



HAL
open science

Constructing Positive Definite Elastic Kernels with Application to Time Series Classification

Pierre-François Marteau, Sylvie Gibet

► **To cite this version:**

Pierre-François Marteau, Sylvie Gibet. Constructing Positive Definite Elastic Kernels with Application to Time Series Classification. 2010. hal-00486916v2

HAL Id: hal-00486916

<https://hal.science/hal-00486916v2>

Preprint submitted on 9 Jul 2010 (v2), last revised 25 May 2014 (v12)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Constructing Positive Definite Elastic Kernels with Application to Time Series Classification

Pierre-François Marteau, *Member, IEEE* and Sylvie Gibet, *Member, IEEE*

Abstract—This paper proposes some extensions to the work on kernels dedicated to string alignment (biological sequence alignment) based on the summing up of scores obtained by local alignments with gaps. The extensions we propose allow to construct, from classical time-warp distances, what we called summative time-warp kernels that are positive definite if some simple sufficient conditions are satisfied. Furthermore, from the same formalism, we derive a time-warp inner product that extends the usual euclidean inner product, providing the capability to handle discrete sequences or time series of variable lengths in an Hilbert space. The classification experiment we conducted, using either first near neighbor classifier or Support Vector Machine classifier leads to conclude that the positive definite elastic kernels we propose outperform the distance substituting kernels for the classical elastic distances we tested. In a similar way, the kernel based on the distance induced by the time-warp inner product outperforms significantly on the considered task the kernel based on the euclidean distance.

Index Terms—Elastic distance, Time warp kernel, Time warp inner product, Definiteness, Time series classification, SVM.



1 INTRODUCTION

ELASTIC similarity measures such as Dynamic Time Warping (DTW) or Edit Distances have proved to be quite efficient compared to non elastic similarity measures such as euclidean measures or LP norms when addressing tasks that require the matching of time series data, in particular time series clustering and classification. A wide scope of applications as in physics, chemistry, finance, bio-informatics, network monitoring, etc, have demonstrated the benefits of using elastic measures. A natural follow-up to the elaboration of elastic measures is to examine whether or not it is possible to construct Reproducing Time Warp Hilbert Spaces (RTWHS) from a given elastic measure, basically vector spaces characterized with inner products having time-warp capabilities. Another intriguing question is to determine whether it is possible or not to define an inner product structure from which a given elastic measure is induced? This question, apart from its theoretical implication, has a great impact when considering the potential application fields, since, if the answer is positive, it provides direct access to the Linear Algebra results and tools.

Unfortunately it seems that common elastic measures that are derived from DTW or Edit Distance are not directly induced by an inner product of any sort, even when such measures are metrics. One can conjecture that it is not possible to embed time series in an Hilbert space having a time-warp capability using these classical elastic measures.

This paper aims at exploring this issue and suggests of Time Warp Kernels (TWK) constructions that try to preserve the properties of elastic measures from which they are derived, while offering the possibility possibility of embedding time series in TWHS. The main contributions of the paper are as follows

- 1) we establish the indefiniteness of the main time-warp measures used in the literature,
- 2) we propose some methods to construct positive definite kernels from classical time-warp measures,
- 3) we define simple Time Warp Inner Product (TWIP) as an extension to the Euclidean Inner Product, and
- 4) we experiment and compare the proposed kernels on time series classification tasks using a large variety of time series datasets.

The paper is organized as follows: the second section of the paper synthesizes the related works; the third section introduces the notation and mathematical backgrounds that are used throughout the paper; the fourth section addresses the non definiteness of classical elastic measures that prevents the direct construction of an inner product from these measures. The fifth section develops the construction of some TWK and TWIP from classical elastic measures and discusses their potential benefits. The sixth section gathers clustering and classification experimentations on a wide range of time series data and compares TWK and TWIP accuracies with classical elastic and non elastic measures. The seventh section proposes a conclusion and some further research perspectives. An appendix gives the proof of our main results.

2 RELATED WORKS

During the last decades, the use of kernel-based methods in pattern analysis has provided numerous results and

• P.F. Marteau and Sylvie Gibet are with the VALORIA Lab., Université Européenne de Bretagne, Université de Bretagne Sud, 56000 Vannes, France.
E-mail: {Pierre-Francois.Marteau, Sylvie.Gibet}@AT.univ-ubs.fr

fruitful applications in various domains such as biology, statistics, networking, signal processing, etc. Some of these domains, such as bioinformatics, or more generally domains that rely on sequence or time series models, require the analysis and processing of variable length vectors, sequences or timestamped data. Various methods and algorithms have been developed to quantify the similarity of such objects. From the original dynamic programming [2] implementation of the symbolic edit distance [12] by Wagner and Fisher [26], the Smith and Waterman (SW) algorithm [22] has been designed to evaluate the similarity between two symbolic sequences by means of a local gap alignment. More efficient local heuristics have since been proposed to meet the massive symbolic data challenge, such as BLAST [1] or FASTA [16]. Similarly, dynamic time warping measures have been developed to evaluate similarity between numeric time series or timestamped data [25], [19], and more recently [6], [14] propose elastic metrics dedicated to such numeric data.

Our capability to construct kernels with elastic or time-warp properties from such an ‘elastic distance’ has attracted attention since significant benefits are expected from potential applications of kernel-based machine learning algorithms to variable length data, or more generally data for which some elastic matching has a meaning. Among the kernel machine algorithms applicable to discrimination or regression tasks, Support Vector Machines (SVM) are reported to yield state-of-the-art performances.

SVM or vast margin classifiers [23], [4], [21] are a set of supervised algorithms that learn how to solve discrimination or regression problems from positive and negative examples. They generalize linear classification algorithms by integrating two concepts: the maximal margin principle and a kernel function that defines the similarity or dissimilarity of any pair of examples, typically such as an inner product between the vector representation of two examples.

The definition of ‘good’ kernels from known elastic or time-warp distances applicable to data objects of variable lengths has been a major challenge since the 1990s. The notion of ‘goodness’ has rapidly been associated to the concept of definiteness. Basically SVM algorithms involve an optimization process whose solution is proved to be uniquely defined if and only if the kernel is positive definite: in that case the objective function to optimize is quadratic and the optimization problem convex. Nevertheless, if the definiteness of kernels is an issue, in practice, many situations exist where definite kernels are not applicable. This seems to be the case for the main elastic measures traditionally used to estimate the similarity of objects of variable lengths. A pragmatic approach consists of using indefinite kernels, although contradictory results have been reported about the impact of definiteness or indefiniteness of kernels on the empirical performances of SVMs. The sub-optimality of the non-convex optimization process is possibly one of

the causes leading to these un-guaranteed performances [27], [8]. Regulation procedures have been proposed to locally approximate indefinite kernel functions by definite ones with some benefits. Among others, some approaches apply direct spectral transformations to indefinite kernels: the methods [28] consist in flipping the negative eigenvalues or shifting the eigenvalues using the minimal shift value required to make the spectrum of eigenvalues positive, and reconstructing the kernel with the original eigenvectors in order to produce a positive semidefinite kernel. Yet, in general, ‘convexification’ procedures are difficult to interpret geometrically and the expected effect of the original indefinite kernel may be lost. Some theoretical highlights have been provided through approaches that consist in embedding the data into a pseudo-Euclidean (pE) space and in formulating the classification problem with an indefinite kernel, such as that of minimizing the distance between convex hulls formed from the two categories of data embedded in the pE space [9]. The geometric interpretation results in a constructive method allowing for the understanding, and in some cases the prediction of the classification behavior of an indefinite kernel SVM in the corresponding pE space.

Our approach is founded on the work of Haussler on convolution kernels [10] defined on a set of discrete structures such as strings, trees, or graphs. The iterative method that is developed is generative, as it allows for the building of complex kernels from the convolution of simple local kernels. Following the work of Haussler [10], Saigo et al [18] define, from the Smith and Waterman algorithm [22], a kernel to detect local alignment between strings by convolving simpler kernels. These authors show that the Smith and Waterman distance measure dedicated to determining similar regions between two nucleotide or protein sequences, that is not definite, is nevertheless connected to the logarithm of a point-wise limit of a series of definite convolution kernels. In fact, these previous studies have very general implications, the first being that classical elastic measures can also be understood as the limit of a series of definite convolution kernels. We generalize to some extent the results presented by Saigo et al. on the Smith and Waterman algorithm and propose extensions to construct time-warp inner products.

3 NOTATIONS AND MATHEMATICAL BACK-GROUNDS

To ensure that this paper is self-contained, we give commonly used definitions, with few details, for metric or quasi metric, inner product, kernel and definiteness, sequence set, and classical elastic measures.

3.1 Premetric, pseudometric and metric

Definition 3.1: A premetric on a set U is a function $\delta : U \times U \rightarrow \mathbb{R}$ which satisfies the following axioms: For all $(x, y) \in U \times U$,

- 1) $\delta(x, y) \geq 0$ (non negativity)
- 2) $\delta(x, x) = 0$ (null if identical)

Definition 3.2: A pseudometric on a set U is a function $\delta : U \times U \rightarrow \mathbb{R}$ which satisfies the following axioms:

For all $(x, y) \in U \times U$,

- 1) $\delta(x, y) \geq 0$ (non negativity)
- 2) $\delta(x, x) = 0$ (null if identical)
- 3) $\delta(x, y) = \delta(y, x)$ (symmetry)
- 4) $\delta(x, z) \leq d(x, y) + d(y, z)$. (subadditivity/triangle inequality)

Definition 3.3: A metric, also called a distance, on a set U is a pseudometric $\delta : U \times U \rightarrow \mathbb{R}$ for which the second axiom rewrites :

For all $(x, y) \in U \times U$, $\delta(x, y) = 0$ if and only if $x=y$. (null iff identical)

3.2 Inner Product

In the following, the field of scalars denoted \mathbb{F} is either the field of real numbers \mathbb{R} or the field of complex numbers \mathbb{C} .

Definition 3.4: An inner product space is a vector space V over the field \mathbb{F} together with an inner product, i.e., with a map

$$\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{F}$$

that satisfies the following three axioms for all vectors $x, y, z \in V$ and all scalars $a \in \mathbb{F}$:

- 1) Conjugate symmetry:
 $\langle x, y \rangle = \overline{\langle y, x \rangle}$.
- 2) Linearity in the first argument:
 $\langle ax, y \rangle = a\langle x, y \rangle$.
 $\langle x + y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$.
- 3) Positive-definiteness:
 $\langle x, x \rangle \geq 0$ with equality only for $x = 0$.

3.3 Kernel and definiteness

Definition 3.5: A kernel on a non empty set U refers to a complex (or real) valued symmetric function $\varphi(x, y) : U \times U \rightarrow \mathbb{C}$ (or \mathbb{R}).

Definition 3.6: Let U be a non empty set. A function $\varphi : U \times U \rightarrow \mathbb{C}$ is called a positive (resp. negative) definite kernel if and only if it is Hermitian (i.e. $\varphi(x, y) = \overline{\varphi(y, x)}$ where the overline stands for the conjugate number) for all x and y in U and $\sum_{i,j=1}^n c_i \bar{c}_j \varphi(x_i, x_j) \geq 0$ (resp. $\sum_{i,j=1}^n c_i \bar{c}_j \varphi(x_i, x_j) \leq 0$), for all n in \mathbb{N} , $\{x_1, x_2, \dots, x_n\} \subseteq U$ and $\{c_1, c_2, \dots, c_n\} \subseteq \mathbb{C}$.

Definition 3.7: Let U be a non empty set. A function $\varphi : U \times U \rightarrow \mathbb{C}$ is called a conditionally positive (resp. conditionally negative) definite kernel if and only if it is Hermitian (i.e. $\varphi(x, y) = \overline{\varphi(y, x)}$ for all x and y in U) and $\sum_{i,j=1}^n c_i \bar{c}_j \varphi(x_i, x_j) \geq 0$ (resp. $\sum_{i,j=1}^n c_i \bar{c}_j \varphi(x_i, x_j) \leq 0$), for all $n \geq 2$ in \mathbb{N} , $\{x_1, x_2, \dots, x_n\} \in U^n$ and $\{c_1, c_2, \dots, c_n\} \in \mathbb{C}^n$ with

$$\sum_{i=1}^n c_i = 0.$$

In the last two above definitions, it is easy to show that it is sufficient to consider mutually different elements in U , i.e. collections of distinct elements x_1, x_2, \dots, x_n . This is what we consider for the remaining of the paper

Definition 3.8: A positive (resp. negative) definite kernel defined on a finite set U is also called a positive (resp. negative) semidefinite matrix. Similarly, a positive (resp. negative) conditionally definite kernel defined on a finite set is also called a positive (resp. negative) conditionally semidefinite matrix.

For convenience sake, we will use PD, ND, CPD and CND for positive definite, negative definite, conditionally positive definite and conditionally negative definite to characterize either kernel or matrix throughout the paper.

Constructing PD kernels from CND kernels is quite straightforward. For instance, if φ is a CND kernel on a set U and $\Omega \in U$ then **[**ref**]** $\varphi'(x, y) = \varphi(x, \Omega) + \varphi(y, \Omega) - \varphi(x, y) - \varphi(\Omega, \Omega)$ is a PD kernel, so are $e^{(\varphi'(x, y))}$ and $e^{-\varphi(x, y)}$. The converse is also true.

3.4 Sequence set

Definition 3.9: Let \mathbb{U} be the set of finite sequences (symbolic sequences or time series): $\mathbb{U} = \{A_1^p \mid p \in \mathbb{N}\}$. A_1^p is a sequence with discrete index varying between 1 and p . We note Ω the empty sequence (with null length) and by convention $A_1^0 = \Omega$ so that Ω is a member of set \mathbb{U} . $|A|$ denotes the length of the sequence A . Let $\mathbb{U}_p = \{A \in \mathbb{U} \mid |A| \leq p\}$ be the set of sequences whose length is shorter or equal to p .

Definition 3.10: Let A be a finite sequence. Let a_i^j be the i^{th} element (symbol or sample) of sequence A . We will consider that $a_i^j \in S \times T$ where S embeds the multidimensional space variables (either symbolic or numeric) and $T \subset \mathbb{R}$ embeds the time stamp variable, so that we can write $a_i^j = (a_i, t_{a_i})$ where $a_i \in S$ and $t_{a_i} \in T$, with the condition that $t_{a_i} > t_{a_j}$ whenever $i > j$ (time stamps strictly increase in the sequence of samples). A_i^j with $i \leq j$ is the subsequence consisting of the i^{th} through the j^{th} element (inclusive) of A . So $A_i^j = a_i^j a_{i+1}^j \dots a_j^j$. Λ denotes the null element. A_i^j with $i > j$ is the null time series, e.g. Ω .

3.5 General Edit/Elastic distance on a sequence set

Definition 3.11: An edit operation is a pair $(a', b') \neq (\Lambda, \Lambda)$ of sequence elements, written $a' \rightarrow b'$. Sequence B results from the application of the edit operation $a \rightarrow b$ into sequence A , written $A \Rightarrow B$ via $a' \rightarrow b'$, if $A = \sigma a' \tau$ and $B = \sigma b' \tau$ for some subsequences

σ and τ . We call $a' \rightarrow b'$ a match operation if $a' \neq \Lambda$ and $b' \neq \Lambda$, a delete operation if $b' = \Lambda$, an insert operation if $a' = \Lambda$.

For any pair of sequences A_1^p, B_1^q , for which we consider the extensions A_0^p, B_0^q whose first element is the null symbol Λ , and for each elementary edit operation related to position $0 \leq i \leq p$ in sequence A and to position $0 \leq j \leq q$ in sequence B is associated to a cost value $\Gamma_{a'_i \rightarrow b'_j}(A_1^p, B_1^q)$, or $\Gamma_{a'_i \rightarrow \Lambda, j}(A_1^p, B_1^q)$ or $\Gamma_{\Lambda, i \rightarrow b'_j}(A_1^p, B_1^q) \in \mathbb{R}$. To simplify this writing we will simply write $\Gamma(a'_i \rightarrow b'_j)$, $\Gamma(a'_i \rightarrow \Lambda)$ or $\Gamma(\Lambda \rightarrow b'_j)$ although this will not be fully appropriate in general.

Definition 3.12: A function $\delta : \mathbb{U} \times \mathbb{U} \rightarrow \mathbb{R}$ is called an edit distance defined on \mathbb{U} if, for any pair of sequences A_1^p, B_1^q , the following recursive equation is satisfied

$$\delta(A_1^p, B_1^q) = \text{Min} \begin{cases} \delta(A_1^{p-1}, B_1^q) + \Gamma(a'_p \rightarrow \Lambda) & \text{delete} \\ \delta(A_1^{p-1}, B_1^{q-1}) + \Gamma(a'_p \rightarrow b'_q) & \text{match} \\ \delta(A_1^p, B_1^{q-1}) + \Gamma(\Lambda \rightarrow b'_q) & \text{insert} \end{cases} \quad (1)$$

Note that not all edit/elastic distances are metric. In particular, the dynamic time warping distance does not satisfy the triangle inequality.

3.5.1 Levenshtein distance

The Levenshtein distance $\delta_{lev}(x, y)$ has been defined for string matching. For this edit distance, the *delete* and *insert* operations induce unitary costs, i.e. $\Gamma(a'_p \rightarrow \Lambda) = \Gamma(\Lambda \rightarrow b'_q) = 1$ while the *match* cost is null if $a'_p = b'_q$ or 1 otherwise.

3.5.2 Dynamic time warping

The DTW similarity measure δ_{dtw} [25][19] is defined according to the previous notations as:

$$\delta_{dtw}(A_1^p, B_1^q) = d_{LP}(a_p, b_q) + \text{Min} \begin{cases} \delta_{dtw}(A_1^{p-1}, B_1^q) \\ \delta_{dtw}(A_1^{p-1}, B_1^{q-1}) \\ \delta_{dtw}(A_1^p, B_1^{q-1}) \end{cases} \quad (2)$$

where $d_{LP}(a_p, b_q)$ is the Lp norm in \mathbb{R}^k , and so for DTW, $\Gamma(a'_p \rightarrow \Lambda) = \Gamma(a'_p \rightarrow b'_q) = \Gamma(\Lambda \rightarrow b'_q) = d_{LP}(a_p, b_q)$. One may note that the time stamp values are not used, therefore the costs of each edit operation involve vectors a and b in S instead of vectors a' and b' in $S \times T$. One of the main restrictions of δ_{dtw} is that it does not comply with the triangle inequality as shown in [6].

3.5.3 Edit Distance with real penalty

$$\delta_{erp}(A_1^p, B_1^q) = \text{Min} \begin{cases} \delta_{erp}(A_1^{p-1}, B_1^q) + \Gamma(a'_p \rightarrow \Lambda) \\ \delta_{erp}(A_1^{p-1}, B_1^{q-1}) + \Gamma(a'_p \rightarrow b'_q) \\ \delta_{erp}(A_1^p, B_1^{q-1}) + \Gamma(\Lambda \rightarrow b'_q) \end{cases} \quad (3)$$

with

$$\begin{aligned} \Gamma(a'_p \rightarrow \Lambda) &= d_{LP}(a_p, g) \\ \Gamma(a'_p \rightarrow b'_q) &= d_{LP}(a_p, b_q) \\ \Gamma(\Lambda \rightarrow b'_q) &= d_{LP}(g, b_q) \end{aligned}$$

where g is a constant in S and $d_{LP}(x, y)$ is the Lp norm of vector $(x - y)$ in S .

Note that the time stamp coordinate is not taken into account, therefore δ_{erp} is a distance on S but not on $S \times T$. Thus, the cost of each edit operation involves vectors a and b in \mathbb{R}^k instead of vectors a' and b' in \mathbb{R}^{k+1} .

According to the authors of ERP [6], the constant g should be set to 0 for some intuitive geometric interpretation and in order to preserve the mean value of the transformed time series when adding gap samples.

3.5.4 Time warp edit distance

Time Warp Edit Distance (TWED) [13], [14] is defined similarly to the edit distance defined for string [12][26]. The similarity between any two time series A and B of finite length, respectively p and q is defined as:

$$\delta_{twed}(A_1^p, B_1^q) = \text{Min} \begin{cases} \delta_{twed}(A_1^{p-1}, B_1^q) + \Gamma(a'_p \rightarrow \Lambda) & \text{delete}_A \\ \delta_{twed}(A_1^{p-1}, B_1^{q-1}) + \Gamma(a'_p \rightarrow b'_q) & \text{match} \\ \delta_{twed}(A_1^p, B_1^{q-1}) + \Gamma(\Lambda \rightarrow b'_q) & \text{delete}_B \end{cases} \quad (4)$$

with

$$\begin{aligned} \Gamma(a'_p \rightarrow \Lambda) &= d(a'_p, a'_{p-1}) + \lambda \\ \Gamma(a'_p \rightarrow b'_q) &= d(a'_p, b'_q) + d(a'_{p-1}, b'_{q-1}) \\ \Gamma(\Lambda \rightarrow b'_q) &= d(b'_q, b'_{q-1}) + \lambda \end{aligned}$$

The time stamps are exploited to evaluate $d(a', b')$. In practice, $d(a', b') = d_{LP}(a, b) + \nu \cdot d_{LP}(t_a, t_b)$ where λ is a positive constant that represents a gap penalty and ν is a non negative constant which characterizes the *stiffness* of the δ_{twed} elastic measure.

4 INDEFINITENESS OF ELASTIC DISTANCE KERNELS

The Levenshtein distance kernel $\varphi(x, y) = \delta_{lev}(x, y)$ is known to be indefinite. Below, we discuss the first known counter-example produced by [7]. Let us consider the subset of sequences $V = \{abc, bad, dab, adc, bcd\}$ that leads to the following distance matrix

$$M_{lev}^V = \begin{pmatrix} 0 & 3 & 2 & 1 & 2 \\ 3 & 0 & 2 & 2 & 1 \\ 2 & 2 & 0 & 3 & 3 \\ 1 & 2 & 3 & 0 & 3 \\ 2 & 1 & 3 & 3 & 0 \end{pmatrix} \quad (5)$$

and consider coefficient vectors C and D in \mathbb{R}^5 such that

$C = [1, 1, -2/3, -2/3, -2/3]$ with $\sum_{i=1}^5 c_i = 0$ and $D = [1/3, 2/3, 1/3, -2/3, -2/3]$ with $\sum_{i=1}^5 d_i = 0$. Clearly $C \cdot M_{lev}^V \cdot C^T = 2/3 > 0$ and $D \cdot M_{lev}^V \cdot D^T = -4/3 < 0$, showing that M_{lev}^V has

no definiteness.

The DTW kernel $\varphi(x, y) = \delta_{dtw}(x, y)$ is also known to be neither CND nor CPD. The following example demonstrates this known result. Let us consider the subset of sequences $V = \{01, 012, 0123, 01234\}$.

Then the DTW empiric gram matrix evaluated on V is

$$M_{dtw}^V = \begin{pmatrix} 0 & 1 & 2 & 3 \\ 1 & 0 & 0 & 0 \\ 2 & 0 & 0 & 0 \\ 3 & 0 & 0 & 0 \end{pmatrix} \quad (6)$$

and consider coefficient vectors C and D in \mathbb{R}^4 such that

$C = [1/4, -3/8, -1/8, 1/4]$ with $\sum_{i=1}^4 c_i = 0$ and $D = [-1/4, -1/4, 1/4, 1/4]$ with $\sum_{i=1}^4 d_i = 0$. Clearly $C \cdot M_{dtw}^V \cdot C^T = 2/32 > 0$ and $D \cdot M_{dtw}^V \cdot D^T = -1/2 < 0$, showing that M_{dtw}^V has no definiteness.

Similarly, it is easy to find simple counter examples that show that neither ERP nor TWED kernels are definite.

Let us consider the subset of sequences $V = \{010, 012, 103, 301, 032, 123, 023, 003, 302, 321\}$.

For the TWED metric, with $\nu = 1.0$ and $\lambda = 0.0$ we get the following matrix:

$$M_{twed}^V = \begin{pmatrix} 0 & 2 & 7 & 9 & 6 & 7 & 5 & 5 & 10 & 9 \\ 2 & 0 & 5 & 9 & 4 & 5 & 3 & 3 & 8 & 9 \\ 7 & 5 & 0 & 6 & 7 & 4 & 6 & 2 & 5 & 10 \\ 9 & 9 & 6 & 0 & 13 & 10 & 12 & 8 & 1 & 4 \\ 6 & 4 & 7 & 13 & 0 & 5 & 3 & 5 & 12 & 9 \\ 7 & 5 & 4 & 10 & 5 & 0 & 2 & 6 & 9 & 6 \\ 5 & 3 & 6 & 12 & 3 & 2 & 0 & 4 & 11 & 8 \\ 5 & 3 & 2 & 8 & 5 & 6 & 4 & 0 & 7 & 10 \\ 10 & 8 & 5 & 1 & 12 & 9 & 11 & 7 & 0 & 5 \\ 9 & 9 & 10 & 4 & 9 & 6 & 8 & 10 & 5 & 0 \end{pmatrix} \quad (7)$$

The eigenvalue spectrum for this matrix is the following:

$\{4.62, 0.04, -2.14, -0.98, -0.72, -0.37, -0.19, -0.17, -0.06, -0.03\}$. This spectrum contains 2 strictly positive eigenvalues, showing that M_{twed}^V has no definiteness.

For the ERP metric, with $g = 0.0$ we get the following matrix:

$$M_{erp}^V = \begin{pmatrix} 0 & 2 & 3 & 3 & 4 & 5 & 4 & 2 & 4 & 5 \\ 2 & 0 & 3 & 5 & 2 & 3 & 2 & 2 & 4 & 5 \\ 3 & 3 & 0 & 4 & 3 & 2 & 3 & 1 & 3 & 4 \\ 3 & 5 & 4 & 0 & 7 & 6 & 7 & 5 & 1 & 2 \\ 4 & 2 & 3 & 7 & 0 & 3 & 2 & 2 & 6 & 5 \\ 5 & 3 & 2 & 6 & 3 & 0 & 1 & 3 & 5 & 4 \\ 4 & 2 & 3 & 7 & 2 & 1 & 0 & 2 & 6 & 5 \\ 2 & 2 & 1 & 5 & 2 & 3 & 2 & 0 & 4 & 5 \\ 4 & 4 & 3 & 1 & 6 & 5 & 6 & 4 & 0 & 1 \\ 5 & 5 & 4 & 2 & 5 & 4 & 5 & 5 & 1 & 0 \end{pmatrix} \quad (8)$$

The eigenvalue spectrum for this matrix is the following:

$\{4.63, 0.02, 1.39e - 17, -2.21, -0.97, -0.56, -0.41, -0.26, -0.17, -0.08\}$. This spectrum contains 3 strictly positive eigenvalues (although the third positive eigenvalue which is very small could be the result of the imprecision of the used diagonalization algorithm), showing that M_{erp}^V has no definiteness.

This shows that the metric properties of a distance defined on \mathbb{U} , in particular the triangle inequality, are not sufficient conditions to establish definiteness (conditionally or not) of the associated distance kernel. One could conjecture that elastic distances cannot be definite (conditionally or not), possibly because of the presence of the max or min operators in the recursive equation. In the following sections, we will see that replacing these min or max operators by a sum operator allows, under some conditions, for the construction of series of positive definite kernels whose limit is quite directly connected to the previously addressed elastic distance kernels.

5 CONSTRUCTING POSITIVE DEFINITE KERNELS FROM ELASTIC DISTANCE

The simple idea leading to the construction of positive definite kernels from a given elastic distance defined on \mathbb{U} is to replace the min or max operator into the recursive equation defining the elastic distance by a \sum operator. Instead of keeping one of the best alignment paths, the new kernel will sum up all the subsequence alignments with some weighting factor that could be optimized. This has been done successfully for the Smith and Waterman symbolic distance that is also known to be indefinite [18] and in the following sub sections, we propose some generalizations and extensions of this result.

5.1 Summative Time Warped Kernels

Definition 5.1: A function $\langle .; . \rangle : \mathbb{U} \times \mathbb{U} \rightarrow \mathbb{R}$ is called a Summative Time Warp Kernel (STWK) if, for any pair of sequences A_1^p, B_1^q , there exists a function $f : \mathbb{R} \rightarrow \mathbb{R}$ such that the following recursive equation is satisfied

$$\langle A_1^p; B_1^q \rangle = \sum \begin{cases} \langle A_1^{p-1}, B_1^q \rangle \star f(\Gamma(a'_p \rightarrow \Lambda)) & \text{delete} \\ \langle A_1^{p-1}, B_1^{q-1} \rangle \star f(\Gamma(a'_p \rightarrow b'_q)) & \text{match} \\ \langle A_1^p, B_1^{q-1} \rangle \star f(\Gamma(\Lambda \rightarrow b'_q)) & \text{insert} \end{cases} \quad (9)$$

Where \star is either the addition or the multiplication. Summative refers to the \sum operator replacing the min or max usually used. The recursion is initialized using $\langle A_1^0, B_1^0 \rangle = \langle \Omega, \Omega \rangle = \xi \in \mathbb{R}$.

This type of kernel sums up the multiplication or addition of the local quantities $f(\Gamma(a' \rightarrow b'))$ for all the possible alignment paths between the two time series.

Definition 5.2: If \star is the addition, the STWK is called additive, otherwise it will be called multiplicative.

The following theorem states necessary and sufficient conditions on $f(\Gamma(a' \rightarrow b'))$ for an STWK to be definite and thus is a basis for the construction of definite STWK.

Theorem 5.3: Definiteness of STWK:

- i) A STWK is positive definite on \mathbb{U} if the local kernel $k(a', b') = f(\Gamma(a' \rightarrow b'))$ is positive definite on $((S \times T) \cup \{\Lambda\})^2$ and $\xi > 0$.
- ii) An additive STWK is negative definite on \mathbb{U} if the local kernel $k(a', b') = f(\Gamma(a' \rightarrow b'))$ is negative definite on $((S \times T) \cup \{\Lambda\})^2$ and $\xi \leq 0$.
- iii) An additive STWK is conditionally positive definite if the local kernel $k(a', b') = f(\Gamma(a' \rightarrow b'))$ is conditionally positive definite on $((S \times T) \cup \{\Lambda\})^2$.

A sketch of proof for theorem 5.3 is given in the appendix.

As the cost function Γ is, in general, conditionally negative definite, choosing $f(h)$ for the exponential ensures that $f(\Gamma(a' \rightarrow b'))$ is a positive definite kernel [20]. Other functions can be used, such as the Inverse Multi Quadric kernel $k(a', b') = \frac{1}{\sqrt{(\Gamma(a' \rightarrow b'))^2 + \theta^2}}$. As with the exponential (Gaussian or Laplace) kernel, the Inverse Multi Quadric kernel results in a positive definite matrix with full rank [15] and thus forms a infinite dimension feature space.

5.2 Some instances of additive and multiplicative STWK

5.2.1 Additive STWK

Definition 5.4:

$$\langle A_1^p, B_1^q \rangle_{twip} = \frac{1}{3} \cdot \sum \begin{cases} \langle A_1^{p-1}, B_1^q \rangle_{twip} & \text{delete} \\ \langle A_1^{p-1}, B_1^{q-1} \rangle_{twip} + e^{-\nu \cdot d(t_{a_p}, t_{b_q})} (a_p \cdot b_q) & \text{match} \\ \langle A_1^p, B_1^{q-1} \rangle_{twip} & \text{insert} \end{cases}$$

where d is a distance, and ν a stiffness parameter. We suggest taking $\xi = 0$ for this additive SWTK.

Proposition 5.5: Definiteness of the additive STWK $\langle \dots \rangle_{twip}$ that has the property of an inner product:

- i) $\langle \dots \rangle_{twip}$ is positive definite.
- ii) Furthermore, $\langle \dots \rangle_{twip}$ is an inner product on $(\mathbb{U}, \oplus, \otimes)$, that we call a Time Warp Inner Product (TWIP), where \oplus and \otimes are defined in Algorithm 1 and definition 5.6 respectively,
- iii) The euclidean inner product $\langle \dots \rangle_{ED}$ on a set of time series of constant lengths and uniformly

sampled is the limit when $\nu \rightarrow \infty$ of $\langle \dots \rangle_{twip}$ on this same set.

Definition 5.6: For all $A \in \mathbb{U}$ and all $\lambda \in \mathbb{R}$, $C = \lambda \otimes A \in \mathbb{U}$ is such that for all $0 \leq i \leq |A|$, $c'_i = (\lambda \cdot a_i, t_{a_i})$ and thus $|C| = |A|$.

Algorithm 1 $A \oplus B$

For all A, B in \mathbb{U} , $C = A \oplus B \in \mathbb{U}$ is given by

```

i ← 1, j ← 1, k ← 1
WHILE i ≤ |A| OR j ≤ |B| DO
  IF i ≤ |A| AND j ≤ |B|
    IF tai = tbj
      c'k = (ai + bj, tai)
      i ← i + 1, j ← j + 1, k ← k + 1
    ELSE IF tai < tbj
      c'k = (ai, tai)
      i ← i + 1, k ← k + 1
    IF tai > tbj
      c'k = (ai, tbj)
      j ← j + 1, k ← k + 1
    END IF
  ENDIF
  IF i ≤ |A|
    c'k = (ai, tai)
    i ← i + 1, k ← k + 1
  END IF
  IF j ≤ |B|
    c'k = (bj, tbj)
    j ← j + 1, k ← k + 1
  END IF
END WHILE

```

Note that any discrete time series spaces of variable lengths and non uniformly sampled, when provided with the metric (norm) induced by a TWIP, is a Hilbert space.

The proof of proposition 5.5 is straightforward and is omitted.

5.2.2 Multiplicative exponentiated STWK

Definition 5.7:

$$\langle A_1^p, B_1^q \rangle_{me} = \frac{1}{3} \cdot \sum \begin{cases} \langle A_1^{p-1}, B_1^q \rangle_{me} \cdot e^{-\nu' \cdot \Gamma(a'_p \rightarrow \Lambda)} & \text{delete} \\ \langle A_1^{p-1}, B_1^{q-1} \rangle_{me} \cdot e^{-\nu' \cdot \Gamma(a'_p \rightarrow b'_q)} & \text{match} \\ \langle A_1^p, B_1^{q-1} \rangle_{me} \cdot e^{-\nu' \cdot \Gamma(\Lambda \rightarrow b'_q)} & \text{insert} \end{cases} \quad (10)$$

where ν' is a stiffness parameter that weighs the contribution of the local elementary costs. The larger ν' is, the more the kernel is selective around the optimal paths. At the limit, when $\nu' \rightarrow \infty$, only the optimal path costs are summed up by the kernel. Note that, as is generally seen, several optimal paths leading to the same global cost exist, $\lim_{\nu' \rightarrow +\infty} -1/\nu' \cdot \log(\langle A, B \rangle_{me})$ does

not coincide with the elastic distance δ that involves the same corresponding elementary costs.

We suggest set $\xi = 1$ for this kind of multiplicative ATWK.

Proposition 5.8: Definiteness of the multiplicative exponentiated STWK $\langle \cdot, \cdot \rangle_{me}$
 $\langle \cdot, \cdot \rangle_{me}$ is positive definite for the cost functions $\Gamma(a'_p \rightarrow \Lambda)$, $\Gamma(a'_p \rightarrow b'_q)$ and $\Gamma(\Lambda \rightarrow b'_q)$ involved in the computation of the δ_{lev} , δ_{dtw} , δ_{erp} and δ_{twed} distances. The multiplicative SWTKs constructed from these distances are referred respectively to $STWK_{lev}$, $STWK_{erp}$, $STWK_{dtw}$, $STWK_{twed}$ in the rest of the paper.

The proof of proposition 5.8 is straightforward and is omitted.

6 CLASSIFICATION EXPERIMENTS

We empirically evaluate the effectiveness of some STWK comparatively to Gaussian Radial Basis Function (RBF) Kernels or elastic distance substituting kernels [8] using some classification tasks on a set of time series coming from quite different application fields. The classification task we have considered consists of assigning one of the possible categories to an unknown time series for the 20 data sets available at the UCR repository [11]. As time is not explicitly given for these datasets, we used the index value of the samples as the time stamps for the whole experiment.

For each dataset, a training subset (TRAIN) is defined as well as an independent testing subset (TEST). We use the training sets to train two kinds of classifiers:

- the first one is a first near neighbor (1-NN) classifier: first we select a training data set containing time series for which the correct category is known. To assign a category to an unknown time series selected from a testing data set (different from the train set), we select its nearest neighbor (in the sense of a distance or similarity measure) within the training data set, then, assign the associated category to its nearest neighbor. For that experiment, a leave one out procedure is performed on the training dataset to optimized the meta parameters of the considered comparability measure.
- the second one is a SVM classifier [4], [24] configured with a Gaussian RBF kernel whose parameters are $C > 0$, a trade off between regularization and constraint violation and σ that determines the width of the Gaussian function. To determine the C and σ hyper parameter values, we adopt a 5-folded cross-validation method on each training subset. According to this procedure, given a predefined training set TRAIN and a test set TEST, we adapt the meta parameters based on the training set TRAIN: we first divide TRAIN into

5 stratified subsets $TRAIN_1, TRAIN_2, \dots, TRAIN_5$; then for each subset $TRAIN_i$ we use it as a new test set, and regard $(TRAIN - TRAIN_i)$ as a new training set; Based on the average error rate obtained on the five classification task, the optimal values of meta parameters are selected as the ones leading to the minimal average error rate.

We have used the LIBSVM library [5] to implement the SVM classifiers.

6.1 Additive STWK

We tested the additive STWK based on the Time Warp Inner Product $\langle A_1^p, B_1^q \rangle_{twip}$ (Eq.10). Precisely, we used the time-warp distance induced by $\langle A_1^p, B_1^q \rangle_{twip}$, basically $\delta_{twip}(A_1^p, B_1^q) = (\langle A_1^p - B_1^q, A_1^p - B_1^q \rangle_{twip})^{1/2}$.

6.1.1 Meta parameters

δ_{twip} is characterized by the meta parameter ν (the *stiffness* parameter) that is optimized for each dataset on the train data by minimizing the classification error rate of a first near neighbor classifier. For this kernel, ν is selected in $\{100, 10, 1, .1, .01, \dots, 1e-5\}$.

To explore the potential benefits of TWIP against the Euclidean inner product, we also tested the Euclidean Distance δ_{ed} which is the limit when $\nu \rightarrow \infty$ of δ_{twip} .

The kernels exploited by the SVM classifiers are the Gaussian kernels $STWK_{twid}(A, B) = e^{\delta_{twid}(A, B)^2 / (2 \cdot \sigma^2)}$ and $K_{ed}(A, B) = e^{\delta_{ed}(A, B)^2 / (2 \cdot \sigma^2)}$. The meta parameters C is selected into the discrete set $\{2^{-5}, 2^{-4}, \dots, 1, 2, \dots, 2^{10}\}$, and σ^2 into $\{2^{-5}, 2^{-4}, \dots, 1, 2, \dots, 2^{10}\}$.

Table 1 gives for each data set and each tested kernels (K_{ed} and $STWK_{twip}$) the corresponding optimized values of the meta parameters.

6.2 Multiplicative STWK

We tested the multiplicative exponentiated STWK based on the $\delta_{erp}, \delta_{dtw}, \delta_{twed}$ distance costs. We consider respectively the positive definite $STWK_{erp}, STWK_{dtw}, STWK_{twed}$ kernels.

Our experiment compares classification errors on the test data for

- the first near neighbor classifiers based on the $\delta_{erp}, \delta_{dtw}, \delta_{twed}$ distance measures (1-NN δ_{erp} , 1-NN δ_{dtw} and 1-NN δ_{twed}),
- the SVM classifiers using Gaussian distance substituting kernels based on the same distances and their corresponding STWK, e.g. SVM δ_{erp} , SVM $STWK_{erp}$, SVM δ_{dtw} , SVM $STWK_{dtw}$, SVM δ_{twed} , SVM $STWK_{twed}$.

For $\delta_{erp}, \delta_{twed}, STWK_{erp}$ and $STWK_{twed}$ we used the L1-norm, while the L2-norm has been implemented for δ_{dtw} and $STWK_{dtw}$, a classical choice for DTW [17].

TABLE 1
Dataset sizes and meta parameters used in conjunction with K_{ed} and $STWK_{twip}$ kernels

DATASET	length	#class	#train	#test	$K_{ed} : C, \sigma$	$STWK_{twip} : \nu, C, \sigma$
Synthetic control	60	6	300	300	1.0;0.125	0.1;2.0;0.0625
Gun-Point	150	2	50	150	256;1.0	1e-5;2048;2.0
CBF	128	3	30	900	8;1.0	0.01;1.0;0.0312
Face (all)	131	14	560	1690	4;0.5	1e-5;16.0;0.062
OSU Leaf	427	6	200	242	2;0.125	0.01;16;0.25
Swedish Leaf	128	15	500	625	128;0.125	1.0;16;0.125
50 Words	270	50	450	455	32;0.5	0.01;32;0.25
Trace	275	4	100	100	8;0.0156	1e-5;512;0.25
Two Patterns	128	4	1000	4000	4.0;0.25	0.01;2.0;0.0625
Wafer	152	2	1000	6174	4.0;0.5	0.1;8.0;0.0625
face (four)	350	4	24	88	8.0;2.0	1000;8.0;2.0
Ligthing2	637	2	60	61	2.0;0.125	.01;1.0;0.125
Ligthing7	319	7	70	73	32.0;256.0	0.1;512.0;8.0
ECCG	96	2	100	100	8.0;0.25	1000.0, 8.0, 0.25
Adiac	176	3	390	391	1024.0;0.125	1000;1024.0;0.125
Yoga	426	2	300	300	64.0;0.125	1.0;16.0;0.0625
Fish	463	7	175	175	64.0;1.0	10.0;64.0;1.0
Coffee	286	2	28	28	128.0;4.0	1000.0;128.0;4.0
OliveOil	570	4	30	30	2.0;0.125	0.01;4.0;0.125
Beef	470	5	30	30	128.0;4.0	1000.0;128.0;4.0

TABLE 2
Meta parameters used in conjunction with δ_{erp} , $STWK_{erp}$, δ_{dtw} , $STWK_{dtw}$, δ_{twd} and $STWK_{twd}$ kernels

DATASET	$\delta_{erp} : g; C; \sigma$	$STWK_{erp} : g; \nu'; C; \sigma$	$\delta_{dtw} : C; \sigma$	$STWK_{dtw} : \nu'; C; \sigma$	$\delta_{twd} : \lambda; \nu; C; \sigma$	$STWK_{twd} : \lambda; \nu; \nu'; C; \sigma$
Synth. cont.	0.0;2.0;0.25	0.0;0.457;256.0;0.062	8.0;4.0	0.047;1024.0;0.062	0.75;0.01;1.0;0.25	0.75;0.01;0.685;8.0;4.0
Gun-Point	-0.35;4.0;0.031	-0.35;0.457;128.0;1.0	16.0;0.0312	0.457;64.0;2.0	0.0;0.001;8.0;1.0	0.0;0.001;0.685;32;32
CBF	-0.11;1.0;1.0	-0.11;0.203;4.0;32.0	1.0;1.0	0.457;2.0;1.0	1.0;0.001;1.0;1.0	1.0;0.00;0.20;4.0;32.0
Face (all)	-1.96;4.0;0.5	-1.96;1.028;8.0;0.62	2.0;0.25	1.028;4.0;0.25	1.0;0.01;8.0;4.0	1.0;0.01;2.312;8.0;4.0
OSU Leaf	-2.25;2.0;0.062	-2.25;1.541;256;0.031	4.0;0.062	1.541;32.0;0.062	1.0;1e-4;8.0;0.25	1.0;1e-4;1.028;64.0;1.0
Swed. Leaf	0.3;8.0;0.125	0.3;0.203;1.0;4.0	4.0;0.031	5.202;0.062;0.5	1.0;1e-4;16.0;0.062	1.0;1e-4;0.304;32.0;1.0
50 Words	-1.39;16.0;0.25	-1.39;0.685;16;0.25	4.0;0.062	1.028;64.0;0.062	1.0;1e-3;8.0;0.5	1.0;1e-3;1.028;32.0;2.0
Trace	0.57;32;0.62	0.57;0.457;256;4.0	4;0.25	0.685;16;0.25	0.25;1e-3;8.0;0.25	0.25;1e-3;300;0.0625;0.25
Two Patt.	-0.89;0.25;0.125	-0.89;0.304;0.004;1.0	0.25;0.125	0.457;2.0;0.125	1.0;1e-3;0.25;0.125	1.0;1e-3;0.685;0.25;0.125
Wafer	1.23;2.0;0.062	1.23;0.685;4.0;0.5	1.0;0.016	1.541;1024;0.031	1.0;0.125;4.0;0.62	1.0;0.125;1.541;1.0;4.0
face (four)	1.97;64;16	1.97;0.685;32;2	16;0.5	0.457;16;2	1.0;0.01;4;2	1.0;0.01;1.027;4;2
Ligthing2	-0.33;2;0.062	-0.33;2.312;128;0.062	2.0;0.031	1.541;32;0.062	0.0;1e-6;8;0.25	0.0;1e-6;1.541;8;8
Ligthing7	-0.40;128;2	-0.40;0.685;32;0.25	4;0.25	0.685;32;0.062	0.25;0.1;4;0.5	0.25;0.1;0.685;4;8
ECCG	1.75;8;0.125	1.75;0.457;16;0.5	2;0.62	1.028;32;0.062	0.5;1.0;4;0.125	0.5;1.0;5.202;8;16
Adiac	1.83;16;0.0156	1.83;2.312;4096;0.031	16;0.0039	1.028;2048;0.031	0.75;1e-4;16;0.016	0.75;1e-4;2.312;128;1
Yoga	0.77;4;0.031	0.77;11.7054096;0.031	4;0.008	26.337;1024;0.031	0.5;1e-5;2;0.125	0.5;1e-5;3.468;256;2
Fish	-0.82;64;0.25	-0.82;0.685;32;0.5	8;0.016	3.468;64;16	0.5;1e-4;4;5	0.5;1e-4;0.457;16;16
Coffee	-3.00;16;0.062	-3.00;26.337;4096;16	8;0.062	5.202;512;4	0;0.1;16;4	0;0.1;300;1024;128
OliveOil	-3.00;8;0.5	-0.82;0.457;256;0.062	2;0.125	0.457;32;0.125	0;0.001;256;32	0;0.001;32;32
Beef	-3.00;128;0.125	-3.00;0.685;0.004;16384	16;0.016	0.457;0.004;16	0;1e-4;2;1	0;1e-4;0.135;0.004;16

TABLE 3

Comparative study using the UCR datasets: classification error rates (in %) obtained using the first near neighbor classification rule and a SVM classifier for the K_{ed} and $STWK_{twip}$ kernels

DATASET	1-NN ED	SVM K_{ed}	1-NN $STWK_{twip}$	SVM $STWK_{twip}$
Synthetic control	12	2	0.67	1.0
Gun-Point	8.67	6	8.67	4.67
CBF	14.8	11.5	2.4	3.66
Face (all)	28.6	16.62	24.44	24.08
OSU Leaf	48.3	43.8	44.63	43.8
Swedish Leaf	21.3	6.24	20	8.8
50 Words	36.9	30.10	31.86	29.23
Trace	24	24	23	8
Two Patterns	9	7.45	4.1	4.1
Wafer	0.52	0.52	0.43	0.54
face (four)	21.59	14.77	21.59	14.77
Ligthing2	24.6	32.78	18.03	26.22
Ligthing7	42.5	36.98	31.51	34.24
ECG	12	9	12	9
Adiac	38.87	24.81	38.87	24.81
Yoga	17	14.86	16.86	14.7
Fish	21.71	13.14	21.71	12.57
Coffee	25	0	17.85	0
OliveOil	13.33	13.33	13.33	13.33
Beef	46.67	30	46.67	30

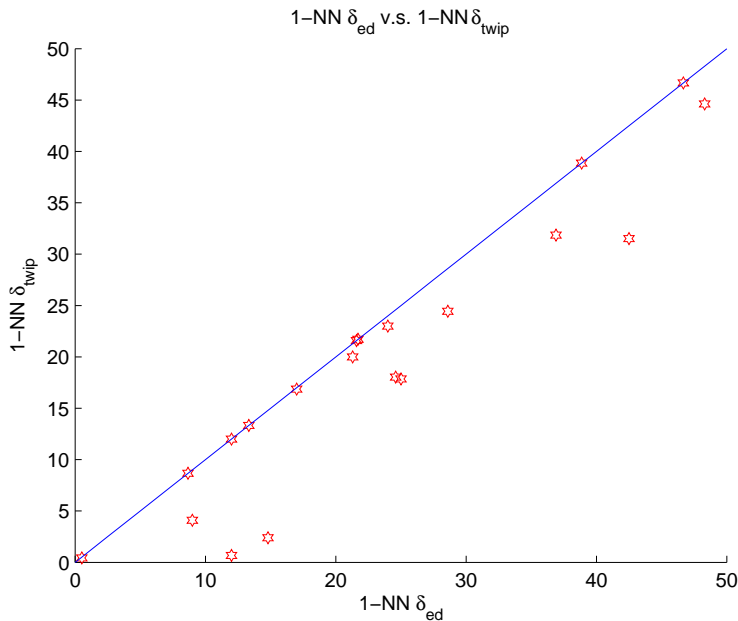


Fig. 1. Comparison of error rates (in %) between two 1-NN classifiers based on the Euclidean Distance (1-NN ED), δ_{ED} , and the distance δ_{TWIP} induced by a time-warp inner product (1-NN TWIP). The straight line has a slope of 1.0 and dots correspond, for the pair of classifiers, to the error rate on the tested data sets. A dot below (resp. above) the straight line indicates that distance δ_{TWIP} has a lower (resp. higher) error rate than distance δ_{ED}

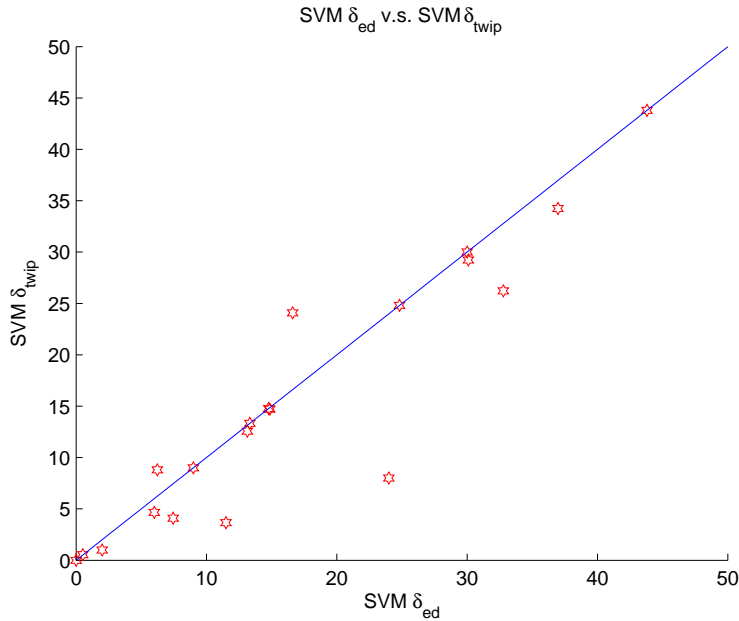


Fig. 2. Comparison of error rates (in %) between two SVM classifiers, the first one based on the Euclidean Distance gaussian kernel (SVM K_{ed}), and the second one based on a gaussian kernel induced by a time-warp inner product (SVM $STWK_{twip}$). The straight line has a slope of 1.0 and dots correspond, for the pair of classifiers, to the error rate on the tested data sets. A dot below (resp. above) the straight line indicates that SVM $STWK_{twip}$ has a lower (resp. higher) error rate than distance SVM K_{ed}

TABLE 4

Comparative study using the UCR datasets: classification error rates (in %) obtained using the first near neighbor classification rule and a SVM classifier for the erp , $STWK_{erp}$, dtw and $STWK_{dtw}$ kernels

DATASET	1-NN δ_{erp}	SVM δ_{erp}	SVM $STWK_{erp}$	1-NN δ_{dtw}	SVM δ_{dtw}	SVM $STWK_{dtw}$
Synthetic control	3.6	1.33	1.33	7	2.33	1
Gun-Point	4	1.33	1.33	9.3	10	1.33
CBF	0.3	3.55	3.22	0.3	1.44	5.44
Face (all)	20.2	18.1	16.86	19.2	15.32	16.98
OSU Leaf	39.7	35.95	30.57	40.9	43.8	23.55
Swedish Leaf	12	7.36	6.24	21	18.56	5.6
50 Words	28.1	24.61	16.04	31	29.45	17.58
Trace	17	1	1	0	0	2
Two Patterns	0	0	0	0	0	0
Wafer	0.9	0.89	0.44	2	2.95	0.39
face (four)	10.2	4.55	3.4	17	12.5	5.68
Ligthing2	14.8	18.03	19.67	13.1	24.59	19.67
Ligthing7	30.1	16.43	17.80	27.4	21.91	16.43
ECG	13	9	13	23	17	13
Adiac	37.8	30.94	24.04	39.6	34.52	25.32
Yoga	14.7	12.1	11.47	16.4	16.87	11.2
Fish	12	9.71	4.57	26.7	19.42	4.57
Coffee	25	17.85	14.29	17.9	7.14	17.85
OliveOil	16.7	16.67	16.67	13.3	16.67	16.67
Beef	50	46.67	50	50	50	42.85

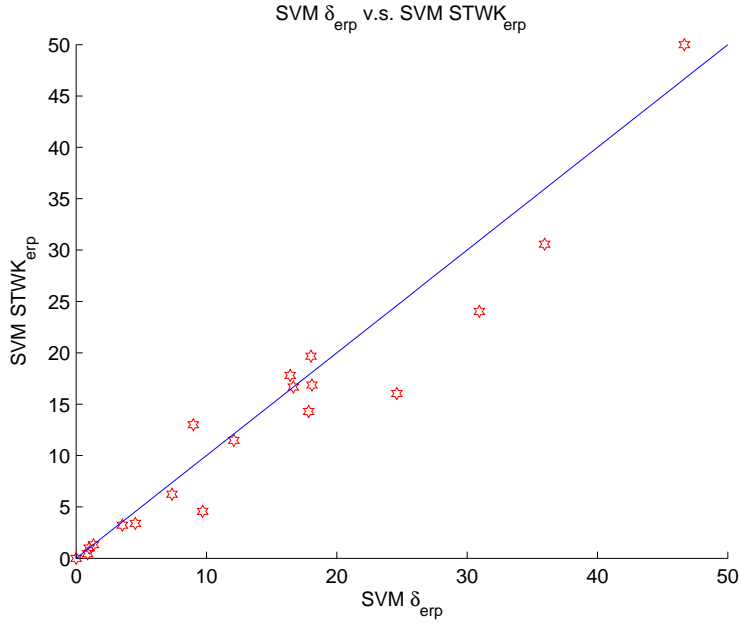


Fig. 3. Comparison of error rates (in %) between two SVM classifiers, the first one based on the δ_{erp} substituting kernel (SVM δ_{erp}), and the second one based on an additive time-warp kernel induced by the ERP distance (SVM $STWK_{erp}$). The straight line has a slope of 1.0 and dots correspond, for the pair of classifiers, to the error rate on the tested data sets. A dot below (resp. above) the straight line indicates that SVM $STWK_{erp}$ has a lower (resp. higher) error rate than distance SVM δ_{erp}

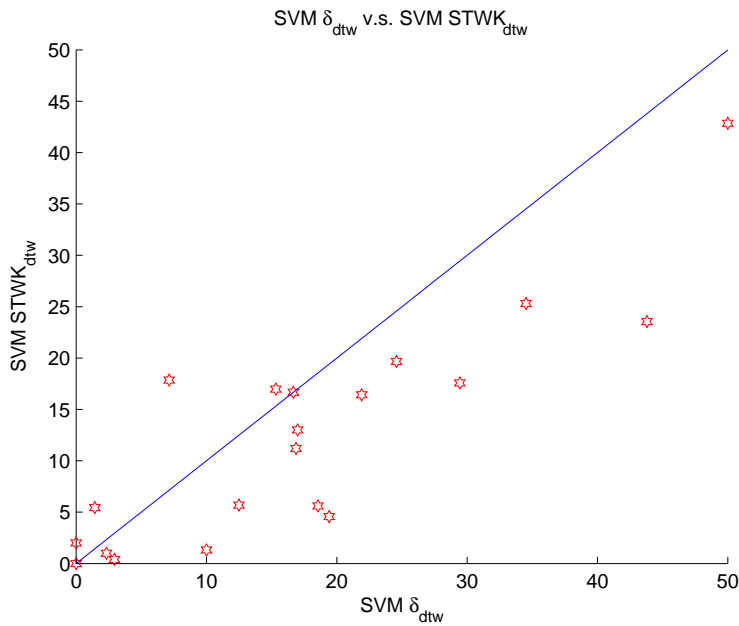


Fig. 4. Comparison of error rates (in %) between two SVM classifiers, the first one based on the δ_{dtw} substituting kernel (SVM δ_{dtw}), and the second one based on an additive time-warp kernel induced by the δ_{dtw} distance (SVM $STWK_{dtw}$). The straight line has a slope of 1.0 and dots correspond, for the pair of classifiers, to the error rate on the tested data sets. A dot below (resp. above) the straight line indicates that SVM $STWK_{dtw}$ has a lower (resp. higher) error rate than distance δ_{dtw}

TABLE 5

Comparative study using the UCR datasets: classification error rates (in %) obtained using the first near neighbor classification rule and a SVM classifier for the δ_{twed} and $STWK_{twed}$ kernels

DATASET	1-NN δ_{twed}	SVM δ_{twed}	SVM $STWK_{twed}$
Synthetic control	0	1.33	1.33
Gun-Point	2	0.067	0
CBF	0.67	3.5	2.44
Face (all)	23.6	16.86	15.09
OSU Leaf	28.1	18.18	22.73
Swedish Leaf	14.6	6.4	5.12
50 Words	18.9	14.51	16.26
Trace	9	0	1
Two Patterns	0.1	0.02	0
Wafer	1	0.4	0.37
face (four)	15.9	2.27	3.4
Ligthing2	19.7	21.31	21.31
Ligthing7	37	21.29	23.28
ECG	11	8	7
Adiac	41.7	31.02	23.52
Yoga	14	9.9	10.83
Fish	8.6	2.86	3.43
Coffee	28.5	28.5	17.86
OliveOil	16.7	13.33	13.33
Beef	33.3	53.33	50

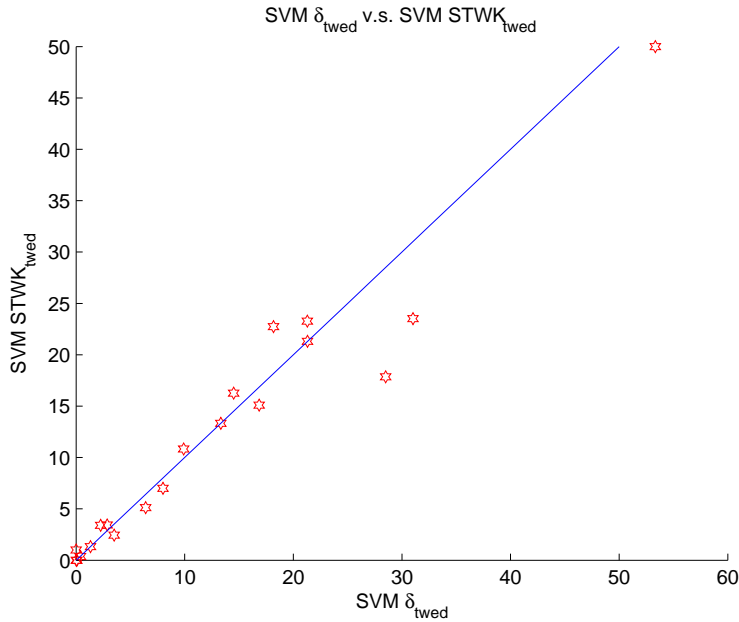


Fig. 5. Comparison of error rates (in %) between two SVM classifiers, the first one based on the δ_{twed} substituting kernel (SVM δ_{twed}), and the second one based on an additive time-warp kernel induced by the ERP distance (SVM $STWK_{twed}$). The straight line has a slope of 1.0 and dots correspond, for the pair of classifiers, to the error rate on the tested data sets. A dot below (resp. above) the straight line indicates that SVM $STWK_{twed}$ has a lower (resp. higher) error rate than distance SVM δ_{twed}

6.2.1 Meta parameters

For δ_{erp} kernel, meta parameter g is optimized for each dataset on the train data by minimizing the classification error rate of a first near neighbor classifier using a Leave One Out (LOO) procedure. For this kernel, g is selected in $\{-3, -2.99, -2.98, \dots, 2.98, 2.99, 3\}$. This optimized value is also used for comparison with the $STWK_{me}(ERP)$ kernel.

For δ_{twed} kernel, meta parameters λ and ν are optimized for each dataset on the train data by minimizing the classification error rate of a first near neighbor classifier. For our experiment, the *stiffness* value (ν) is selected from $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$ and λ is selected from $\{0, .25, .5, .75, 1.0\}$. If different (ν, λ) values lead to the minimal error rate estimated for the training data then the pairs containing the highest ν value are selected first, then the pair with the highest λ value is finally selected. These optimized (λ, ν) values are also used for comparability purposes with the $STWK_{twed}$ kernel.

The kernels exploited by the SVM classifiers are the Gaussian Radial Basis Function (RBF) kernels $K(A, B) = e^{\delta(A, B)^2 / (2 \cdot \sigma^2)}$ where δ stands for $\delta_{erp}, \delta_{dtw}, \delta_{twed}, STWK_{erp}(ERP), STWK_{dtw}, STWK_{twed}$. Meta parameter C is selected from $\{2^{-5}, 2^{-4}, \dots, 1, 2, \dots, 2^{10}\}$, and σ^2 from $\{2^{-5}, 2^{-4}, \dots, 1, 2, \dots, 2^{10}\}$. The best values are obtained using a cross validation procedure.

For the $STWK_{erp}, STWK_{dtw}$ and $STWK_{twed}$ kernels, meta parameter $1/\nu'$ is selected from the discrete set $S = \{10^{-5}, 10^{-4}, \dots, 1, 10, 100\}$.

The optimization procedure is as follows:

- for each value in S , we train a SVM $STWK_*$ classifier on the training dataset using the previously described 5-folded cross validation procedure to select the SVM meta parameters $cost$ and σ and the average of the classification error is recorded.
- the best σ, C and ν' values are the ones that lead to the minimal average error.

Table 2 gives for each data set and each tested kernel ($\delta_{erp}, \delta_{dtw}, \delta_{twed}, STWK_{erp}, STWK_{dtw}$ and $STWK_{twed}$) the corresponding optimized values of the meta parameters.

6.3 Discussion

6.3.1 Additive STWK experiment analysis

Table 3 shows the classification error rates obtained for the tested methods, e.g. the first near neighbor classifier based on the Euclidean Distance and the distance induced by the time-warp inner product (1-NN ED and 1-NN δ_{twip}), the Gaussian RBF kernel SVM based on the euclidean distance and the distance induced by the time-warp inner product (SVM K_{ed} and SVM $STWK_{twip}$).

This experiment shows that the time-warp inner product is significantly more effective for the considered tasks comparatively to the edit distance measure, since it exhibits, on average, the lowest error rates for the testing data for both the 1-NN and SVM classifiers, as shown in Table 3 and Figures 1 and 2. The stiffness parameter in δ_{twip} seems to play a significant role in these classification tasks, and this for a quite large majority of data sets.

6.3.2 Multiplicative STWK experiment analysis

Tables 4 and 5 show the classification error rates obtained for the tested methods, e.g. the first near neighbor classifier based on the $\delta_{erp}, \delta_{dtw}$ and δ_{twed} distances (1-NN δ_{erp} , 1-NN δ_{dtw} and 1-NN δ_{twed}), the Gaussian RBF kernel SVM based on the same distances (SVM δ_{erp} , SVM δ_{dtw} and SVM δ_{twed}) and euclidean distance and the Gaussian RBF kernel SVM based on the STWK kernels (SVM $STWK_{erp}$, SVM $STWK_{dtw}$ and SVM $STWK_{twed}$).

In this experiment, we show that the SVM classifiers clearly outperform the 1-NN classifiers. But the interesting results reported in tables 4 and 5 and figures 3, 4 and 5 is that SVM $STWK_{erp}$ and SVM $STWK_{erp}$ perform slightly better than SVM δ_{erp} and SVM δ_{twed} respectively, and the SVM $STWK_{dtw}$ is clearly much better than the SVM δ_{dtw} . This could come from the fact that δ_{erp} and δ_{twed} are metrics but not δ_{dtw} . SVM δ_{dtw} behaves poorly compared to the other tested classifiers probably because the SVM optimization process does not perform well. Nevertheless, the $STWK_{dtw}$ kernel based on δ_{dtw} seems to correct greatly its drawbacks.

7 CONCLUSION

Following the work on convolution kernels [10] and local alignment kernels defined for string processing around the Smith and Waterman algorithm [22] [18], we propose summative time-warp kernels (STWK) applicable for string and time series processing. We give some simple sufficient conditions to build positive definite STWK. Our generalization leads us to propose additive and multiplicative STWK. For multiplicative STWK, we show that, for the exponentiated version we have tested, the sufficient conditions are basically satisfied by classical elastic distances defined by a recursive equation. In particular this is true for the edit distance, the well known Dynamic Time Warping measure and for some variants such as the Edit Distance With Real penalty and the Time Warp Edit Distance, the latter two being metrics as well as the symbolic edit distance. From the general additive STWK definition we have suggested a time-warp inner product (TWIP) from which a metric (or norm) that generalizes the euclidean distance (or euclidean norm) is induced. The experiments conducted on a variety of time series datasets show that both the multiplicative positive definite STWKs outperform the indefinite elastic

distances they are derived from when considering 1-NN and SVM classification tasks. Our experiments also show that the additive STWK we constructed from the proposed instance of TWIP significantly outperforms the kernels derived from the euclidean inner product.

This time-warp inner product opens some interesting perspectives since it leads to reconsidering the notion of orthogonality in discrete time series spaces. In particular, in such spaces provided with a TWIP the discrete sine and cosine waveforms are no longer orthogonal. If so, what may look like a discrete elastic Fourier transform ?

APPENDIX A

A.1 Proof of theorem 5.3

i) Let us show that if the function $f(\Gamma(a' \rightarrow b')) : (S \times T) \cup \{\Lambda\} \rightarrow \mathbb{R}$ is positive definite and if $\xi > 0$, then an additive or multiplicative STWK is definite positive.

To that end, we consider the set $\mathbb{V}_r = \{V \subset \mathbb{U} \text{ such that } \text{Max}_{(A_k, A_l) \in V^2} (|A_k| + |A_l|) \leq r\}$, and show by induction on r that, **P1**: if the previous conditions are satisfied, for all $r \in \mathbb{R}^+ \cup \{0\}$, all $V = \{A_1, A_2, \dots, A_{|V|}\} \in \mathbb{V}_r$, all $(\alpha_1, \alpha_2, \dots, \alpha_{|V|}) \in \mathbb{R}^{|V|}$, $\sum_{k,l} \alpha_k \alpha_l \langle A_k, A_l \rangle \geq 0$.

Base case (BC): Proposition P1 is obviously true for $r = 0$ since we restrict the sequence set to $\mathbb{U}_0 = \{\Omega\}$ since $\xi > 0$.

Inductive Hypothesis (IH): Let us suppose that proposition P1 is true for any integer k such that $k \leq r$ where $r \geq 0$ and let show that it is verified for any k such that $k \leq r + 1$.

Let us first note $\langle A_k, A_l \rangle_{i,j} = \langle A_{k,1}^i, A_{l,1}^j \rangle$ the restriction of the STWK up to index $0 \leq i \leq |A_k|$ in A_k and index $0 \leq j \leq |A_l|$ in A_l .

For all finite subset $V = \{A_1, A_2, \dots, A_{|V|}\} \in \mathbb{V}_{r+1}$ and all $(\alpha_1, \alpha_2, \dots, \alpha_{|V|}) \in \mathbb{R}^{|V|}$ we have

$$\begin{aligned} & \sum_{k,l} \alpha_k \alpha_l \langle A_k, A_l \rangle = \\ & \sum_{m,n/|A_m|,|A_n| \geq 1} \alpha_m \alpha_n \langle A_m, A_n \rangle \\ & + \sum_{p,q/|A_p|=0,|A_q| \geq 1} \alpha_p \alpha_q \langle A_p, A_q \rangle \\ & + \sum_{p,q/|A_p|=0,|A_q| \geq 1} \alpha_p \alpha_q \langle A_q, A_p \rangle \\ & + \sum_{r,s/|A_r|=|A_s|=0} \alpha_r \alpha_s \langle A_r, A_s \rangle \end{aligned}$$

thus, by definition of the STWK

$$\Gamma_{a'_i \rightarrow b'_j}(A_1^p, B_1^q)$$

$$\sum_{k,l} \alpha_k \alpha_l \langle A_k, A_l \rangle =$$

$$\begin{aligned} & \sum_{m,n/|A_m| \geq 1, |A_n| \geq 1} \alpha_m \alpha_n \langle A_m, A_n \rangle_{|A_m|, |A_n|-1} \star \\ & f(\Gamma_{\Lambda \rightarrow a'_{|A_n|}}(A_m, A_n)) \\ & + \sum_{m,n/|A_m| \geq 1, |A_n| \geq 1} \alpha_m \alpha_n \langle A_m, A_n \rangle_{|A_m|-1, |A_n|-1} \star \\ & f(\Gamma_{a'_{|A_m|} \rightarrow a'_{|A_n|}}(A_m, A_n)) \\ & + \sum_{m,n/|A_m| \geq 1, |A_n| \geq 1} \alpha_m \alpha_n \langle A_m, A_n \rangle_{|A_m|-1, |A_n|} \star \\ & f(\Gamma_{a'_{|A_m|} \rightarrow \Lambda}(A_m, A_n)) \\ & + \sum_{p,q/|A_p|=0, |A_q| \geq 1} \alpha_p \alpha_q \langle \Omega, A_q \rangle_{0, |A_q|-1} \star \\ & f(\Gamma_{\Lambda \rightarrow a'_{|A_q|}}(A_p, A_q)) \\ & + \sum_{p,q/|A_p|=0, |A_q| \geq 1} \alpha_p \alpha_q \langle A_q, \Omega \rangle_{|A_q|-1, 0} \star \\ & f(\Gamma_{a'_{|A_q|} \rightarrow \Lambda}(A_q, A_p)) \\ & + \sum_{r,s/|A_r|=|A_s|=0} \alpha_r \alpha_s \cdot \xi \end{aligned}$$

Within each term (except the last one) of the sum constituting the right hand side of the previous equality, the \star operator associates two kernels whose arguments are identical. The restricted STWK appearing within these terms are positive definite since they all apply on subsets that belong to some \mathbb{V}_s with $s \leq r$ and which are by IH positive definite. As \star is either the addition or the multiplication and as positive definite kernels are closed under summation or multiplication [3], all these terms are positive. The last term being obviously positive, we establish that $\sum_{k,l} \alpha_k \alpha_l \langle A_k, A_l \rangle \geq 0$. Thus proposition P1

is true at $r + 1$. By induction, proposition P1 is true for all $r \in \mathbb{R}^+ \cup \{0\}$. Finally, we have established that for all finite subset $V = \{A_1, A_2, \dots, A_{|V|}\} \in \mathbb{V}_{r+1}$ and all $(\alpha_1, \alpha_2, \dots, \alpha_{|V|}) \in \mathbb{R}^{|V|}$ we have $\sum_{k,l} \alpha_k \alpha_l \langle A_k, A_l \rangle \geq 0$, e.g. the STWK $\langle \dots \rangle$ is positive definite \square

ii) and iii) are proved in a very similar way.

A.2 Proof of proposition 5.5

The proofs of i) and ii) are obtained using a similar recursion as the one used to prove theorem 5.3. iii) is immediate.

ACKNOWLEDGMENTS

REFERENCES

- [1] S. Altschul, W. Gish, W. Miller, E. Myers, and D. Lipman. A basic local alignment search tool. *Journal of Molecular Biology*, 215:403–410, 1990.

- [2] R. Bellman. *Dynamic Programming*. Princeton Univ Press, 1957. New Jersey.
- [3] Christian Berg, Jens Peter Reus Christensen, and Paul Ressel. *Harmonic Analysis on Semigroups: Theory of Positive Definite and Related Functions*, volume 100 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, April 1984.
- [4] Bernhard E. Boser, Isabelle Guyon, and Vladimir Vapnik. A training algorithm for optimal margin classifiers. In *COLT*, pages 144–152, 1992.
- [5] Chih-Chung Chang and Chih-Jen Lin. *LIBSVM: a library for support vector machines*, 2001.
- [6] L. Chen and R. Ng. On the marriage of lp-norm and edit distance. In *Proceedings of the 30th International Conference on Very Large Data Bases*, pages 792–801, 2004.
- [7] Corinna Cortes, Patrick Haffner, and Mehryar Mohri. Positive Definite Rational Kernels. In *Proceedings of COLT'03*, volume 2777 of *Lecture Notes in Computer Science*, pages 41–56, Washington D.C., August 2003. Springer, Heidelberg, Germany.
- [8] B. Haasdonk and C. Bahlmann. Learning with distance substitution kernels. In *DAGM-Symposium*, pages 220–227, 2004.
- [9] Bernard Haasdonk. Feature space interpretation of svms with indefinite kernels. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27:482–492, 2005.
- [10] D. Haussler. Convolution kernels on discrete structures. Technical report, University of California, Santa Cruz, 2008. Technical Report.
- [11] E. J. Keogh, X. Xi, L. Wei, and C.A. Ratanamahatana. The ucr time series classification-clustering datasets, 2006. http://www.cs.ucr.edu/~eamonn/time_series_data/.
- [12] V. I. Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Doklady Akademii Nauk SSSR*, 163(4):845–848, 1965 (Russian), pages 707–710, 1966. English translation in *Soviet Physics Doklady*, 10(8).
- [13] P. F. Marteau. Time warp edit distance. Technical report, VALORIA, Universite de Bretagne Sud, 2008. Technical Report valoriaUBS-2008-3v, <http://hal.archives-ouvertes.fr/hal-00258669/fr/>.
- [14] P. F. Marteau. Time warp edit distance with stiffness adjustment for time series matching. *IEEE Trans. Pattern Anal. Mach. Intell.*, 31(2):306–318, 2009.
- [15] C. A. Micchelli. Interpolation of scattered data: distance matrices and conditionally positive functions. *Constructive Approximation*, 2:11–22, 1986.
- [16] W. Pearson. Rapid and sensitive sequence comparisons with fastp and fasta. *Methods Enzymol*, 183:63–98, 1990.
- [17] C. A. Ratanamahatana and E. J. Keogh. Making time-series classification more accurate using learned constraints. In *Proceedings of the Fourth SIAM International Conference on Data Mining (SDM'04)*, pages 11–22, 2004.
- [18] H. Saigo, J.P. Vert, N. Ueda, and T. Akutsu. Protein homology detection using string alignment kernels. *Bioinformatics*, 20:1682–1689, 2004.
- [19] H. Sakoe and S. Chiba. A dynamic programming approach to continuous speech recognition. In *Proceedings of the 7th International Congress of Acoustic*, pages 65–68, 1971.
- [20] I. J. Schoenberg. Metric spaces and positive definite functions. *Transactions of the American Mathematical Society*, 44(3):522–536, nov 1938.
- [21] Bernhard Scholkopf and Alexander J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge, MA, USA, 2001.
- [22] T. Smith and Waterman M. Identification of common molecular subsequences. *Journal of Molecular Biology*, 147:195–197, 1981.
- [23] Vladimir Vapnik. *Statistical Learning Theory*. Wiley-Interscience, ISBN 0-471-03003-1, 1989.
- [24] Vladimir Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag, 1995.
- [25] V. M. Velichko and N. G. Zagoruyko. Automatic recognition of 200 words. *International Journal of Man-Machine Studies*, 2:223–234, 1970.
- [26] R. A. Wagner and M. J. Fischer. The string-to-string correction problem. *Journal of the ACM (JACM)*, 21:168–173, 1973.
- [27] Adam Woznica, Alexandros Kalousis, and Melanie Hilario. Distances and (indefinite) kernels for sets of objects. *Data Mining, IEEE International Conference on*, 0:1151–1156, 2006.
- [28] Gang Wu, Edward Y. Chang, and Zhihua Zhang. Learning with non-metric proximity matrices. In *MULTIMEDIA '05: Proceedings of the 13th annual ACM international conference on Multimedia*, pages 411–414, New York, NY, USA, 2005. ACM.