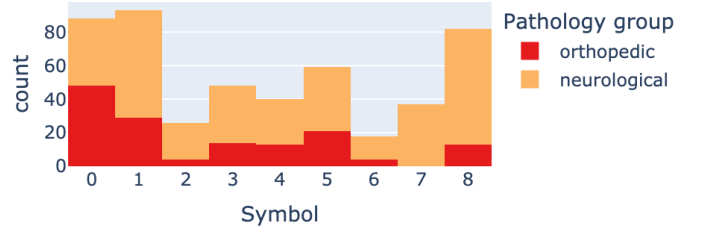
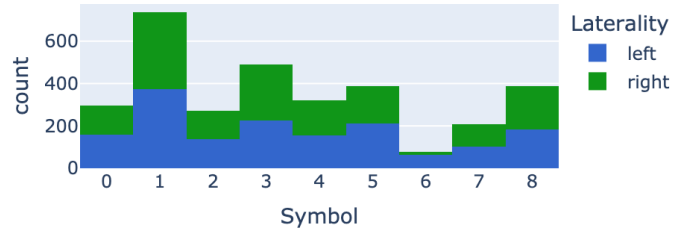


Fig. 7. Histograms of the symbols obtained throughout all 442 symbolic sequences, with an emphasis on the laterality (top) and on the pathology group (bottom). For the bottom histogram, the healthy subjects are not displayed because they are the majority group and would alter the visualization.

has only been assigned to the neurological group.

In conclusion, it appears that the proposed symbolization enables us to find the structure of the signals, and to characterize each segment according to its type (walk, U-turn, etc.), and its links with a specific group. As  $d_{symb}$  is based on these symbolic sequences, this suggests that all these phenomena will be taken into account when comparing multivariate signals.

### B. Interpretation of the $d_{symb}$ distance measure

Based on the previously described symbolizations, the distance matrix (according to the proposed distance measure  $d_{symb}$ ) between all 442 multivariate signals are computed. For sake of comparison, the distance matrix for the dependent DTW (DTW-D) and the independent DTW (DTW-I) (see Section II) are computed. Note that the DTW distances are computed directly on the spectrograms (and not on the symbolic sequences).

In a first experiment, for each distance measure, the silhouette coefficient is calculated using the distance matrix and the ground truth labels corresponding to the pathological group

TABLE I

SUMMARY OF THE SILHOUETTE SCORES FOR EACH DISTANCE MEASURE, AVERAGED OVER ALL SIGNALS.

Distance measure	Mean Silhouette score	Median Silhouette score
DTW-D	0.15	0.18
DTW-I	0.15	0.19
$d_{symb}$	0.33	0.40

(healthy, neurological, or orthopedic). The score is bounded between  $-1$  for incorrect grouping and  $+1$  for highly dense and well-separated groups. The obtained Silhouette scores are given in Table I. According to these results,  $d_{symb}$  provides groups that are denser and better separated than DTW-D and DTW-I. It seems that  $d_{symb}$  better captures the properties of each group and is able to detect subtle differences between subjects. This can probably be explained by the symbolization process and the relevance of the different symbols already highlighted in Section V-A.

For the second experiment, let us consider two scaled univariate gait signals of different lengths. These signals are displayed on Figure 8 along with the symbolization of their spectrograms. The number of samples of the top signal is 2,550 and for the bottom, it is 1,994. Their difference in length is 556, while the mean difference of lengths out of all 442 signals is 468, hence these signals can be considered to be of different lengths. The  $d_{symb}$  distance between these two signals is 70, while the mean of  $d_{symb}$  out of all 442 signals is 155. Thus, these signals are considered similar by  $d_{symb}$ . Therefore, despite being of dissimilar lengths, signals can be considered as similar by  $d_{symb}$ . This suggests that  $d_{symb}$  is robust to the difference in lengths and focuses on the phenomena of interest in the signals.

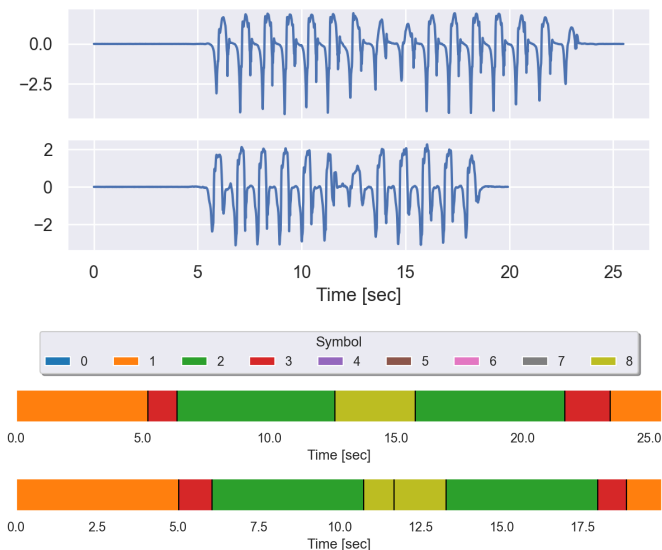


Fig. 8. Two scaled univariate gait signals of different lengths, along with the symbolization of their spectrograms.

For the third experiment, the distance distribution between the right and left feet for a given recording is studied. As

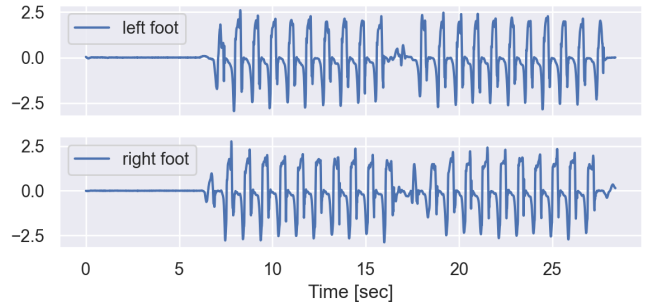


Fig. 9. Left and right foot activities of a subject with a neurological pathology.

TABLE II

MEAN RANK OF QUERYING A SIGNAL'S OPPOSITE LATERALITY.

Distance measure	Pathology group		
	Healthy	Neurological	Orthopedic
DTW-D	15.8	19.9	6.4
DTW-I	6.7	12.6	3.9
$d_{symb}$	9.3	5.1	2.0

can be seen on Figure 3, the structure of both signals are rather similar for these healthy subjects, except for the U-turn segment. Indeed, during a healthy U-turn, one foot turns around, while the other serves as a support foot, thus creating an asymmetry in the recorded signals. On the contrary, for pathological subjects whose gait is affected, right and left foot movements tend to be more symmetrical, even in the U-turn segment, as can be seen in Figure 9. Then, although it may appear counter-intuitive, a *good* distance between the right and left feet is expected to be large for healthy subjects and small for pathological ones. To investigate this question, a query-by-content task is conducted. Given a recording of a right foot, we compute the rank at which the corresponding left foot is found according to  $d_{symb}$ , DTW-D, and DTW-I. The average ranks for each distance measure and each group are displayed on Table II. We observe that the results obtained with  $d_{symb}$  are coherent with the medical considerations: the average ranks are large for healthy subjects and small for pathological subjects. On the contrary, both DTW-D and DTW-I fail to capture this property: this can easily be shown by considering the ranks obtained on the neurological group.

## VI. CONCLUSION

We have introduced  $d_{symb}$ , a novel distance measure on multivariate and non-stationary signals.  $d_{symb}$  leverages the general edit distance. It uses a novel symbolization scheme for multivariate signals to transform the real-valued multivariate signals into symbolic sequences, that are then fed to the edit distance. Thanks to the adaptive segmentation of the proposed symbolization,  $d_{symb}$  can handle the non-stationarity, which is remarkable.

$d_{symb}$  has been applied to spectrograms of gait signals, which are multivariate and non-stationary. Experiments have shown how interpretable the symbolization is. Indeed, each



symbol corresponds to a specific regime of human locomotion. Moreover,  $d_{symb}$  is more suitable than DTW-D and DTW-I in order to compare multivariate non-stationary signals. For instance, it creates groups that are more dense and better separated. Please note that this article primarily addresses the use case of human locomotion to exemplify the various properties of  $d_{symb}$ . However,  $d_{symb}$  is not limited to Gait analysis and can be effectively utilized for analyzing any other type of data as well.

#### ACKNOWLEDGMENT

Sylvain W. Combettes is supported by the Industrial Data Analytics and Machine Learning chair of ENS Paris-Saclay, and by a public grant overseen by the French National Research Agency (ANR) through the program UDOPIA, project funded by the ANR-20-THIA-0013-01. Charles Truong is funded by the PhLAMES chair of ENS Paris-Saclay. Part of the computations has been executed on Atos Edge computer, funded by the Industrial Data Analytics and Machine Learning chair of ENS Paris-Saclay.

#### REFERENCES

- [1] A. Shifaz, C. Pelletier, F. Petitjean, and G. I. Webb, "Elastic similarity and distance measures for multivariate time series," *Knowl Inf Syst.*, 2023.
- [2] D. J. Berndt and J. Clifford, "Using dynamic time warping to find patterns in time series," In *KDD workshop* (Vol. 10, No. 16, pp. 359-370), July 1994.
- [3] M. Shokooi-Yekta, B. Hu, H. Jin, J. Wang, and E. Keogh, "Generalizing DTW to the multi-dimensional case requires an adaptive approach," *Data Min Knowl Disc* 31(1):1–31, 2017.
- [4] E. J. Keogh, M. J. Pazzani, "Derivative dynamic time warping," In: *Proceedings of the 2001 SIAM international conference on data mining*, SIAM, pp 1–11, 2001.
- [5] Y.-S. Jeong, M. Jeong, O. Omiaomu, "Weighted dynamic time warping for time series classification," *Pattern Recognit* 44(9):2231–2240, 2011.
- [6] J. Lin, E. Keogh, S. Lonardi, and B. Chiu, "A symbolic representation of time series, with implications for streaming algorithms," In *Proceedings of the 8th ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery, DMKD '03*, page 2–11, New York, NY, USA, 2003.
- [7] T. Le Nguyen, S. Gsponer, I. Ilie, M. O'Reilly, and G. Ifrim, "Interpretable time series classification using linear models and multi-resolution multi-domain symbolic representations," *Data Min Knowl Disc* 33: 1183-1222, 2019.
- [8] J. Lin, E. Keogh, L. Wei, and S. Lonardi, "Experiencing sax: a novel symbolic representation of time series," *Data Min Knowl Disc*, 15:107–144, 2007.
- [9] P. Schäfer and M. Höggqvist, "Sfa: A symbolic fourier approximation and index for similarity search in high dimensional data sets," In *Proceedings of the 15th International Conference on Extending Database Technology, EDBT '12*, page 516–527, New York, NY, USA, 2012.
- [10] S. Elsworth and S. Güttel, "Abba: adaptive brownian bridge-based symbolic aggregation of time series," *Data Min Knowl Disc*, 34:1175–1200, 2020.
- [11] T. L. Nguyen, S. Gsponer, I. Ilie, M. O'Reilly, and G. Ifrim, "Interpretable time series classification using linear models and multi-resolution multi-domain symbolic representations," *Data Min Knowl Disc*, 33:1183–1222, 2019.
- [12] B. Esmael, A. Arnaout, R. K. Fruhwirth, and G. Thonhauser. "Multivariate time series classification by combining trend-based and value-based approximations." In *Computational Science and Its Applications – ICCSA 2012*, pages 392–403, Berlin, Heidelberg, 2012.
- [13] H. Park and J.-Y. Jung, "Sax-arm: Deviant event pattern discovery from multivariate time series using symbolic aggregate approximation and association rule mining," *Expert Systems with Applications*, 141, p. 112950, 2020.
- [14] A. Camera, T. Palpanas, J. Shieh and E. Keogh, "isax 2.0: indexing and mining one billion time series," *2010 IEEE International Conference on Data Mining*, Sydney, NSW, Australia, pp. 58-67, 2010.
- [15] C. Truong, L. Oudre, and N. Vayatis, "Selective review of offline change point detection methods," *Signal Processing*, 167:107299, 2020.
- [16] R. Killick, P. Fearnhead, and I. A. Eckley, "Optimal detection of changepoints with a linear computational cost," *Journal of the American Statistical Association*, 2012.
- [17] G. Navarro, "A guided tour to approximate string matching," *ACM Comput. Surv.*, 33(1):31–88, 2001.
- [18] V. I. Levenshtein et al., "Binary codes capable of correcting deletions, insertions, and reversals," In *Soviet Physics Doklady*, volume 10, pages 707–710, 1966.
- [19] D. S. Hirschberg, "Algorithms for the longest common subsequence problem." *J ACM* 24(4):664–675 21, 1977.
- [20] M. Vlachos, G. Kollios, and D. Gunopulos, "Discovering similar multi-dimensional trajectories," In: *Proceedings 18th international conference on data engineering*, IEEE, pp 673–684, 2002.
- [21] C. Truong, R. Barrois-Müller, T. Moreau, C. Provost, A. Vienne-Jumeau, A. Moreau, P.-P. Vidal, N. Vayatis, S. Buffat, A. Yelnik, D. Ricard, and L. Oudre, "A data set for the study of human locomotion with inertial measurements units," *Image Processing On Line*, 9, pp. 381–390, 2019.