

# An Empirical Evaluation of Similarity Measures for Time Series Classification

Joan Serrà, Josep Ll. Arcos

*Artificial Intelligence Research Institute (IIA-CSIC),  
Spanish National Research Council,  
08193 Bellaterra, Barcelona, Spain.*

---

## Abstract

Time series are ubiquitous, and a measure to assess their similarity is a core part of many computational systems. In particular, the similarity measure is the most essential ingredient of time series clustering and classification systems. Because of this importance, countless approaches to estimate time series similarity have been proposed. However, there is a lack of comparative studies using empirical, rigorous, quantitative, and large-scale assessment strategies. In this article, we provide an extensive evaluation of similarity measures for time series classification following the aforementioned principles. We consider 7 different measures coming from alternative measure ‘families’, and 45 publicly-available time series data sets coming from a wide variety of scientific domains. We focus on out-of-sample classification accuracy, but in-sample accuracies and parameter choices are also discussed. Our work is based on rigorous evaluation methodologies and includes the use of powerful statistical significance tests to derive meaningful conclusions. The obtained results show the equivalence, in terms of accuracy, of a number of measures, but with one single candidate outperforming the rest. Such findings, together with the followed methodology, invite researchers on the field to adopt a more consistent evaluation criteria and a more informed decision regarding the baseline measures to which new developments should be compared.

*Keywords:* Time Series, Similarity, Classification, Evaluation

---

## 1. Introduction

2 Data in the form of time series pervades a large number of scientific do-  
3 mains (Keogh, 2011; Keogh et al., 2011). Observations that unfold over time

4 usually represent valuable information subject to analysis, classification, in-  
5 dexing, prediction, or interpretation (Kantz and Schreiber, 2004; Han and  
6 Kamber, 2005; Liao, 2005; Fu, 2011). Real-world examples include finan-  
7 cial data (e.g., stock market fluctuations), medical data (e.g., electrocardio-  
8 grams), computer data (e.g., log sequences), or motion data (e.g., location  
9 of moving objects). Even object shapes or handwriting can be effectively  
10 transformed into time series, facilitating their analysis and retrieval (Keogh  
11 et al., 2011, 2009).

12 A core issue when dealing with time series is determining their pair-  
13 wise similarity, i.e., the degree to which a given time series resembles an-  
14 other. In fact, a time series similarity (or dissimilarity) measure is central to  
15 many mining, retrieval, clustering, and classification tasks (Han and Kam-  
16 ber, 2005; Liao, 2005; Fu, 2011; Keogh and Kasetty, 2003). Furthermore,  
17 there is evidence that simple approaches to such tasks exploiting generic  
18 time series similarity measures usually outperform more elaborate, some-  
19 times specifically-targeted strategies. This is the case, for instance, with  
20 time series classification, where a one-nearest neighbor approach using a  
21 well-known time series similarity measure was found to outperform an ex-  
22 haustive list of alternatives (Xi et al., 2006), including decision trees, multi-  
23 scale histograms, multi-layer perceptron neural networks, order logic rules  
24 with boosting, or multiple classifier systems.

25 Deriving a measure that correctly reflects time series similarities is not  
26 straightforward. Apart from dealing with high dimensionality (time series  
27 can be roughly considered as multi-dimensional data; Han and Kamber,  
28 2005), the calculation of such measures needs to be fast and efficient (Keogh  
29 and Kasetty, 2003). Indeed, with better information gathering tools, the size  
30 of time series data sets may continue to increase in the future. Moreover,  
31 there is the need for generic/multi-purpose similarity measures, so that they  
32 can be readily applied to any data set, whether this application is the final  
33 goal or just an initial approach to a given task. This last aspect highlights  
34 another desirable quality for time series similarity measures: their robustness  
35 to different types of data (cf. Keogh and Kasetty, 2003; Wang et al., 2012).

36 Over the years, several time series similarity measures have been pro-  
37 posed (for pointers to such measures see, e.g., Liao, 2005; Fu, 2011; Wang  
38 et al., 2012). Nevertheless, few quantitative comparisons have been made  
39 in order to evaluate their efficacy in a multiple-data framework (Keogh and  
40 Kasetty, 2003). Apart from being an interesting and important task by it-  
41 self (Keogh, 2011), and as opposed to clustering (Liao, 2005), time series  
42 classification offers the possibility to straightforwardly assess the merit of  
43 time series similarity measures under a controlled, objective, and quantita-

44 tive framework.

45 In a recent study, Wang et al. (2012) perform an extensive comparison of  
46 classification accuracies for 9 measures (plus 4 variants) across 38 data sets  
47 coming from various scientific domains. One of the main conclusions of the  
48 study is that, even though the newly proposed measures can be theoretically  
49 attractive, the efficacy of some common and well-established measures is,  
50 in the vast majority of cases, very difficult to beat. Specifically, dynamic  
51 time warping (DTW; Berndt and Clifford, 1994) is found to be consistently  
52 superior to the other studied measures (or, at worst, for a few data sets,  
53 equivalent). In addition, the authors emphasize that the Euclidean distance  
54 remains a quite accurate, robust, simple, and efficient way of measuring the  
55 similarity between two time series. Finally, by looking in detail at the results  
56 presented by Wang et al. (2012), we can spot a group of time series similarity  
57 measures that seems to have an efficacy comparable to DTW: those based  
58 on edit distances. In particular, the edit distance for real sequences (EDR;  
59 Chen et al., 2005) seems to be very competitive, if not slightly better than  
60 DTW. Interestingly, none of the three measures above was initially targeted  
61 to generic time series data, but were introduced with hindsight (Agrawal  
62 et al., 1993; Berndt and Clifford, 1994; Chen et al., 2005). The intuition  
63 behind Euclidean distance relates to spatial proximity, DTW was initially  
64 devised for the specific task of spoken word recognition (Sakoe and Chiba,  
65 1978), and edit distances were introduced for measuring the dissimilarity  
66 between two strings (Levenshtein, 1966).

67 The study by Wang et al. (2012) is, to the best of our knowledge, the  
68 only comparative study dealing with time series classification using multiple  
69 similarity measures and a large collection of data. In general, the studies  
70 introducing a new measure only compare against a few other measures<sup>1</sup>, and  
71 usually using a reduced data set corpus (cf. Keogh and Kasetty, 2003). Fur-  
72 thermore, there is a lack of agreement in the literature regarding evaluation  
73 methodologies. Besides, statistical significance is usually not studied or, at  
74 best, improperly evaluated. This is very inconvenient, as robust evaluation  
75 methodologies and statistical significance are the principal tools by which  
76 we can establish, in a formal and rigorous way, differences across the consid-  
77 ered measures (Salzberg, 1997; Hollander and Wolfe, 1999; Demšar, 2006).  
78 In addition, the optimal parameter values for every measure are rarely dis-  
79 cussed. All these issues impact the scientific development of the field as one  
80 is never sure, e.g., of which measure should be used as a baseline for future

---

<sup>1</sup>In the majority of cases, as our results will show, not the most appropriate ones.

81 developments, or of which parameters are the most sensible choice.

82 In this work, we perform an empirical evaluation of similarity measures  
83 for time series classification. We follow the initiative by Wang et al. (2012),  
84 and consider a big pool of publicly-available time series data sets (45 in our  
85 case). However, instead of additionally focusing on representation meth-  
86 ods, computational/storage demands, or more theoretical issues, we here  
87 take a pragmatic approach and restrict ourselves to classification accuracy.  
88 We believe that this is the most important aspect to be considered in a  
89 first stage and that, in contrast to the other aforementioned issues, it is  
90 not sufficiently well-covered in the existing literature. As for the consid-  
91 ered measures, we decide to include DTW and EDR, as these were found  
92 to generally achieve the highest accuracies among all measures compared  
93 in Wang et al. (2012). Apart from these two, we choose the Euclidean dis-  
94 tance plus 4 different measures not considered in such study, making up to  
95 a total of 7. Further important contributions that differentiate the current  
96 work from previous studies include (a) an extensive summary and back-  
97 ground of the considered measures, with basic formulations, applications,  
98 and references, (b) the formalization of a robust evaluation methodology,  
99 exploiting standard out-of-sample cross-validation strategies, (c) the use of  
100 rigorous statistical significance tests in order to assess the superiority of a  
101 given measure, (d) the evaluation of both train and test accuracies, and (e)  
102 the assessment of the optimal parameters for each measure and data set.

103 The rest of the paper is organized as follows. Firstly, we provide the  
104 background on time series similarity measures, outline some of their appli-  
105 cations, and detail their calculation (Sec. 2). Next, we explain the proposed  
106 evaluation methodology (Sec. 3). Subsequently, we report the obtained re-  
107 sults (Sec. 4). A conclusion section ends the paper (Sec. 5).

## 108 2. Time series similarity measures

109 The list of approaches for dealing with time series similarity is vast,  
110 and a comprehensive enumeration of them all is beyond the scope of the  
111 present work (for that, the interested reader is referred to Gusfield, 1997;  
112 Wang et al., 2012; Han and Kamber, 2005; Liao, 2005; Marteau, 2009; Fu,  
113 2011). In this section, we present several representative examples of different  
114 ‘families’ of time series similarity measures (Liao, 2005; Wang et al., 2012):  
115 lock-step measures (Euclidean distance), feature-based measures (Fourier  
116 coefficients), model-based measures (auto-regressive), and elastic measures  
117 (DTW, EDR, TWED, and MJC). An effort has been made in selecting the  
118 most standard measures of each group, emphasizing the approaches that are

reported to have good performance. We also try to avoid measures with too many parameters, since such parameters may be difficult to learn in small training data sets and, furthermore, could lead to over-fitting. Alternative measures found to be consistently less accurate than DTW or EDR are not considered (Wang et al., 2012). Apart from all the aforementioned measures, we also include a random measure, consisting of a uniformly distributed random number between 0 and 1. This will act as our random baseline.

### 2.1. Euclidean distance

The simplest way to estimate the dissimilarity between two time series is to use any  $L_n$  norm such that

$$d_{L_n}(\mathbf{x}, \mathbf{y}) = \left( \sum_{i=1}^M (x_i - y_i)^n \right)^{\frac{1}{n}}, \quad (1)$$

where  $n$  is a positive integer,  $M$  is the length of the time series, and  $x_i$  and  $y_i$  are the  $i$ -th element of time series  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. Measures based on  $L_n$  norms correspond to the group of so-called lock-step measures (Wang et al., 2012), which compare samples that are at exactly the same temporal location (Fig. 1, top). Notice that in case the time series  $\mathbf{x}$  and  $\mathbf{y}$  not being of the same length, one can always re-sample one to the length of the other, an approach that works well for a number of data sources (Keogh and Kasetty, 2003).

Using Eq. 1 with  $n = 2$  we obtain the Euclidean distance, one of the most used time series dissimilarity measures, favored by its computational simplicity and indexing capabilities. Applications range from early classification of time series (Xing et al., 2011) to rule discovery in economic, communications, and ecological time series (Das et al., 1998). Some authors state that the accuracy of the Euclidean distance can be very difficult to beat, specially for large data sets containing many time series (cf. Wang et al., 2012). To the best of our knowledge, these claims are only quantitatively supported by one-nearest neighbor classification experiments using two artificially-generated/synthetic data sets (Geurts, 2002). We believe that such claims need to be carefully assessed with extensive experiments and under broader conditions, considering multiple measures, different distance-exploiting algorithms, and real-world data sets.

### 2.2. Fourier coefficients

A simple extension of the Euclidean distance is not to compute it directly using the raw time series, but using features extracted from it. For instance,

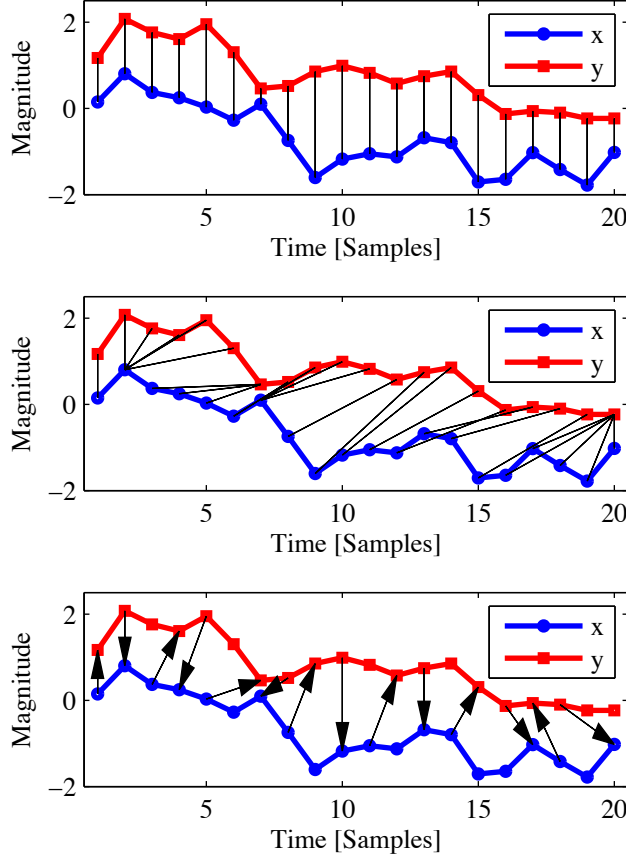


Figure 1: Examples of dissimilarity calculations between time series  $\mathbf{x}$  and  $\mathbf{y}$ : Euclidean distance (top), DTW alignment (center), and MJC (bottom). See text for details.

153 by first representing the time series by their Fourier coefficients (FC), one  
 154 uses

$$d_{\text{FC}}(\mathbf{x}, \mathbf{y}) = \left( \sum_{i=1}^{\theta} (\hat{x}_i - \hat{y}_i)^2 \right)^{\frac{1}{2}}, \quad (2)$$

155 where  $\hat{x}_i$  and  $\hat{y}_i$  are complex value pairs denoting the  $i$ -th Fourier coefficient  
 156 of  $\hat{\mathbf{x}}$  and  $\hat{\mathbf{y}}$ , the discrete Fourier transforms (DFT) of the raw time series (Op-  
 157 penheim et al., 1999). Notice that in Eq. 2 we introduce the parameter  $\theta$ ,  
 158 the actual number of considered coefficients. Because of the symmetry of  
 159 the DFT, the sum only needs to be performed, at most, over half of the

coefficients, so that  $\theta = M/2$ . Notice that, by the Parseval theorem (Oppenheim et al., 1999), the Euclidean distance between FCs is equivalent to the standard Euclidean distance between the raw time series (see, e.g., Agrawal et al., 1993). However, having parameter  $\theta$ , one usually takes the opportunity to filter out high-frequency coefficients, i.e., coefficients  $\hat{x}_i$  and  $\hat{y}_i$  whose  $i$  is close to  $M/2$ . This has the (sometimes desired) effect of removing rapidly-fluctuating components of the signal. Hence, if high frequencies are not relevant for the intended analysis or we have some high-frequency noise, this operation will usually carry some increase in accuracy. Furthermore, if  $\theta$  is relatively small, similarity computations can be substantially accelerated.

Computing the Euclidean distance on a reduced set of features is an extremely common approach in literature. FCs are the standard choice for efficient time series retrieval, exploiting the aforementioned acceleration capabilities. Pioneering work includes Agrawal et al. (Agrawal et al., 1993) and Faloutsos et al. (Faloutsos et al., 1994) dealing with synthetic and financial data. More recent works use FCs with data from other domains. For instance, the case-based reasoning system of Montani et al. (Montani et al., 2006) uses FCs to compare medical time series. Apart from FCs, wavelet coefficients have been extensively used (Chan and Fu, 1999). For instance, Olsson et al. (Olsson et al., 2004) use a wavelet analysis to remove noise and extract features in their system of fault diagnosis in industrial equipment. Research suggests that, although they provide some advantages, wavelet coefficients do not generally outperform FCs for the considered task (Wu et al., 2000). Comparatively less used time series features are based on singular value decomposition (Wu et al., 1996), piece-wise aggregate approximations (Keogh et al., 2001), or the coefficients of fitted polynomials (Cai and Ng, 2004) among others.

### 2.3. Auto-regressive models

A further option for computing similarities between time series using features extracted from them is to employ time series models (Liao, 2005; Fu, 2011). The main idea behind model-based measures is to learn a model of the two time series and then use its parameters for computing a similarity value. In the literature, several approaches follow this idea. For instance, Maharaj (2000) uses the  $p$ -value of a chi-square statistic to cluster auto-regressive coefficients representing stationary time series. Ramoni et al. (2002) present a Bayesian algorithm for clustering time series. They transform each series into a Markov chain and then cluster similar chains to discover the most probable set of generating processes. Pavinelli et al.

(2004) use Gaussian mixture models of reconstructed phase spaces to classify time series of different sources. Serrà et al. (2012a) study the use of the error of several learned models to identify similar time series corresponding to musical information.

In the present study we consider the use of auto-regressive (AR) models for time series feature extraction (Piccolo, 1990). Given an AR model of the form

$$x_i = a_0 + \sum_{j=1}^{\eta} a_j x_{i-j}, \quad (3)$$

where  $a_j$  denotes the  $j$ -th regression coefficient and  $\eta$  is the order of the model, we can estimate its coefficients, e.g., by the Yule-Walker function (Marple, 1987). Then, the dissimilarity between two time series can be calculated, for instance, using the Euclidean distance between their estimated coefficients (Piccolo, 1990), analogously as in Eq. 2. The number of AR coefficients is controlled by the parameter  $\eta$  which, similarly to  $\theta$  with FCs, directly affects the final speed of similarity calculations (AR and FCs are usually estimated offline, prior to similarity calculations).

#### 2.4. Dynamic time warping

Dynamic time warping (DTW; Sakoe and Chiba, 1978; Berndt and Clifford, 1994) is a classic approach for computing the dissimilarity between two time series. It has been exploited in countless works: to construct decision trees (Rodríguez and Alonso, 2004), to retrieve similar shapes from large image databases (Bartolini et al., 2005), to match incomplete time series in medical applications (Tormene et al., 2009), to align signatures in an identity authentication task (Kholmatov and Yanikoglu, 2005), etc. In addition, several extensions for speeding up its calculations exist (Keogh and Ratanamahatana, 2005; Salvador and Chan, 2007; Lemire, 2009).

DTW belongs to the group of so-called elastic dissimilarity measures (Wang et al., 2012), and works by optimally aligning (or ‘warping’) the time series in the temporal domain so that the accumulated cost of this alignment is minimal (Fig. 1, center). In its canonical form, this accumulated cost can be obtained by dynamic programming, recursively applying

$$D_{i,j} = f(x_i, y_j) + \min \{D_{i,j-1}, D_{i-1,j}, D_{i-1,j-1}\} \quad (4)$$

for  $i = 1, \dots, M$  and  $j = 1, \dots, N$ , being  $M$  and  $N$  the lengths of time series  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. Except for the first cell, which is initialized to  $D_{0,0} = 0$ , the matrix  $D$  is initialized to  $D_{i,j} = \infty$  for  $i = 0, 1, \dots, M$  and  $j = 0, 1, \dots, N$ . In the case of dealing with uni-dimensional time series, the local



cost function  $f()$ , also called sample dissimilarity function, is usually taken to be the square of the difference between  $x_i$  and  $y_j$  (Berndt and Clifford, 1994), i.e.,  $f(x_i, y_j) = (x_i - y_j)^2$ . In the case of dealing with multidimensional time series or having some domain-specific knowledge, the local cost function  $f()$  must be chosen appropriately, although the Euclidean distance is often used. The final DTW dissimilarity measure typically corresponds to the total accumulated cost, i.e.,  $d_{\text{DTW}}(\mathbf{x}, \mathbf{y}) = D_{M,N}$ . A normalization of  $d_{\text{DTW}}$  can be performed on the basis of the alignment of the two time series, which is found by backtracking from  $D_{M,N}$  to  $D_{0,0}$  (Rabiner and Juang, 1993). However, in preliminary analysis we found the normalized variant to be equivalent, or sensibly less accurate, than the unnormalized one.

The canonical form of DTW presented in Eq. 4 can incorporate many variants. In particular, several constraints can be applied to the computation of  $D$ . A common constraint (Sakoe and Chiba, 1978) is to introduce a window parameter  $\omega \in [0, N]$ , such that the recursive formula of Eq. 4 is only applied for  $i = 1, \dots, M$  and

$$j = \max\{1, i' - \omega\}, \dots, \min\{N, i' + \omega\}, \quad (5)$$

where  $i'$  is progressively adjusted for dealing with different time series lengths, i.e.,  $i' = \lfloor iN/M \rfloor$ , using  $\lfloor \cdot \rfloor$  as the round-to-the-nearest-integer operator. Notice that if  $\omega = 0$  and  $N = M$ ,  $d_{\text{DTW}}$  will correspond to the squared Euclidean distance (the value in  $D_{M,N}$  will be the sum of the squared differences, see Eqs. 1 and 4). Notice furthermore that, when  $\omega = N$ , we are using the unconstrained version of DTW (the constraints in Eq. 5 have no effect). Thus, we include two DTW variants in a single formulation. In general, the introduction of constraints, and specially of the window parameter  $\omega$ , carries some advantages (Keogh and Kasetty, 2003; Rabiner and Juang, 1993; Wang et al., 2012). For instance, constraints prevent ‘pathological alignments’ and, therefore, usually provide better similarity estimates (pathological alignments typically go beyond the main diagonal of  $D$ ). Moreover, constraints allow for reduced computational costs, since only a percentage of the cells in  $D$  needs to be examined (Sakoe and Chiba, 1978; Rabiner and Juang, 1993).

DTW currently stands as the main benchmark against which new similarity measures need to be compared (Xi et al., 2006; Wang et al., 2012). Very few measures have been proposed that systematically outperform DTW for a number of different data sources. These measures are usually more complex than DTW, sometimes requiring extensive tuning of one or more parameters. Additionally, it is often the case that no careful, rigorous, and

extensive evaluation of the accuracy of such measures is done, and further studies fail to assess the statistical significance of their improvement. Thus we could say that the superiority of such measures is, at best, unclear. In this paper, we pay special attention to all these aspects in order to formally assess the considered measures under a common framework. As it will be shown, there exists a similarity measure outperforming DTW for a statistically significant margin (Sec. 4).

## 2.5. Edit distance on real sequences

Turning to previous evidence (Wang et al., 2012), we observe that perhaps the only measure able to seriously challenge DTW is the edit distance on real sequences (EDR; Chen et al., 2005). The EDR corresponds to the extension of the original edit or **Levenshtein distance** (Levenshtein, 1966) to real-valued time series. Such extensions are not commonplace, but recent research is starting to focus on them (Morse and Patel, 2007; Marteau, 2009). As noted by Chen et al. (2005), EDR outperformed previous edit distance variants for time series similarity.

The computation of the EDR can be formalized by a dynamic programming approach. Specifically, we compute

$$D_{i,j} = \begin{cases} D_{i-1,j-1} & \text{if } m(x_i, y_j) = 1 \\ 1 + \min \{D_{i,j-1}, D_{i-1,j}, D_{i-1,j-1}\} & \text{if } m(x_i, y_j) = 0, \end{cases} \quad (6)$$

for  $i = 1, \dots, M$  and  $j = 1, \dots, N$ . The match function used is

$$m(x_i, y_j) = \Theta(\varepsilon - f(x_i, y_j)), \quad (7)$$

where  $\Theta()$  is the Heaviside step function such that  $\Theta(z) = 1$  if  $z \geq 0$  and 0 otherwise, and  $\varepsilon \in [0, \infty)$  is a suitably chosen threshold parameter that controls the degree of resemblance between two time series samples being considered as a match. The first row of  $D$  is initialized to  $D_{i,0} = i$  for  $i = 0, 1, \dots, M$  and the first column of  $D$  to  $D_{0,j} = j$  for  $j = 0, 1, \dots, N$ . Following Chen et al. (2005), who initially reported some accuracy improvements of EDR over DTW, we set the local cost function  $f()$  to the absolute difference between the sample values, i.e.,  $f(x_i, y_j) = |x_i - y_j|$ . This has the additional advantage that we can easily relate  $\varepsilon$  to the standard deviation of the time series (Sec. 3.5).

299 *2.6. Time-warped edit distance*

300 The time-warped edit distance (TWED; Marteau, 2009) is perhaps the  
 301 most interesting extension of dynamic programming algorithms like DTW  
 302 and EDR. In a sense, it is a combination of these two. Like edit dis-  
 303 tances, TWED comprises a mismatch penalty  $\lambda$  and, like dynamic time  
 304 warping, it introduces a so-called stiffness parameter  $\nu$ , controlling its ‘elas-  
 305 ticity’ (Marteau, 2009). For uniformly-sampled time series, the formulation  
 306 of TWED corresponds to

$$D_{i,j} = \min \{D_{i,j} + \Gamma_{\mathbf{xy}}, D_{i-1,j} + \Gamma_{\mathbf{x}}, D_{i,j-1} + \Gamma_{\mathbf{y}}\}, \quad (8)$$

307 for  $i = 1, \dots, M$  and  $j = 1, \dots, N$ , with

$$\begin{aligned} \Gamma_{\mathbf{xy}} &= f(x_i, y_j) + f(x_{i-1}, y_{j-1}) + 2\nu|i - j|, \\ \Gamma_{\mathbf{x}} &= f(x_i, x_{i-1}) + \nu + \lambda, \\ \Gamma_{\mathbf{y}} &= f(y_j, y_{j-1}) + \nu + \lambda, \end{aligned} \quad (9)$$

308 where  $f()$  can be any  $L_n$  metric (Eq. 1). Following Marteau (2009), and as  
 309 done for EDR as well, we choose  $f(x_i, y_j) = |x_i - y_j|$ . Together with DTW  
 310 and EDR, the final dissimilarity value is taken to be  $d_{\text{TWED}}(\mathbf{x}, \mathbf{y}) = D_{M,N}$ .

311 An interesting aspect of TWED is that, in its original formulation (Marteau,  
 312 2009), it takes time stamp differences into account. Therefore, it is able to  
 313 cope with time series of different sampling rates, including down-sampled  
 314 time series. A further interesting aspect, and contrasting to DTW and other  
 315 measures, is that TWED is a metric (Marteau, 2009). Thus, one can exploit  
 316 the triangular inequality to speed up the search in the metric space. Finally,  
 317 it is worth mentioning that the combination of the two previous characteris-  
 318 tics results in a lower bound of the TWED dissimilarity, which can be used  
 319 to speed up nearest neighbor retrieval.

320 *2.7. Minimum jump costs dissimilarity*

321 The main idea behind the minimum jump costs dissimilarity measure (MJC;  
 322 Serra and Arcos, 2012) is that, if a given time series  $\mathbf{x}$  resembles  $\mathbf{y}$ , the cu-  
 323 mulative cost of iteratively ‘jumping’ between their samples should be small<sup>2</sup>  
 324 (Fig. 1, bottom). Supposing that for the  $i$ -th jump we are at time step  $t_x$

---

<sup>2</sup>An implementation of MJC is made available online by the authors: [http://www.iiia.csic.es/~jserra/downloads/2012\\_SerraArcos\\_MJC-Dissim.tar.gz](http://www.iiia.csic.es/~jserra/downloads/2012_SerraArcos_MJC-Dissim.tar.gz) (last accessed on September 15, 2013).

of time series  $\mathbf{x}$ , and that we previously visited time step  $t_y - 1$  of  $\mathbf{y}$ , the minimum jump cost is expressed as

$$c_{\min}^{(i)} = \min \left\{ c_{t_x}^{t_y}, c_{t_x}^{t_y+1}, c_{t_x}^{t_y+2}, \dots \right\}, \quad (10)$$

where  $c_{t_x}^{t_y+\Delta}$  is the cost of jumping from  $x_{t_x}$  to  $y_{t_y+\Delta}$  and  $\Delta = 0, 1, 2, \dots$  is an integer time step increment such that  $t_y + \Delta \leq N$ . After a jump is made,  $t_x$  and  $t_y$  are updated accordingly:  $t_x$  becomes  $t_x + 1$  and  $t_y$  becomes  $t_y + \Delta + 1$ . In case we want to jump from  $\mathbf{y}$  to  $\mathbf{x}$ , only  $t_x$  and  $t_y$  need to be swapped (Serrà and Arcos, 2012).

To define a jump cost  $c_{t_x}^{t_y+\Delta}$ , the temporal and the magnitude dimensions of the time series are considered:

$$c_{t_x}^{t_y+\Delta} = (\phi\Delta)^2 + f(x_{t_x}, y_{t_y+\Delta}), \quad (11)$$

where  $\phi$  represents the cost of advancing in time and  $f()$  is the local cost function, which we take to be  $f(x_{t_x}, y_{t_y+\Delta}) = (x_{t_x} - y_{t_y+\Delta})^2$ , similarly to what is done with DTW (Eq. 4). Notice that, akin to the general formulation of TWED, the term  $(\phi\Delta)^2$  introduces a nonlinear penalty that depends on the temporal gap. Here, the value of  $\phi$  is set proportional to the standard deviation  $\sigma$  expected for the time series and, at the same time, proportional to the real-valued parameter  $\beta \in [0, \infty)$ , which controls how difficult is to advance in time (for more details see Serrà and Arcos, 2012). To obtain a symmetric dissimilarity measure,  $d_{\text{MJC}}(\mathbf{x}, \mathbf{y}) = \min \{d_{\text{XY}}, d_{\text{YX}}\}$  can be used, where  $d_{\text{XY}}$  and  $d_{\text{YX}}$  are the cumulative MJCs obtained by starting at  $x_1$  and  $y_1$ , respectively.

### 3. Evaluation methodology

#### 3.1. Classification scheme

The efficacy of a time series similarity measure is commonly evaluated by the classification accuracy it achieves (Keogh and Kasetty, 2003; Wang et al., 2012). For that, the error ratio of a distance-based classifier is calculated for a given labeled data set, understanding the error ratio as the number of wrongly classified items divided by the total number of tested items. The standard choice for the classifier is the one-nearest neighbor (1NN) classifier. Following Wang et al. (2012), we can enumerate several advantages of using this approach. First, the error of the 1NN classifier critically depends on the similarity measure used. Second, the 1NN classifier is parameter-free and easy to implement. Third, there are theoretical results relating the error

357 of an 1NN classifier to errors obtained with other classification schemes.  
358 Fourth, some works suggest that the best results for time series classification  
359 come from simple nearest neighbor methods. For more details on these  
360 aspects we refer to Mitchell (1997); Hastie et al. (2009), and the references  
361 provided by Wang et al. (2012).

### 362 3.2. Data sets

363 We perform experiments with 45 publicly-available time series data sets  
364 from the UCR time series repository (Keogh et al., 2011). This is the world’s  
365 biggest time series repository, and some authors estimate that it makes up to  
366 more than 90% of all publicly-available, labeled data sets (Wang et al., 2012).  
367 The repository comprises synthetic, as well as real-world data sets, and  
368 also includes one-dimensional time series extracted from two-dimensional  
369 shapes (Keogh et al., 2011). The 45 data sets considered here correspond  
370 to the totality of the UCR repository, as by March 2013. Within such data  
371 sets, the number of classes ranges from 2 to 50, the number of time series  
372 per data set ranges from 56 to 9,236, and time series lengths go from 24  
373 to 1,882 samples. For further details on these data sets we refer to (Keogh  
374 et al., 2011).

### 375 3.3. Cross-validation

376 To properly assess a classifier’s error, out-of-sample validation needs to  
377 be done (Salzberg, 1997). In our experiments, we follow a standard 3-  
378 fold cross-validation scheme using balanced data sets (Mitchell, 1997; Hastie  
379 et al., 2009), i.e., using the same number of items per class. We repeat the  
380 validation 20 times and report average error ratios. Balancing the data sets  
381 allows for balanced error estimations regarding the class distribution, and  
382 repeating cross-fold validation several times allows for more precise estima-  
383 tions (Mitchell, 1997; Hastie et al., 2009). The use of a cross-fold validation  
384 scheme is essential for avoiding the bias that a particular split of the data  
385 could introduce (Salzberg, 1997; Hastie et al., 2009).

386 We also computed error ratios for the original splits provided in the  
387 UCR time series repository (Keogh et al., 2011). This allowed us to confirm  
388 that the 1NN error ratios from our implementations of DTW and Euclidean  
389 distance agree with the values reported there. In addition, we observed that  
390 the error ratios obtained by such splits were substantially different from the  
391 ones obtained by cross-validation, up to the point of even modifying the  
392 ranking of some algorithms with respect to those error ratios in some data  
393 sets. This indicates a potential bias in such individual splits, an aspect that  
394 is well-known in the machine learning community (Salzberg, 1997; Mitchell,

1997; Hastie et al., 2009). We refer the interested reader to any machine learning textbook for a more in-depth discussion of cross-fold validation schemes and their appropriateness over individual splits. Besides, using a single split difficults statistical significance assessment (see below). A full account of the raw error ratios for all measures and data sets is available online<sup>3</sup>, including the error ratios for the aforementioned original splits.

### 3.4. Statistical significance

To assess the statistical significance of the difference between two error ratios we employ the well-known Wilcoxon signed-rank test (Hollander and Wolfe, 1999). The Wilcoxon signed-rank test is a non-parametric statistical hypothesis test used when comparing two repeated measurements (or related samples, or matched samples) in order to assess whether their population mean ranks differ. It is the natural alternative to the Student’s  $t$ -test for dependent samples when the population distribution cannot be assumed to be normal (Hollander and Wolfe, 1999). For a given data set, we use as input the  $20 \times 3$  accuracy values obtained for each classifier (i.e., the test fold accuracies). Besides, for comparing similarity measures on a more global basis using all data sets, we employ as input the 45 average accuracy values obtained for each data set. Following common practice (Salzberg, 1997; Hollander and Wolfe, 1999), the threshold significance level is set to 5%. Additionally, to compensate for multiple pairwise comparisons, we apply the Holm-Bonferroni method (Holm, 1979), a post-hoc statistical analysis method controlling the so-called family-wise error rate that is more powerful than the usual Bonferroni correction (Demšar, 2006).

### 3.5. Parameter choices

Before performing the experiments, all time series from all data sets were z-normalized so that each individual time series had zero mean and unit variance. Furthermore, we optimized the measures’ parameters in the training phase of our cross-validation. This optimization step consisted of a grid search within a suitable range of parameter values, forcing the same number of parameter combinations per algorithm (Table 1). The values of the grid are chosen according to common practice and the specifications given in the papers introducing each measure (Sec. 2). Specifically, for FC we used 25 linearly-spaced integer values of  $\theta \in [2, N/2]$ . For AR we

---

<sup>3</sup>[http://www.iiia.csic.es/~jserra/downloads/2013\\_SerraArcos\\_AnEmpiricalEvaluation.tar.gz](http://www.iiia.csic.es/~jserra/downloads/2013_SerraArcos_AnEmpiricalEvaluation.tar.gz) (last accessed on September 15, 2013).

Measure	Parameter	Minimum value	Maximum value	Number of steps	Extra value
FC	$\theta$	2	$0.5N$	25	-
AR	$\eta$	1	$0.25N$	25	-
DTW	$\omega$	0	$0.25N$	24	$N$
EDR	$\varepsilon$	$0.02\sigma$	$\sigma$	25	-
TWED	$\nu$	$10^{-5}$	1	5	-
TWED	$\lambda$	0	1	5	-
MJC	$\beta$	0	25	24	$10^{10}$

Table 1: Parameter grid for the considered similarity measures (recall that  $N$  corresponds to the length of the time series and, since we z-normalize all time series,  $\sigma = 1$ ). For DTW and MJC we consider an extra value corresponding to unconstrained DTW and to the Euclidean configuration of MJC, respectively. All parameter values were linearly spaced except  $\nu$ , which was logarithmically spaced.

used 25 linearly-spaced integer values of  $\eta \in [1, 0.25N]$  (because of the z-normalization, we remove  $a_0$  in Eq. 3). For DTW we used 24 linearly-spaced integer values of  $\omega \in [0, 0.25N]$  plus  $w = N$  (the unconstrained DTW variant). For EDR we used 25 linearly-spaced real values of  $\varepsilon \in [0.02\sigma, \sigma]$ ,  $\sigma$  being the standard deviation of the time series (because of the z-normalization  $\sigma = 1$ ). For TWED we used all possible 25 combinations for  $\nu = [10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1]$  and  $\lambda = [0, 0.25, 0.5, 0.75, 1]$ . For MJC we used 24 linearly-spaced real values of  $\beta \in [0, 25]$  plus  $\beta = 10^{10}$  (in practice corresponding to the squared Euclidean distance variant, Eq. 11). After the grid search, the parameter value yielding to the lowest leave-one-out error ratio for the training set was kept for out-of-sample testing.

## 4. Results

### 4.1. Classification performance: test

If we look at the overall results, we see that all considered measures clearly outperform the random baseline for practically all the 45 data sets (Table 2). Furthermore, we see that some of them achieve near-perfect accuracies for a number of data sets (e.g., *CBF*, *CinC\_ECG\_torso*, *ECGFiveDays*, *Two\_Patterns*, or *TwoLeadECG*). However, no single measure achieves the best performance for all the data sets. The Euclidean distance is found to be the best-performing measure in 2 data sets, FC is the best-performing in 4 data sets, AR in 1, DTW in 6, EDR in 7, TWED in 20, and MJC in 5. If we count only the data sets where one measure statistically significantly outperforms the rest, the numbers reduce to 0 for Euclidean, 2 for FC, 1 for

AR, 2 for DTW, 2 for EDR, 6 for TWED, and 0 for MJC. Thus, interestingly, there are some data sets where choosing a specific similarity measure can make a difference.

#	Data set	Random	Euc	FC	AR	DTW	EDR	TWED	MJC
1	<i>50words</i>	0.969	0.503	0.685	0.867	0.332	0.289	<b>0.237*</b>	0.319
2	<i>Adiac</i>	0.970	0.345	<b>0.266*</b>	0.725	0.355	0.423	0.335	0.346
3	<i>Beef</i>	0.763	0.417	<b>0.390</b>	0.504	0.472	0.439	0.506	0.448
4	<i>CBF</i>	0.655	0.013	0.358	0.432	0.000	0.002	<b>0.000</b>	0.001
5	<i>ChlorineConcentration</i>	0.673	0.071	0.063	<b>0.038*</b>	0.072	0.094	0.093	0.070
6	<i>CinC_ECG_torso</i>	0.749	0.002	0.008	0.102	0.001	<b>0.000*</b>	0.001	0.002
7	<i>Coffee</i>	0.394	0.019	0.024	0.139	<b>0.014</b>	0.031	0.021	0.023
8	<i>Cricket_X</i>	0.913	0.378	0.348	0.713	0.209	0.237	<b>0.190*</b>	0.253
9	<i>Cricket_Y</i>	0.928	0.423	0.411	0.814	0.222	0.224	<b>0.209</b>	0.267
10	<i>Cricket_Z</i>	0.920	0.380	0.353	0.731	0.212	0.235	<b>0.194*</b>	0.254
11	<i>DiatomSizeReduction</i>	0.744	0.008	0.011	0.222	0.010	0.016	0.012	<b>0.007</b>
12	<i>ECG200</i>	0.515	0.130	0.145	0.227	0.139	0.148	<b>0.109</b>	0.130
13	<i>ECGFiveDays</i>	0.505	0.007	<b>0.000</b>	0.072	0.003	0.003	0.005	0.001
14	<i>FaceAll</i>	0.931	0.139	0.152	0.649	0.053	0.019	<b>0.019</b>	0.034
15	<i>FaceFour</i>	0.679	0.111	0.149	0.545	0.069	0.028	0.025	<b>0.024</b>
16	<i>FacesUCR</i>	0.929	0.138	0.148	0.648	0.052	0.019	<b>0.018</b>	0.041
17	<i>Fish</i>	0.871	0.183	0.234	0.617	0.184	<b>0.084</b>	0.094	0.114
18	<i>Gun_Point</i>	0.506	0.058	0.031	0.149	0.023	<b>0.010</b>	0.017	0.014
19	<i>Haptics</i>	0.793	0.604	0.610	0.678	0.554	0.611	<b>0.544</b>	0.563
20	<i>InlineSkate</i>	0.862	0.524	0.601	0.497	0.462	0.456	0.416	<b>0.411</b>
21	<i>ItalyPowerDemand</i>	0.489	0.035	0.083	0.261	<b>0.033</b>	0.042	0.036	0.034
22	<i>Lighting2</i>	0.488	0.297	0.281	0.450	0.162	0.220	<b>0.161</b>	0.254
23	<i>Lighting7</i>	0.817	0.371	0.463	0.707	<b>0.252</b>	0.362	0.256	0.336
24	<i>Mallat</i>	0.870	0.018	0.020	0.058	0.015	0.006	<b>0.006</b>	0.014
25	<i>MedicalImages</i>	0.912	0.313	0.455	0.458	0.247	0.330	<b>0.228</b>	0.305
26	<i>MoteStrain</i>	0.513	0.087	0.162	0.336	0.058	0.024	<b>0.021</b>	0.034
27	<i>NonInvasiveFetalECG1</i>	0.978	0.171	0.213	0.401	0.175	0.186	0.182	<b>0.169</b>
28	<i>NonInvasiveFetalECG2</i>	0.975	<b>0.106</b>	0.146	0.296	0.107	0.118	0.108	0.110
29	<i>OliveOil</i>	0.644	<b>0.104</b>	0.185	0.663	0.154	0.194	0.146	0.127
30	<i>OSULeaf</i>	0.832	0.409	0.306	0.617	0.359	<b>0.191*</b>	0.232	0.256
31	<i>SonyAIBORobotSurface</i>	0.510	0.017	0.040	0.079	0.018	0.026	0.017	<b>0.015</b>
32	<i>SonyAIBORobotSurfaceII</i>	0.489	0.018	0.032	0.113	0.021	0.023	<b>0.016</b>	0.019
33	<i>StarLightCurves</i>	0.671	0.124	<b>0.070*</b>	0.274	0.083	0.107	0.097	0.109
34	<i>SwedishLeaf</i>	0.932	0.196	0.142	0.376	0.129	0.101	<b>0.094</b>	0.100
35	<i>Symbols</i>	0.838	0.038	0.074	0.260	0.019	<b>0.015</b>	0.016	0.018
36	<i>Synthetic_control</i>	0.834	0.087	0.393	0.511	<b>0.009*</b>	0.047	0.014	0.034
37	<i>Trace</i>	0.757	0.169	0.117	0.117	<b>0.000*</b>	0.034	0.011	0.038
38	<i>Two_Patterns</i>	0.743	0.020	0.491	0.724	0.000	0.000	<b>0.000</b>	0.001
39	<i>TwoLeadECG</i>	0.507	0.006	0.012	0.202	0.001	0.002	<b>0.001</b>	0.003
40	<i>UWaveGestureLibrary_X</i>	0.872	0.234	0.566	0.694	0.199	0.214	<b>0.192*</b>	0.203
41	<i>UWaveGestureLibrary_Y</i>	0.876	0.288	0.631	0.645	<b>0.263</b>	0.280	0.265	0.267
42	<i>UWaveGestureLibrary_Z</i>	0.879	0.298	0.546	0.678	0.265	0.271	<b>0.250*</b>	0.261
43	<i>Wafer</i>	0.497	0.004	0.003	0.013	0.005	<b>0.002</b>	0.003	0.005
44	<i>WordsSynonyms</i>	0.960	0.496	0.675	0.855	0.327	0.304	<b>0.251*</b>	0.310
45	<i>Yoga</i>	0.500	0.070	0.108	0.333	0.061	<b>0.034</b>	0.037	0.047
Average rank		7.99	4.40	5.07	6.80	3.00	3.42	<b>2.29</b>	3.04

Table 2: Error ratios for all considered measures and data sets. The symbol \* denotes a statistically significant difference with respect to the other measures for a given data set ( $p < 0.05$ , Sec. 3.4). The last row contains the average rank of each measure across all data sets (i.e., the average position after sorting the errors for a given data set in ascending order).



455 Beyond accuracies, this latter aspect can potentially highlight inherent  
 456 data set qualities. For instance, the fact that a feature/model-based measure  
 457 clearly outperforms the others for a particular data set indicates that such  
 458 time series may be very well characterized by the extracted features/fitted  
 459 model (e.g., FC with *Adiac* for features and AR with *ChlorineConcentra-*  
 460 *tion* for models). In addition, the good or bad performance of Euclidean  
 461 and elastic measures gives us an intuition of the importance of alignments,  
 462 warping, or sample correspondences (e.g., these may be very important for  
 463 *Trace* and the three *Face* data sets, where there is an order of magnitude  
 464 difference between Euclidean and warping-based measures, but not much  
 465 for *DiatomSizeReduction* or *NonInvasiveFetalECG2*, where Euclidean gets  
 466 numbers that are very close, or even better than the ones obtained by the  
 467 warping-based measures).

468 In general, we see that TWED outperforms the other measures in several  
 469 data sets, with an average rank of 2.29 (Table 2). In fact, if we compare  
 470 the considered measures on a more global scale, taking the matched error  
 471 ratios across data sets (Sec. 3.4), we obtain that TWED is statistically  
 472 significantly superior to the rest (Fig. 2). Next, we see that DTW, MJC, and  
 473 EDR form a group of equivalent measures, with no statistically significant  
 474 difference between them. The performed statistical analysis also separates  
 475 the remaining measures from these and also between themselves. Apart  
 476 from this more global analysis, further pairwise comparisons can be made,  
 477 confirming the aforementioned global tendencies (Fig. 3).

#### 478 4.2. Classification performance: test vs. train

479 For choosing the most optimal parameters for a given measure and  
 480 data set we solely dispose of the training data. Hence, it is important  
 481 to know whether the error ratios for training and testing sets are similar,  
 482 otherwise one could be incurring into the so-called “Texas sharpshooter fal-  
 483 lacy” (Batista et al., 2011), i.e., one could not predict a measure’s utility  
 484 ahead of time by just looking at training data. For comparing train and test  
 485 error ratios, we can compute an error gain value for a couple of measures on  
 486 each data set and check whether such values for train and test agree. To do  
 487 so, a kind of real-valued contingency table can be plotted, called the “Texas  
 488 sharpshooter plot” by Batista et al. (2011). Due to space reasons, we here  
 489 only show such contingency tables for TWED against DTW and Euclidean  
 490 distance (Fig. 4). The results show that error gains between TWED and  
 491 DTW/Euclidean mostly agree between training and testing. As mentioned  
 492 in Sec. 3.3, a full, raw account of train and test errors is available online.  
 493 Having a close look at those full results, we can see that, in general, the

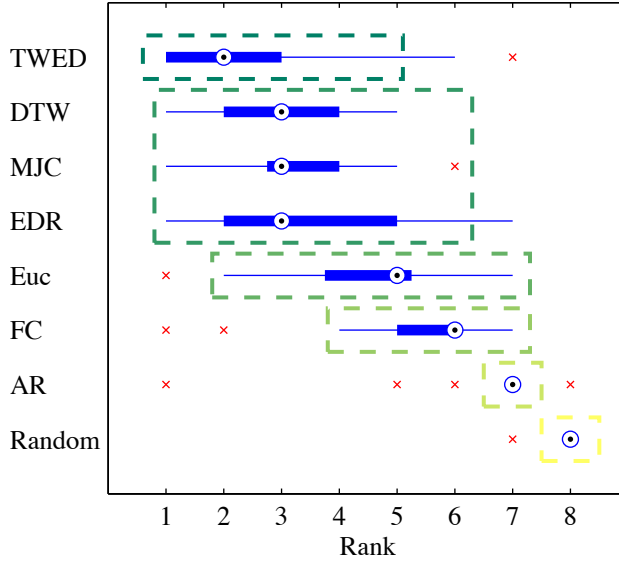


Figure 2: Box plot for the distribution of performance ranks of each measure across data sets. The dashed lines denote statistically significantly equivalent groups of measures ( $p < 0.05$ , Sec. 3.4).

best-performing measure at the training stage is also the best-performing measure at the testing stage. The few exceptions can be easily listed (Table 3). The relative rankings for the measures that do not perform best also mostly agree between train and test.

#### 4.3. Parameter assessment

We finally report on the parameters chosen for each measure after training with 66% of balanced data (Fig. 5). Firstly, we observe that, in the vast majority of cases, a specific value for a given parameter is consistently chosen across the  $20 \times 3$  performed iterations (we see clear peaks in the distributions of Fig. 5). Among these consistent choices, perhaps TWED and MJC present the most spread distributions. Such aspect, together with the fairly good accuracies obtained for these two specific measures (Sec. 4.1), indicates a certain degree of robustness against specific parameter choices. This is a very desirable quality of a time series similarity measure, even more if we have to train a classifier with a potentially incomplete set of training instances.

Next, we see that the selected parameters are generally not in the borders

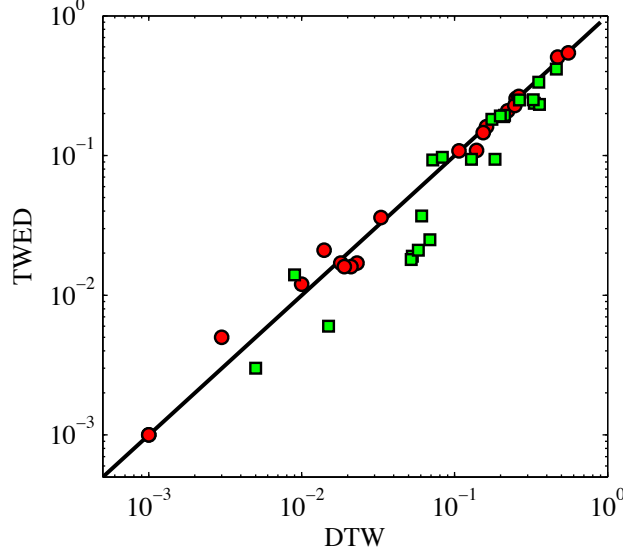


Figure 3: Error ratios comparison between DTW and TWED (notice the logarithmic axes). The lower-right triangular part corresponds to TWED outperforming DTW, whereas the upper-left part corresponds to the opposite case. The green squares indicate statistically significant performance differences ( $p < 0.05$ , Sec. 3.4).

of the specified ranges, thus indicating that a reasonable choice has been made (Fig 5). This is particularly true for DTW and EDR. The only measure that could potentially benefit from reconsidering the parameters' range is TWED. As it can be seen,  $\nu$  and  $\lambda$  seem to be consistently chosen in the lower and upper parts of the specified ranges, respectively. This suggests that the best combination for some data sets could lie outside the parameter space outlined by Marteau (2009), i.e., in  $0 < \nu < 10^{-4}$  and/or  $\lambda > 1$ . If that was the case, TWED could potentially achieve even much higher accuracies. Interestingly, TWED is not the best-performing measure for some of the data sets where 'border' parameter values are chosen (e.g., *CBF*, *Fish*, *StarLightCurves*, *TwoPatterns*).

Finally, we can comment on the particularities of some data sets with relation to classification. For instance, we see that a relatively large window parameter  $\omega$  (DTW) is chosen for data sets 36 to 39 (i.e., *Synthetic\_control*, *Trace*, *Two\_Patterns*, and *TwoLeadECG*). This denotes that tracking alignments or warping paths beyond the main diagonal of  $D$  (Eq. 4) might be advantageous for classification in these data sets. Interestingly, the stiffness

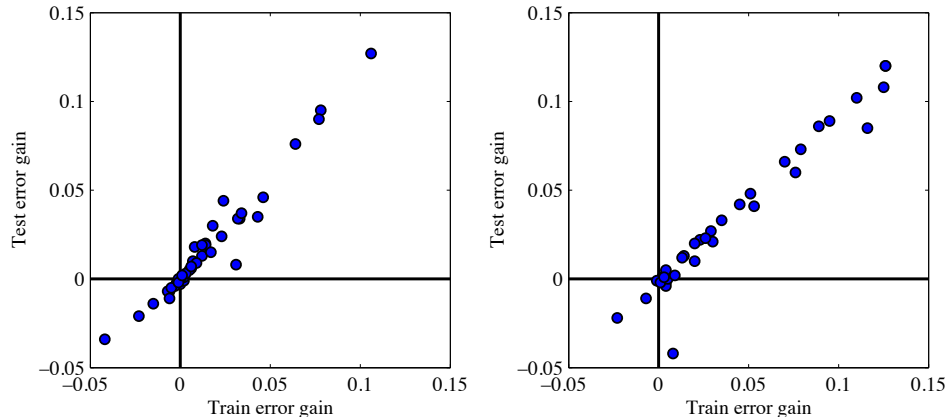


Figure 4: Texas sharpshooter plots for TWED against DTW (left) and Euclidean distance (right). Here, error gain is measured by subtracting the TWED error ratio from the one of DTW/Euclidean. Dots around the diagonal indicate agreement of error gain for train and test. False positives, i.e., dots in the lower-right quadrant, indicate that TWED, being the best measure after training, does not reach the lowest error at testing. For instance, in the case of TWED vs. Euclidean (right), the *OliveOil* data set false positive stands out at coordinates  $(0.008, -0.042)$  (see also Table 3). For further details on the construction of Texas sharpshooter plots we refer to Batista et al. (2011).

parameter  $\nu$  (TWED), which accounts for a similar but opposite concept (Sec. 2.6), takes relatively small values. Such agreement across different measures reinforces the hypothesis that tracking intricate alignments or strongly warped paths may be advantageous for these data sets. Analogous and complementary conclusions can be derived for other data sets. For instance, in data sets 11 (*DiatomSizeReduction*) and 13 (*ECGFiveDays*), a small number of both FCs  $\theta$  and AR coefficients  $\eta$  is chosen. As FC and AR achieve competitive accuracies in those specific data sets, we could suspect that low-frequency components are important for correctly classifying the instances in those data sets (Secs. 2.2 and 2.3).

## 5. Conclusion

From a general perspective, the obtained results show that there is a group of equivalent similarity measures, with no statistically significant differences among them (DTW, EDR, and MJC). The existing literature suggests that some longest common sub-sequence approaches (Gusfield, 1997), together with alternative variants of DTW and EDR (e.g., Sakoe and Chiba,

#	Data set	Measure	Outperf. by	Gain
3	<i>Beef</i>	FC	EDR	0.049
4	<i>CBF</i>	TWED	DTW	<0.001
7	<i>Coffee</i>	DTW	FC	0.004
12	<i>ECG200</i>	TWED	MJC	0.002
15	<i>FaceFour</i>	MJC	EDR	0.007
18	<i>Gun_Point</i>	EDR	MJC	0.004
19	<i>Haptics</i>	TWED	MJC	0.009
21	<i>ItalyPowerDemand</i>	DTW	EDR	0.001
28	<i>NonInvasiveFetalECG2</i>	Euclidean	TWED	0.001
29	<i>OliveOil</i>	Euclidean	TWED	0.008
39	<i>TwoLeadECG</i>	TWED	DTW	<0.001

Table 3: List of best-performing measures in testing (the column “Measure”) but actually outperformed by others in training (the column “Outperf. by”). The column “Gain” corresponds to the absolute value of the train error gain, i.e., the absolute difference between error ratios at training stage (see also Fig. 4).

1978; Morse and Patel, 2007), could potentially join this group (Marteau, 2009; Wang et al., 2012). However, according to the results reported here, the TWED measure originally proposed by Marteau (2009) seems to consistently outperform all the considered distances, including DTW, EDR, and MJC. Thus, we believe this often unconsidered measure should take a baseline role in future evaluations of time series similarity measures (beyond accuracy, additional properties enumerated in Sec. 2.6 make it also very attractive). The Euclidean distance, although somehow competitive, generally performs statistically significantly worse than TWED, DTW, MJC, and EDR. Its accuracy on large data sets was also not very impressive. Below Euclidean distance, but statistically significantly above the random baseline, we find FC and AR measures. Of course, the general statements above do not exclude the possibility that a particular measure or variant could be very well-suited for a specific data set and statistically significantly outperform the rest (cf. Keogh and Kasetty, 2003). In Sec. 4.1 we have enumerated several examples of that.

When comparing train and test errors, we have seen that these mostly agree, with train errors generally providing a good guess of the test errors on unseen data. We have listed some notable exceptions to this rule and used Texas sharpshooter plots to further assess this aspect for TWED vs. DTW and Euclidean. When assessing the best parameter choices for each measure,

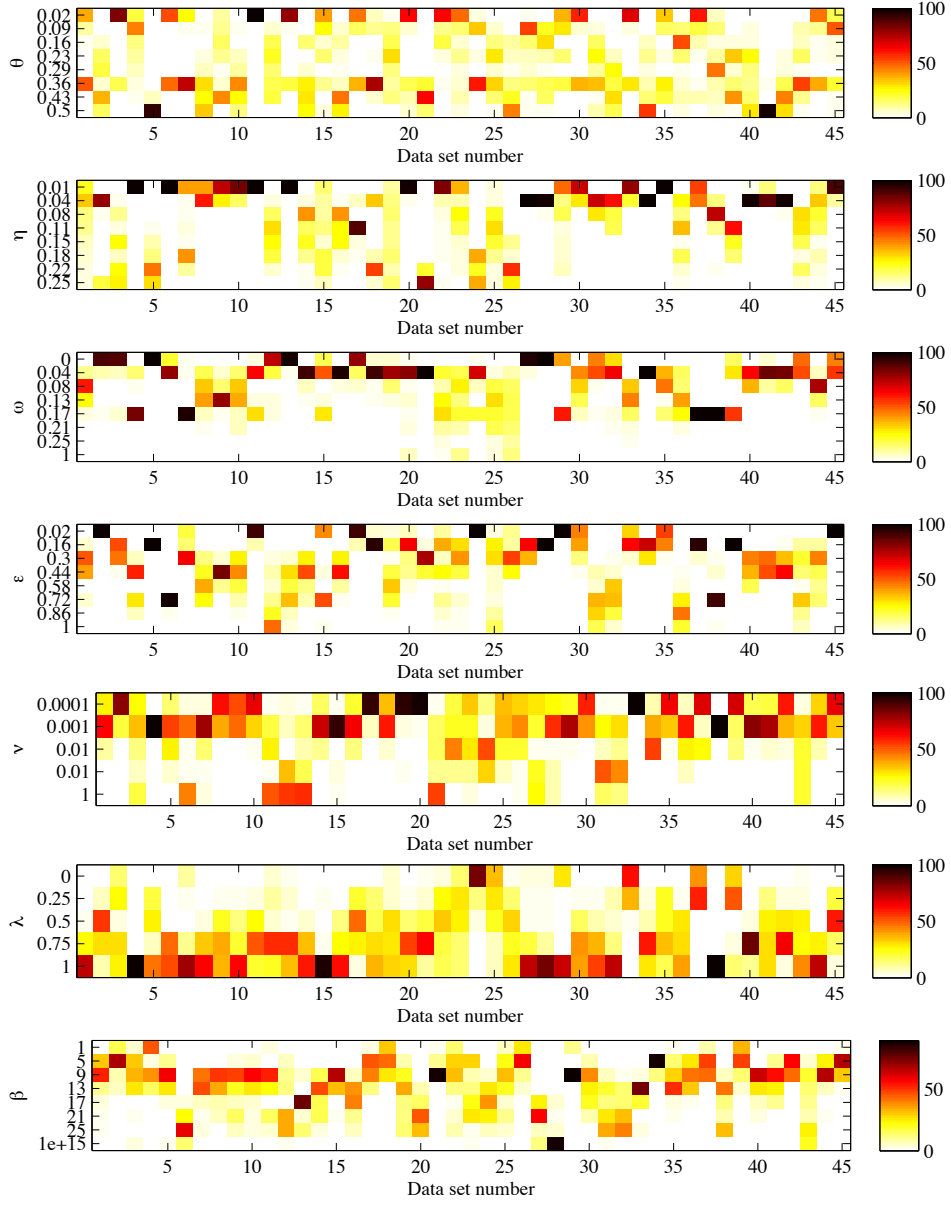


Figure 5: Percentage of times (color code) that a given parameter value (vertical axis) is chosen for each data set (horizontal axis; for the names behind each number see Table 2). From top to bottom, the plots correspond to FC ( $\theta$ ), AR ( $\eta$ ), DTW ( $\omega$ ), EDR ( $\varepsilon$ ), TWED ( $\nu$ ), TWED ( $\lambda$ ), and MJC ( $\beta$ ).

565 we have seen that the considered ranges are typically suitable for the task  
566 at hand. We have also discussed some particularities regarding parameter  
567 choices and the nature of a few data sets.

568 The similarity measure is a crucial step in computational approaches  
569 dealing with time series. However, there are some additional issues worth  
570 mentioning, in particular with regard to post-processing steps focused on  
571 improving similarity assessments (pre-processing steps are sufficiently well-  
572 discussed in the existing literature (see, e.g., Keogh and Kasetty, 2003; Han  
573 and Kamber, 2005; Wang et al., 2012, and references therein). A very in-  
574 teresting post-processing step is the complexity-invariant correction factor  
575 introduced by Batista et al. (2011). Such correction factor prevents from  
576 assigning low dissimilarity values to time series of different complexity, thus  
577 preventing the inclusion of time series of different nature in the same clus-  
578 ter. The way to assess complexity depends on the situation, but Batista  
579 et al. (2011) introduce a quite straightforward way: the  $L_2$  norm of the  
580 sample-based derivative of a time series. Overall, considering different types  
581 of ‘invariance’ is a sensible approach (Batista et al., 2011, provide a good  
582 overview). Here, we have already implicitly considered a number of them,  
583 although more as a pre-processing or method-specific strategy: global ampli-  
584 tude and scale invariance (z-normalization), warping invariance (any elastic  
585 measure, in our case DTW, EDR, TWED, and MJC), phase invariance  
586 ( $AR^4$ ), and occlusion invariance (EDR and TWED).

587 Another interesting post-processing step is the hubness correction for  
588 time series classification introduced by Radovanović et al. (2010). Based on  
589 the finding that some instances in high-dimensional spaces tend to become  
590 hubs by being unexpectedly (and usually wrongly) considered nearest neigh-  
591 bors of several other instances, a correction factor can be introduced. This  
592 usually does not harm classification accuracy and can definitely improve per-  
593 formance for some data sets (Radovanović et al., 2010). A further strategy  
594 for enhancing time series similarity and potentially reducing hubness is the  
595 use of unsupervised clustering algorithms to prune nearest neighbor candi-  
596 dates (Serrà et al., 2012b). Future work should focus on the real quantitative  
597 impact of strategies for enhancing time series similarity like the ones above,  
598 with a special emphasis on its impact to different measures and classification  
599 schemes.

600 The empirical comparison of multiple approaches across a large-scale case  
601 basis is an important and necessary step towards any mature research field.

---

<sup>4</sup>For FC we use both phase and magnitude (Sec. 2.2).

Besides getting a more global picture and highlighting relevant approaches, it pushes towards unified validation procedures and analysis tools. It is hoped that this article will serve as a steppingstone for those interested in advancing in time series similarity, clustering, and classification.

## Acknowledgements

We thank the people who made available or contributed to the UCR time series repository. This research has been funded by 2009-SGR-1434 from Generalitat de Catalunya, JAEDOC069/2010 from Consejo Superior de Investigaciones Científicas, and TIN2009-13692-C03-01 and TIN2012-38450-C03-03 from the Spanish Government.

## References

- Agrawal, R., Faloutsos, C., Swami, A., 1993. Efficient similarity search in sequence databases, in: Proc. of the Int. Conf. on Foundations of Data Organization and Algorithms, pp. 69–84.
- Bartolini, I., Ciaccia, P., Patella, M., 2005. WARP: accurate retrieval of shapes using phase of Fourier descriptors and time warping distance. IEEE Trans. on Pattern Analysis and Machine Intelligence 27, 142–147.
- Batista, G.E.A.P.A., Wang, X., Keogh, E.J., 2011. A complexity-invariant distance measure for time series, in: Proc. of the SIAM Int. Conf. on Data Mining, pp. 699–710.
- Berndt, D.J., Clifford, J., 1994. Using dynamic time warping to find patterns in time series, in: Proc. of the AAAI Workshop on Knowledge Discovery in Databases, pp. 359–370.
- Cai, Y., Ng, R., 2004. Indexing spatio-temporal trajectories with Chebyshev polynomials, in: Proc. of the ACM SIGMOD Int. Conf. on Management of Data, pp. 599–610.
- Chan, K.P., Fu, A.C., 1999. Efficient time series matching by wavelets, in: Proc. of the IEEE Int. Conf. on Data Engineering, pp. 126–133.
- Chen, L., Öszu, M.T., Oria, V., 2005. Robust and fast similarity search for moving object trajectories, in: Proc. of the ACM SIGMOD Int. Conf. on Management of Data, pp. 491–502.



- 633 Das, G., Lin, K.I., Mannila, H., Renganathan, G., Smyth, P., 1998. Rule  
634 discovery from time series, in: Proc. of the AAAI Int. Conf. on Knowledge  
635 Discovery and Data Mining, pp. 16–22.
- 636 Demšar, J., 2006. Statistical comparison of classifiers over multiple data  
637 sets. *Journal of Machine Learning Research* 7, 1–30.
- 638 Faloutsos, C., Ranganathan, M., Manolopoulos, Y., 1994. Fast subsequence  
639 matching in time-series databases, in: Proc. of the ACM SIGMOD Int.  
640 Conf. on Management of Data, pp. 419–429.
- 641 Fu, T.C., 2011. A review on time series data mining. *Engineering Applica-*  
642 *tions of Artificial Intelligence* 24, 164–181.
- 643 Geurts, P., 2002. Contributions to decision tree induction: bias/variance  
644 tradeoff and time series classification. Ph.D. thesis. University of Liège,  
645 Liège, Belgium.
- 646 Gusfield, D., 1997. Algorithms on strings, trees, and sequences: computer  
647 science and computational biology. Cambridge University Press, Cam-  
648 bridge, UK.
- 649 Han, J., Kamber, M., 2005. Data mining: concepts and techniques. Morgan  
650 Kaufmann, Waltham, USA.
- 651 Hastie, T., Tibshirani, R., Friedman, J., 2009. The elements of statistical  
652 learning. 2nd ed., Springer, Berlin, Germany.
- 653 Hollander, M., Wolfe, D.A., 1999. Nonparametric statistical methods. 2nd  
654 ed., Wiley, New York, USA.
- 655 Holm, S., 1979. A simple sequentially rejective multiple test procedure.  
656 *Scandinavian Journal of Statistics* 6, 65–70.
- 657 Kantz, H., Schreiber, T., 2004. Nonlinear time series analysis. Cambridge  
658 University Press, Cambridge, UK.
- 659 Keogh, E.J., 2011. Machine learning in time series databases (and everything  
660 is a time series!). Tutorial at the AAAI Int. Conf. on Artificial Intelligence.
- 661 Keogh, E.J., Chakrabarti, K., Pazzani, M., Mehrotra, S., 2001. Dimension-  
662 ality reduction for fast similarity search in large time series databases.  
663 *Knowledge and Information Systems* 3, 263–286.

- 664 Keogh, E.J., Kasetty, S., 2003. On the need for time series data mining  
665 benchmarks: a survey and empirical demonstration. *Data Mining and*  
666 *Knowledge Discovery* 7, 349–371.
- 667 Keogh, E.J., Ratanamahatana, C.A., 2005. Exact indexing of dynamic time  
668 warping. *Knowledge and Information Systems* 7, 358–386.
- 669 Keogh, E.J., Xi, X., Wei, L., Ratanamahatana, C.A., 2009. Supporting  
670 exact indexing of arbitrarily rotated shapes and periodic time series under  
671 Euclidean and warping distance measures. *VLDB Journal* 11, 611–630.
- 672 Keogh, E.J., Zhu, Q., Hu, B., Hao, Y., Xi, X., Wei, L., Ratanamahatana,  
673 C.A., 2011. The UCR time series classification/clustering homepage.  
674 URL: [http://www.cs.ucr.edu/%7eeamonn/time\\_series\\_data](http://www.cs.ucr.edu/%7eeamonn/time_series_data).
- 675 Kholmatov, A., Yanikoglu, B., 2005. Identity authentication using improved  
676 online signature verification method. *Pattern Recognition Letters* 26,  
677 2400–2408.
- 678 Lemire, D., 2009. Faster retrieval with a two-pass dynamic time warping  
679 lower bound. *Pattern Recognition* 42, 2169–2180.
- 680 Levenshtein, V.I., 1966. Binary codes capable of correcting deletions, inser-  
681 tions, and reversals. *Soviet Physics Doklady* 10, 707–710.
- 682 Liao, T.W., 2005. Clustering of time series data: a survey. *Pattern Recog-*  
683 *nition* 38, 1857–1874.
- 684 Maharaj, E.A., 2000. Clusters of time series. *Journal of Classification* 17,  
685 297–314.
- 686 Marple, S.L., 1987. *Digital spectral estimation*. Prentice-Hall, Englewood  
687 Cliffs, USA.
- 688 Marteau, R.F., 2009. Time warp edit distance with stiffness adjustment  
689 for time series matching. *IEEE Trans. on Pattern Analysis and Machine*  
690 *Intelligence* 31, 306–318.
- 691 Mitchell, T.M., 1997. *Machine Learning*. McGraw-Hill, New York, USA.
- 692 Montani, S., Portinale, L., Leonardi, G., Bellazzi, R., 2006. Case-based  
693 retrieval to support the treatment of end stage renal failure patients. *Ar-*  
694 *tificial Intelligence in Medicine* 37, 31–42.

- 695 Morse, M.D., Patel, J.M., 2007. An efficient and accurate method for eval-  
696 uating time series similarity, in: Proc. of the ACM SIGMOD Int. Conf.  
697 on Management of Data, pp. 569–580.
- 698 Olsson, E., Funk, P., Xiong, N., 2004. Fault diagnosis in industry using  
699 sensor readings and case-based reasoning. *Journal of Intelligent Fuzzy*  
700 *Systems* 15, 41–46.
- 701 Oppenheim, A.V., Schafer, R.W., Buck, J.R., 1999. Discrete-time signal  
702 processing. 2nd ed., Prentice-Hall, Upper Saddle River, USA.
- 703 Piccolo, D., 1990. A distance measure for classifying ARMA models. *Journal*  
704 *of Time Series Analysis* 11, 153–163.
- 705 Povinelli, R.J., Johnson, M.T., Lindgren, A.C., Ye, J., 2004. Time se-  
706 ries classification using Gaussian mixture models of reconstructed phase  
707 spaces. *IEEE Trans. on Knowledge Discovery and Data Engineering* 16,  
708 779–783.
- 709 Rabiner, L.R., Juang, B., 1993. Fundamentals of speech recognition.  
710 Prentice-Hall, Upper Saddle River, USA.
- 711 Radovanović, M., Nanopoulos, A., Ivanovic, M., 2010. Time-series classifi-  
712 cation in many intrinsic dimensions, in: Proc. of the SIAM Int. Conf. on  
713 Data Mining, pp. 677–688.
- 714 Ramoni, M., Sebastiani, P., Cohen, P., 2002. Bayesian clustering by dynam-  
715 ics. *Machine Learning* 47, 91–121.
- 716 Rodríguez, J.J., Alonso, C., 2004. Interval and dynamic time warping-based  
717 decision trees, in: Proc. of the ACM Symp. on Applied Computing (SAC),  
718 pp. 548–552.
- 719 Sakoe, H., Chiba, S., 1978. Dynamic programming algorithm optimization  
720 for spoken word recognition. *IEEE Trans. on Acoustics, Speech, and*  
721 *Language Processing* 26, 43–50.
- 722 Salvador, S., Chan, P., 2007. Toward accurate dynamic time warping in  
723 linear time and space. *Intelligent Data Analysis* 11, 561–580.
- 724 Salzberg, S.L., 1997. On comparing classifiers: pitfalls to avoid and a rec-  
725 ommended approach. *Data Mining and Knowledge Discovery* 1, 317–328.

- 726 Serrà, J., Arcos, J.L., 2012. A competitive measure to assess the similar-  
727 ity between two time series, in: Proc. of the Int. Conf. on Case-Based  
728 Reasoning (ICCBR), pp. 414–427.
- 729 Serrà, J., Kantz, H., Serra, X., Andrzejak, R.G., 2012a. Predictability of  
730 music descriptor time series and its application to cover song detection.  
731 IEEE Trans. on Audio, Speech and Language Processing 20, 514–525.
- 732 Serrà, J., Zanin, M., Herrera, P., Serra, X., 2012b. Characterization and ex-  
733 ploitation of community structure in cover song networks. Pattern Recog-  
734 nition Letters 33, 1032–1041.
- 735 Tormene, P., Giorgino, T., Quaglini, S., Stefanelli, M., 2009. Matching  
736 incomplete time series with dynamic time warping: an algorithm and an  
737 application to post-stroke rehabilitation. Artificial Intelligence in Medicine  
738 45, 11–34.
- 739 Wang, X., Mueen, A., Ding, H., Trajcevski, G., Scheuermann, P., Keogh,  
740 E.J., 2012. Experimental comparison of representation methods and dis-  
741 tance measures for time series data. Data Mining and Knowledge Discov-  
742 ery In press. URL: <http://dx.doi.org/10.1007/s10618-012-0250-5>.
- 743 Wu, D., Agrawal, D., El Abbadi, A., Singh, A., Smith, T.R., 1996. Efficient  
744 retrieval for browsing large image databases, in: Proc. of the Int. Conf.  
745 on Knowledge Information, pp. 11–18.
- 746 Wu, Y.L., Agrawal, D., El Abbadi, A., 2000. A comparison of DFT and  
747 DWT based similarity search in time-series databases, in: Proc. of the  
748 Int. Conf. on Information and Knowledge Management (CIKM), pp. 488–  
749 495.
- 750 Xi, X., Keogh, E.J., Shelton, C.R., Wei, L., Ratanamahatana, C.A., 2006.  
751 Fast time series classification using numerosity reduction, in: Proc. of the  
752 Int. Conf. on Machine Learning, pp. 1033–1040.
- 753 Xing, Z., Pei, J., Yu, P.S., 2011. Early classification on time series. Knowl-  
754 edge and Information Systems 31, 105–127.