




ABBA: adaptive Brownian bridge-based symbolic aggregation of time series

Steven Elsworth¹ · Stefan Güttel¹ 

Received: 29 May 2019 / Accepted: 14 May 2020 / Published online: 3 June 2020
© The Author(s) 2020

Abstract

A new symbolic representation of time series, called ABBA, is introduced. It is based on an adaptive polygonal chain approximation of the time series into a sequence of tuples, followed by a mean-based clustering to obtain the symbolic representation. We show that the reconstruction error of this representation can be modelled as a random walk with pinned start and end points, a so-called Brownian bridge. This insight allows us to make ABBA essentially parameter-free, except for the approximation tolerance which must be chosen. Extensive comparisons with the SAX and 1d-SAX representations are included in the form of performance profiles, showing that ABBA is often able to better preserve the essential shape information of time series compared to other approaches, in particular when time warping measures are used. Advantages and applications of ABBA are discussed, including its in-built differencing property and use for anomaly detection, and Python implementations provided.

Keywords Time series · Symbolic aggregation · Dimension reduction · Brownian bridge

1 Introduction

Symbolic representations of time series are an active area of research, being useful for many data mining tasks including dimension reduction, motif and rule discovery, prediction, and clustering of time series. Symbolic time series representations allow for the use of algorithms from text processing and bioinformatics, which often take

Responsible editor: Panagiotis Papapetrou.

✉ Stefan Güttel
stefan.guettel@manchester.ac.uk

Steven Elsworth
steven.elsworth@manchester.ac.uk

¹ Department of Mathematics, The University of Manchester, Manchester M13 9PL, UK

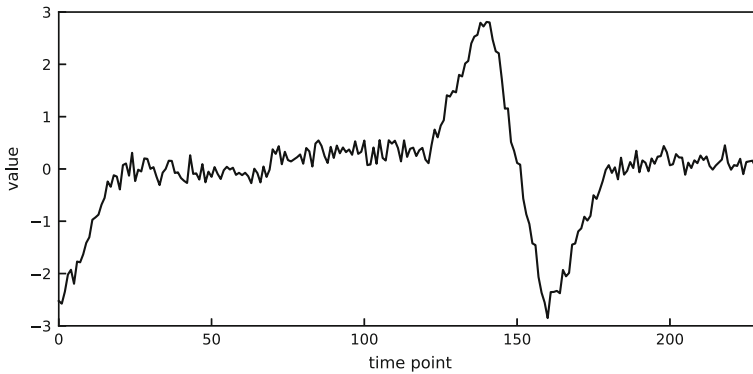


Fig. 1 Illustrative example time series T used throughout the paper

advantage of the discrete nature of the data. Our focus in this work is to develop a symbolic representation which is dimension reducing whilst preserving the essential *shape* of the time series. Our definition of shape is different from the one commonly implied in the context of time series: we focus on representing the peaks and troughs of the time series in their correct order of appearance, but we are happy to slightly stretch the time series in *both* the time and value directions. In other words, our focus is not necessarily on approximating the time series values at the correct time points, but on representing the local up-and-down behavior of the time series and identifying repeated motifs. This is obviously not appropriate in all applications, but we believe it is close to how humans summarize the overall behavior of a time series, and in that our representation might be useful for trend prediction, anomaly detection, and motif discovery.

To illustrate, let us consider the time series shown in Fig. 1. This series is sampled at equidistant time points with values $t_0, t_1, \dots, t_N \in \mathbb{R}$, where $N = 230$. There are various ways of describing this time series, for example:

- (a) It is exactly representable as a high-dimensional vector $T = [t_0, t_1, \dots, t_N] \in \mathbb{R}^{N+1}$.
- (b) It starts at a value of about -3 , then climbs up to a value of about 0 within 25 time steps, then it stays at about 0 for 100 time steps, after which it goes up to a value of about 3 within 25 time steps, and so on.
- (c) It starts at a value of about -3 , then goes up rapidly by about 3 units, followed by a longer period with almost no change in value, after which it *again* goes up rapidly by about 3 units, and so on.

Note how in (a) and (b) the emphasis is on the actual values of the time series, whereas in (c) we mainly refer to trends in the time series in relation to previously observed trends. High-level information might be difficult to extract from (a) directly, while (b) could be seen as putting too much emphasis on the time series values instead of the overall shape. The symbolic representation developed in this paper, called *adaptive Brownian bridge-based aggregation* (ABBA), adaptively reduces T to a shorter sequence of symbols with an emphasis on the shape information. The resulting description will be conceptually similar to (c) from the examples above.

To formalize the discussion and introduce notation, we consider the problem of aggregating a time series $T = [t_0, t_1, \dots, t_N] \in \mathbb{R}^{N+1}$ into a *symbolic representation* $S = [s_1, s_2, \dots, s_n] \in \mathbb{A}^n$, where each s_j is an element of an alphabet $\mathbb{A} = \{a_1, a_2, \dots, a_k\}$ of k symbols. The sequence S should be of considerably lower dimension than the original time series T , that is $n \ll N$, and it should only use a small number of meaningful symbols, that is $k \ll n$. The representation should also allow for the approximate reconstruction of the original time series with a controllable error, with the shape of the reconstruction suitably close to that of the original. Both n , the length of the symbolic representation, and k , the number of symbols, should be chosen automatically without parameter tuning required.

This paper is organized as follows. In Sect. 2 we give an overview of existing symbolic representations and other algorithms which are conceptually similar to ABBA. To evaluate the approximation accuracy of ABBA, we must compare the shape of the original time series and the reconstruction from its symbolic representation. Section 3 reviews existing distance measures for this purpose and discusses how well they perform in measuring shape. Sections 4–7 contain the key contributions of this paper:

- Section 4 introduces ABBA, our novel dimension-reducing symbolic time series representation which aims to preserve the shape of the original time series. We explain in detail how ABBA's compression and reconstruction procedures work.
- In Sect. 5 we show that the error of the ABBA reconstruction behaves like a random walk with pinned start and end values. This observation appears to be novel in itself and allows us to balance the error of the piecewise linear approximation with that of the digitization procedure, thereby allowing the method to choose the number of symbols k automatically.
- Section 6 contains performance comparisons of ABBA with other popular symbolic representations using various distance measures, with a particular emphasis on the compression versus accuracy relation. Aside from verifying that ABBA can represent time series to higher accuracy than SAX and 1d-SAX using a comparable number of symbols k and string length n , we also find that SAX outperforms 1d-SAX when the same number of symbols k is used for both.
- In Sect. 7 we discuss some practical applications of ABBA including the handling of linear trends, anomaly detection, and VizTree visualization.

Finally, we conclude in Sect. 8 with an outlook on future work.

2 Background and related work

Despite the large number of dimension-reducing time series representations in the literature, very few are *symbolic*. Most techniques are *numeric* in the sense that they reduce a time series to a lower-dimensional vector with its components taken from a continuous range; see Bettaiah and Ranganath (2014), Fu (2011), Lin et al. (2007) for reviews. Here we provide an overview of existing symbolic representations relevant to ABBA.

The construction of symbolic time series representations typically consists of two parts. First, the time series is segmented, with the length of each segment being either

specified by the user or found adaptively via a bottom-up, top-down, or sliding window approach (Keogh et al. 2001). The segmentation procedure intrinsically controls the degree of dimension reduction. The second part, the discretization process, assigns a symbol to each segment.

Symbolic Aggregate approxImation (SAX), a very popular symbolic representation, consists of a piecewise approximation of the time series followed by a symbolic conversion using Gaussian breakpoints (Lin et al. 2007). SAX starts by partitioning T into segments of constant length $l \in \mathbb{N}$, and then represents each segment by the mean of its values (i.e., a piecewise constant approximation). The means are converted into symbols using breakpoints that partition a Gaussian bell curve into k equally-sized areas. In addition to its simplicity, an attractive feature of SAX is the existence of distance measures that serve as lower bounds for the Euclidean distance between the original time series. On the other hand, both the segment length $l \in \mathbb{N}$ and the number of symbols k must be specified in advance. SAX is designed such that each symbol appears with equal probability, which works best when the time series values are approximately normally distributed.

The literature on applications of SAX is extensive and many variants have been proposed. Most variants modify the symbolic representation to incorporate the slope of the time series on each segment. This is often justified by applications in finance, where the extreme values of time series provide valuable information which is lost with the piecewise constant approximation used in SAX. The modifications often come at the cost of losing the lower bounds on distance measures. We now provide a brief overview of some of these variants.

Trend-based and Valued-based Approximation (TVA) uses SAX to symbolically represent the time series values, enhanced with U, D, or S symbols to represent an upwards, downwards, or straight trend, respectively (Esmael et al. 2012). The TVA representation alternates between value symbols and slope symbols, making the symbolic representation twice as long as a SAX representation with the same number of segments. A similar approach is *Trend-based SAX* (TSAX) which uses two trend symbols per segment (Zhang et al. 2018).

Extended SAX (ESAX) represents each segment by the minimum, maximum, and mean value of the time series ordered according to their appearance in the segment, defining the mean to appear in the center of the segment (Lkhagva et al. 2006). This results in a symbolic representation three times longer than the corresponding SAX representation with the same number of segments. *Enhanced SAX* (EN-SAX) forms a vector for each segment consisting of the minimum, maximum and mean value. The vectors are then clustered and a symbol is allocated to each cluster (Barnaghi et al. 2012). *Time-Weighted Average for SAX* (TWA_SAX) uses the time weighted average for each segment instead of the mean (Benyahmed et al. 2015). This can encapsulate important patterns which are missed by the mean.

Trend-based Symbolic approximation (TSX) represents each segment by four symbols (Li et al. 2012). The first symbol corresponds to the SAX representation. The following three symbols correspond to the slopes between the first, last, most peak and most dip points, which are defined in terms of vertical distance from the trend line (the straight line connecting the end point values of a segment). The slopes are

converted to symbols using a lookup table. This results in a symbolic representation four times longer than the SAX representation with the same number of segments.

The *1d-SAX algorithm* uses linear regression to fit a straight line to each segment (Malinowski et al. 2013). Each segment is then represented by the gradient and the average value of the line. Two sets of Gaussian breakpoints are used to provide symbols for both the averages and the slopes. It is unclear how many breakpoints should be allocated for the averages, and how many should be allocated for the slopes. The total number of symbols is the product of the respective number of breakpoints.

Using the same number of segments, the above SAX variants result in an increase in the length of the symbolic representation by some factor. It is unclear whether any of these approaches performs better than SAX when the SAX segment length len is decreased by the same factor (keeping the overall length of the symbolic representation constant). As with the original SAX approach, all of these variants require the user to specify the segment length len and the number of symbols k in advance.

In many time series applications, the assumption that the values of the normalized time series follow a normal distribution is a strong one. To overcome this, the *adaptive SAX algorithm* (aSAX) uses k -means clustering to find the breakpoints for the symbolic conversion (Pham et al. 2010). However, as piecewise constant approximations are used, the aSAX approach fails to represent the extreme points of the time series.

SAX's digitization procedure based on Gaussian breakpoints allows its extension to a multi-resolution symbolic representation known as indexable SAX (iSAX) (Shieh and Keogh 2008). This clever indexing procedure allows mining of datasets containing millions of time series. At the heart of the algorithm is a SAX representation where each window uses Gaussian breakpoints with 2^c regions, where c can change from segment to segment.

The *sensorPCA algorithm* overcomes the fixed window length problem by using a sliding window to start a new segment when the standard deviation of the approximation exceeds some prespecified tolerance (Ganz et al. 2013). However, Ganz et al. (2013) does not provide a method to convert the mean values and window lengths to a symbolic representation.

Symbolic Aggregate approximation Optimized by data (SAXO) is a data-driven approach based on a regularized Bayesian coclustering method called minimum optimized description length (Bondu et al. 2016; Boullé 2006). The discretization of the time series is optimized using Bayesian statistics. The number of symbols and the underlying distribution change for each time interval. The computational complexity of SAXO is far greater than that of SAX.

Mörchen and Ultsch (2006) take a completely different approach based on the persistence of a time series. A persistent time series is one where the value at a certain point is closely related to the previous value; see also Kim (2000). The authors provide “persist”, a symbolic representation based on the Kullback–Leibler divergence between the marginal and the self-transition probability distributions of the discretization symbols.

Symbolic Polynomial (SP) is a symbolic representation designed to detect local patterns (Grabocka et al. 2014). It is constructed by an overlapping sliding window of length w and stepsize 1. For each window, one computes the coefficients of a regression polynomial of degree d . The coefficients of each order are collected and allocated a

symbol using an equi-area discretization. This symbolic representation provides no dimensional reduction as each window is represented by d symbols.

Baydogan and Runger (2015) introduce a symbolic representation of multivariate time series called SMTS. They construct a data table consisting of time index, time values, and first differences of the time series. A tree learner is trained on the data and each of the leaf nodes is allocated a symbol. Their approach allows multiple tree learners, which in the univariate case results in a symbolic representation much larger than the original.

Piecewise linear approximations of time series have been used for many years. The lengths of the linear pieces (segments) can be prespecified or chosen adaptively. Each segment is approximated using either linear interpolation or linear regression (Keogh et al. 2001). Luo et al. (2015) describe how the linear segments can be stitched so that each piece is represented by two parameters rather than three. An example of a piecewise linear approximation algorithm is the Ramer–Douglas–Peucker algorithm, an iterative endpoint fitting procedure which uses adaptive linear interpolation with a prespecified tolerance. These methods provide an effective shape-preserving and dimension-reducing representation but not a symbolic representation.

3 Distance measures

The accuracy of a symbolic time series representation S can be assessed by the distance between the original time series T and its reconstruction \hat{T} from S . We note that the original time series should first be normalized to have zero mean and unit variance. This ensures that distance measures are comparable across different time series; see Keogh and Kasetty (2003) for a discussion of the importance of normalization.

A detailed overview of time series distance measures and their applications can be found in Aghabozorgi et al. (2015). Distance measures for time series typically fall into two main categories: *lock-step alignment* and *elastic alignment* (Abanda et al. 2019). Lock-step alignment refers to the element-wise comparison of time series, i.e., the i -th element of one time series is compared to the i -th element of another. Such measures can only compare time series of equal length. The most popular lock-step distance is the Euclidean distance. The Euclidean distance is a poor measure of shape similarity in two particular cases: if the time series have the same shape but are stretched in value (see Fig. 2a), or if the time series have the same shape but are warped in time (see Fig. 2b). The first issue can be mitigated by differencing the time series before measuring the distance. The second issue is intrinsic to lock-step alignment distance measures.

Elastic alignment distance measures construct a nonlinear mapping between time series elements, effectively allowing for one value in a time series to be compared to multiple consecutive values in another. The most popular elastic alignment method is Dynamic Time Warping (DTW), originally proposed in Berndt and Clifford (1994). The DTW distance measure corresponds to the Euclidean distance between two DTW-aligned time series. This distance measure can be used to compare time series of different lengths but it has a quadratic computational complexity in both time and space; for further details see Keogh and Ratanamahatana (2005). Many methods have

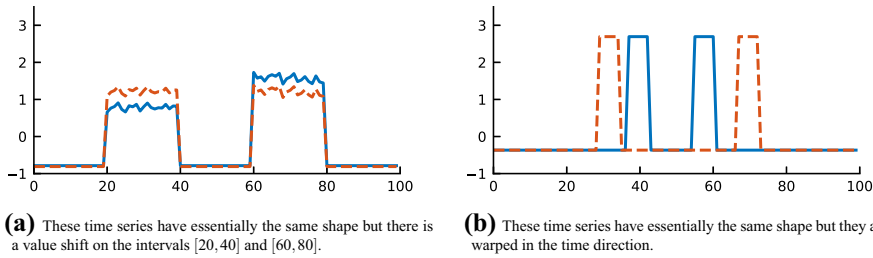


Fig. 2 The time series in these plots have the same essential shape according to our interpretation. Euclidean distance is a poor measure of shape for **a** and **b**, whereas DTW distance is a poor measure of shape for **(a)**. A differencing of the time series in **(a)** would make DTW a suitable shape distance

been proposed to either approximate the DTW distance at a reduced cost or calculate bounds to avoid computing the DTW alignment altogether. Keogh and Pazzani (2001) notice that DTW may pair a rising trend in one time series with a falling trend in another, and they overcome this problem by a variant known as Derivative Dynamic Time Warping (DDTW). The elastic alignment allows DTW to overcome the issues when two time series have the same shape but are warped in time (see Fig. 2b), but DTW is still a poor measure of shape similarity if the time series have the same shape but are vertically stretched (see Fig. 2a). Again, this can be mitigated by differencing the time series before measuring their DTW distance.

It is because of these advantages and drawbacks of the Euclidean and DTW distance measures and their differenced counterparts that we will test the performance of ABBA with all these distance measures in Sect. 6.

4 Adaptive Brownian bridge-based aggregation

We now introduce ABBA, a symbolic representation of time series where the symbolic length n and the number of symbols k are chosen adaptively. The ABBA representation is computed in two stages.

1. *Compression* The original time series T is approximated by a piecewise linear and continuous function, with each linear piece being chosen adaptively based on a user-specified tolerance. The result is a sequence of tuples (len, inc) consisting of the length of each piece and its increment in value.
2. *Digitization* A near-optimal alphabet \mathbb{A} is identified via mean-based clustering, with each cluster corresponding to a symbol. Each tuple (len, inc) is assigned a symbol corresponding to the cluster in which it belongs.

The reconstruction of a time series from its ABBA representation involves three stages.

1. *Inverse-digitization* Each symbol of the symbolic representation is replaced with the center of the associated cluster. The length values of the centers may not necessarily be integers.
2. *Quantization* The lengths of the reconstructed segments are re-aligned with an integer grid.

Table 1 Summary of notation

Original time series:	$T = [t_0, t_1, \dots, t_N] \in \mathbb{R}^{N+1}$
After compression:	$[(\text{len}_1, \text{inc}_1), (\text{len}_2, \text{inc}_2), \dots, (\text{len}_n, \text{inc}_n)] \in \mathbb{R}^{2 \times n}$
After digitization:	$S = [s_1, s_2, \dots, s_n] \in \mathbb{A}^n$ with $\mathbb{A} = \{a_1, a_2, \dots, a_k\}$
After inverse-digitization:	$[(\widetilde{\text{len}}_1, \widetilde{\text{inc}}_1), (\widetilde{\text{len}}_2, \widetilde{\text{inc}}_2), \dots, (\widetilde{\text{len}}_n, \widetilde{\text{inc}}_n)] \in \mathbb{R}^{2 \times n}$
After quantization:	$[(\widehat{\text{len}}_1, \widehat{\text{inc}}_1), (\widehat{\text{len}}_2, \widehat{\text{inc}}_2), \dots, (\widehat{\text{len}}_n, \widehat{\text{inc}}_n)] \in \mathbb{R}^{2 \times n}$
After inverse-compression:	$\widehat{T} = [\widehat{t}_0, \widehat{t}_1, \dots, \widehat{t}_N] \in \mathbb{R}^{N+1}$

3. *Inverse-compression* The piecewise linear continuous approximation is converted back to a pointwise time series representation using a stitching procedure.

Both the computation of the ABBA representation and the reconstruction are inexpensive. It is essential that the digitization process uses incremental changes in value rather than slopes. This way, ABBA consistently works with increments in both the time and value coordinates. Only in this case a mean-based clustering algorithm will identify meaningful clusters in both coordinate directions. As we will explain in Sect. 5, the error of the ABBA reconstruction behaves like a random walk pinned at zero for both the start and the end point of the time series. But first, we provide a more detailed explanation of the key parts of ABBA. For clarity, we summarize the notation used throughout this section in Table 1.

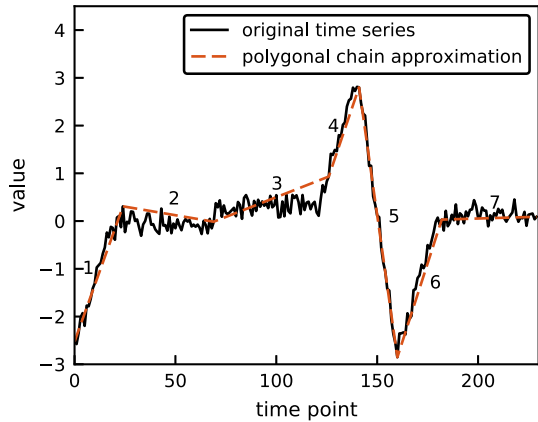
4.1 Compression

The ABBA compression is achieved by an adaptive piecewise linear continuous approximation of T . Given a tolerance tol , the method adaptively selects $n+1$ indices $i_0 = 0 < i_1 < \dots < i_n = N$ so that the time series $T = [t_0, t_1, \dots, t_N]$ is approximated by a polygonal chain going through the points (i_j, t_{i_j}) for $j = 0, 1, \dots, n$. This gives rise to a partition of T into n pieces $P_j = [t_{i_{j-1}}, t_{i_{j-1}+1}, \dots, t_{i_j}]$, each of length $\text{len}_j := i_j - i_{j-1} \geq 1$ in the time direction. We ensure that the squared Euclidean distance of the values in P_j from the straight polygonal line is bounded by $(\text{len}_j - 1) \cdot \text{tol}^2$. More precisely, starting with $i_0 = 0$ and given an index i_{j-1} , we find the largest possible i_j such that $i_{j-1} < i_j \leq N$ and

$$\sum_{i=i_{j-1}}^{i_j} \left(\underbrace{t_{i_{j-1}} + (t_{i_j} - t_{i_{j-1}}) \cdot \frac{i - i_{j-1}}{i_j - i_{j-1}}}_{\text{straight line approximation}} - \underbrace{t_i}_{\text{actual value}} \right)^2 \leq (i_j - i_{j-1} - 1) \cdot \text{tol}^2. \tag{1}$$

Note that the first and the last values $t_{i_{j-1}}$ and t_{i_j} are not counted in the distance measure as the straight line approximation passes exactly through them. If required, one can restrict the maximum length of each segment by imposing an upper bound $i_j \leq i_{j-1} + \text{max_len}$ with a given integer $\text{max_len} \geq 1$.

Fig. 3 Result of the ABBA compression. The time series is now represented by $n = 7$ tuples of the form (inc, len) and the starting value t_0



Each linear piece P_j of the resulting polygonal chain \tilde{T} is described by a tuple (len_j, inc_j) , where $inc_j = t_{i_j} - t_{i_{j-1}}$ is the increment in value (not the slope!). As the polygonal chain is continuous, the first value of a segment can be inferred from the end value of the previous segment. Hence the whole polygonal chain can be recovered exactly from the first value t_0 and the tuple sequence

$$(len_1, inc_1), (len_2, inc_2), \dots, (len_n, inc_n) \in \mathbb{R}^2. \tag{2}$$

An example of the ABBA compression procedure applied to the time series in Fig. 1 is shown in Fig. 3. Here a tolerance of $\tau_{ol} = 0.4$ has been used, resulting in $n = 7$ pieces. As the approximation error on each piece P_j satisfies (1), the polygonal chain \tilde{T} also has a bounded Euclidean distance from T :

$$\begin{aligned} \text{euclid}(T, \tilde{T})^2 &\leq [(i_1 - i_0 - 1) + (i_2 - i_1 - 1) + \dots \\ &\quad + (i_n - i_{n-1} - 1)] \cdot \tau_{ol}^2 \\ &= (N - n) \cdot \tau_{ol}^2. \end{aligned} \tag{3}$$

Hence we are sure that the ABBA approximation \tilde{T} (red dashed curve) in Fig. 3 has a Euclidean distance of at most $\sqrt{223} \times 0.4 \approx 6.0$ from the original time series T (black solid curve).

Let us comment on the computational complexity of this compression phase. Assuming that all pieces have an average length \overline{len} , the evaluation of (1) for each $i_j = i_{j-1} + 2, i_{j-1} + 3, \dots, i_{j-1} + \overline{len}$ involves a sum of $1, 2, \dots, \overline{len} - 1$ terms, respectively. (Recall that the straight line approximation passes through the end points and hence the corresponding terms in the sum need not be evaluated.) Therefore, the overall number of summands to be evaluated and added up is $1 + 2 + \dots + \overline{len} - 1 = \mathcal{O}(\overline{len}^2)$ per piece. Under the natural assumption that the average piece length \overline{len} is independent of the length N of the time series, we have $N \propto n \cdot \overline{len}$. Hence, the compression phase has a linear complexity of $\mathcal{O}(N) = \mathcal{O}(n)$.

4.2 Digitization

Digitization refers to the assignment of the tuples in (2) to k clusters S_1, S_2, \dots, S_k . Before clustering, we separately normalize the tuple lengths and increments by their standard deviations σ_{len} and σ_{inc} , respectively. We use a further scaling parameter scl to assign different weight (“importance”) to the length of each piece in relation to its increment value. Hence, we effectively cluster the *scaled tuples*

$$\left(\text{scl} \frac{\text{len}_1}{\sigma_{\text{len}}}, \frac{\text{inc}_1}{\sigma_{\text{inc}}} \right), \left(\text{scl} \frac{\text{len}_2}{\sigma_{\text{len}}}, \frac{\text{inc}_2}{\sigma_{\text{inc}}} \right), \dots, \left(\text{scl} \frac{\text{len}_n}{\sigma_{\text{len}}}, \frac{\text{inc}_n}{\sigma_{\text{inc}}} \right) \in \mathbb{R}^2. \quad (4)$$

If $\text{scl} = 0$, then clustering is performed on the increments alone, while if $\text{scl} = 1$, we cluster in both the length and increment dimension with equal weighting. The cluster assignment is performed by (approximately) minimizing the within-cluster-sum-of-squares

$$\text{WCSS} = \sum_{i=1}^k \sum_{(\text{len}, \text{inc}) \in S_i} \left\| \left(\text{scl} \frac{\text{len}}{\sigma_{\text{len}}}, \frac{\text{inc}}{\sigma_{\text{inc}}} \right) - \bar{\mu}_i \right\|^2,$$

with each 2d cluster center $\bar{\mu}_i = (\bar{\mu}_i^{\text{len}}, \bar{\mu}_i^{\text{inc}})$ corresponding to the mean of the scaled tuples associated with the cluster S_i . In certain situations one may want to cluster only on the lengths of the pieces and ignore their increments, formally setting $\text{scl} = \infty$. In this case, the cluster assignment is performed by (approximately) minimizing

$$\text{WCSS} = \sum_{i=1}^k \sum_{(\text{len}, \text{inc}) \in S_i} \left| \frac{\text{len}}{\sigma_{\text{len}}} - \bar{\mu}_i^{\text{len}} \right|^2,$$

where $\bar{\mu}_i^{\text{len}}$ is the mean of the scaled lengths in the cluster S_i .

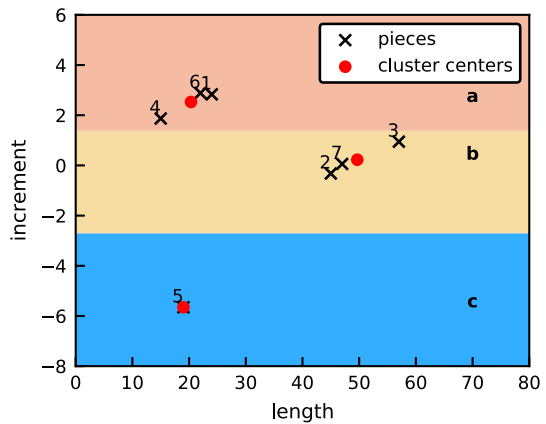
Given a clustering of the n tuples into clusters S_1, \dots, S_k we use the *unscaled* cluster centers μ_i

$$\mu_i = (\mu_i^{\text{len}}, \mu_i^{\text{inc}}) = \frac{1}{|S_i|} \sum_{(\text{len}, \text{inc}) \in S_i} (\text{len}, \text{inc})$$

to define the maximal cluster variances in the length and increment directions as

$$\begin{aligned} \text{Var}_{\text{len}} &= \max_{i=1, \dots, k} \frac{1}{|S_i|} \sum_{(\text{len}, \text{inc}) \in S_i} |\text{len} - \mu_i^{\text{len}}|^2, \\ \text{Var}_{\text{inc}} &= \max_{i=1, \dots, k} \frac{1}{|S_i|} \sum_{(\text{len}, \text{inc}) \in S_i} |\text{inc} - \mu_i^{\text{inc}}|^2, \end{aligned}$$

Fig. 4 Result of the ABBA digitization with scaling parameter $s_{cl} = 0$. The tuples (len, inc) are converted to the symbol sequence $abbacab$



respectively. Here, $|S_i|$ is the number of tuples in cluster S_i . We seek the smallest number of clusters k such that

$$\max(s_{cl} \cdot \text{Var}_{len}, \text{Var}_{inc}) \leq \text{tol}_s^2 \tag{5}$$

with a tolerance tol_s . This tolerance will be specified in Sect. 5 as a function of the user-specified tolerance tol and is therefore not a free parameter. (In the case of $s_{cl} = \infty$, we seek the smallest k such that $\text{Var}_{len} \leq \text{tol}_s^2$.) Once the optimal k has been found, each cluster S_1, \dots, S_k is assigned a symbol a_1, \dots, a_k , respectively. Finally, each tuple in the sequence (2) is replaced by the symbol of the cluster it belongs to, resulting in the symbolic representation $S = [s_1, s_2, \dots, s_n]$.

If $s_{cl} = 0$ or $s_{cl} = \infty$, a 1d clustering method can be used which takes advantage of sorting algorithms; see the review by Grønlund et al. (2017). We use the `ckmeans` algorithm (Wang and Song 2011), an order $\mathcal{O}(n \log n + kn)$ dynamic programming algorithm which optimally clusters the data by minimizing the WCSS in just one dimension. We have modified the algorithm to choose the smallest k such that the maximal cluster variance is bounded by tol_s^2 .

For nonzero finite values of s_{cl} , k -means clustering is used. This algorithm has an average complexity of $\mathcal{O}(kn)$ per iteration [see also Arthur and Vassilvitskii (2006) for an analysis of the worst case complexity] and might of course result in a sub-optimal clustering. In our ABBA implementation the user can specify an interval $[\text{min_k}, \dots, \text{max_k}]$ and we search for the smallest k in that interval such that (5) holds. If no such k exists, we set $k = \text{max_k}$.

By default, we set $s_{cl} = 0$ as we believe this corresponds most naturally to preserving the up-and-down behavior of the time series. In other words, we ignore the lengths of the pieces and only cluster the value increments. With the value increments represented accurately, the errors in lengths correspond to horizontal stretching in the time direction.

An illustration of the digitization process on the pieces from Fig. 3 can be seen in Fig. 4 with $s_{cl} = 0$ (our default parameter choice), Fig. 5 with $s_{cl} = 1$, and Fig. 6 with $s_{cl} = \infty$.

Fig. 5 Result of the ABBA digitization with $scl = 1$. The tuples (len, inc) are converted to the symbol sequence `abbacab`

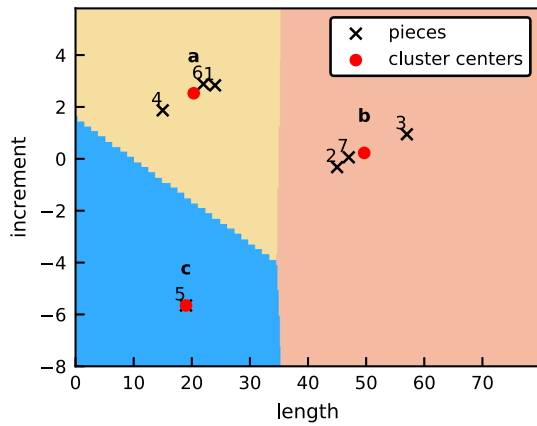
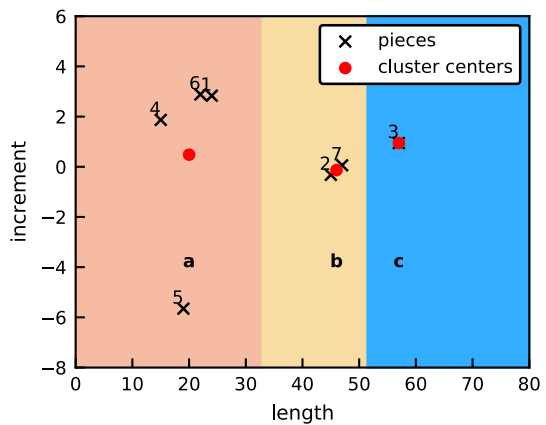


Fig. 6 Result of the ABBA digitization with $scl = \infty$. The tuples (len, inc) are converted to the symbol sequence `abcaaab`



4.3 Inverse digitization and quantization

When reversing the digitization process, each symbol of the alphabet is replaced by the center $(\overline{len}_i, \overline{inc}_i)$ of the corresponding cluster given as

$$(\overline{len}_i, \overline{inc}_i) = \frac{1}{|S_i|} \sum_{(len, inc) \in S_i} (len, inc).$$

Note that the mean-based clustering for digitization is performed on the scaled tuples (4), but the cluster centers used for the inverse digitization are computed with the unscaled tuples (2). The inverse digitization process results in a sequence of n tuples

$$(\widetilde{len}_1, \widetilde{inc}_1), (\widetilde{len}_2, \widetilde{inc}_2), \dots, (\widetilde{len}_n, \widetilde{inc}_n) \in \mathbb{R}^2,$$

where each tuple is a cluster center, that is $(\widetilde{len}_i, \widetilde{inc}_i) \in \{(\overline{len}_1, \overline{inc}_1), (\overline{len}_2, \overline{inc}_2), \dots, (\overline{len}_k, \overline{inc}_k)\}$.

The lengths $\widetilde{\text{len}}_i$ obtained from this averaging are not necessarily integer values as they were in the compressed representation (2). We therefore perform a simple quantization procedure which realigns the cumulated lengths with their closest integers. We start with rounding the first length, $\widehat{\text{len}}_1 := \text{round}(\widetilde{\text{len}}_1)$, keeping track of the rounding error $e := \widetilde{\text{len}}_1 - \widehat{\text{len}}_1$. This error is added to the second length $\widetilde{\text{len}}_2 := \widetilde{\text{len}}_2 + e$, which is then rounded to $\widehat{\text{len}}_2 := \text{round}(\widetilde{\text{len}}_2)$ with error $e := \widetilde{\text{len}}_2 - \widehat{\text{len}}_2$, and so on. As a result we obtain a sequence of n tuples

$$(\widehat{\text{len}}_1, \widehat{\text{inc}}_1), (\widehat{\text{len}}_2, \widehat{\text{inc}}_2), \dots, (\widehat{\text{len}}_n, \widehat{\text{inc}}_n) \in \mathbb{R}^2 \tag{6}$$

with integer lengths $\widehat{\text{len}}_i$. (The increments remain unchanged but we rename them for consistency: $\widehat{\text{inc}}_i := \widetilde{\text{inc}}_i$.)

5 Error analysis

During the compression procedure, we construct a polygonal chain \widetilde{T} going through selected points $\{(i_j, t_{i_j})\}_{j=0}^n$ of the original time series T , with a controllable Euclidean distance (3). After the digitization, inverse digitization, and quantization, we obtain a new tuple sequence (6) which can be stitched together to a polygonal chain \widehat{T} going through the points $\{(\widehat{i}_j, \widehat{t}_j)\}_{j=0}^n$, with $(\widehat{i}_0, \widehat{t}_0) = (0, t_0)$. Our aim is to analyze the distance between \widehat{T} and \widetilde{T} , and then balance it with the distance between \widetilde{T} and T .

We first note that

$$(\widehat{i}_j, \widehat{t}_j) = \left(\sum_{\ell=1}^j \widehat{\text{len}}_\ell, t_0 + \sum_{\ell=1}^j \widehat{\text{inc}}_\ell \right), \quad j = 0, \dots, n.$$

As all the lengths $\widehat{\text{len}}_\ell$ and increments $\widehat{\text{inc}}_\ell$ correspond to cluster centers (averages of all the points in a cluster, consistently rounded during quantization), we have the interesting property that the accumulated deviations from the true lengths and increments exactly cancel out at the right endpoint of the last piece P_n , that is: $(\widehat{i}_n, \widehat{t}_n) = (i_n, t_{i_n}) = (N, t_N)$. In other words, the polygonal chain \widehat{T} starts and ends at the same values as T (and hence T).

We now analyze the behavior of \widehat{T} in between the start and endpoints, focusing on the case that $\text{scl} = 0$ and assuming for simplicity that all cluster centers S_i have the same mean length $\mu_i^{\text{len}} = N/n$. (This is not a strong assumption as in the dynamic time warping distance the lengths of the pieces is irrelevant.) We compare \widehat{T} with the polygonal chain \widetilde{T} time-warped to the same regular length grid as \widehat{T} , which will give an upper bound on $\text{dtw}(\widehat{T}, \widetilde{T})$. Denoting by $d_\ell := \widehat{\text{inc}}_\ell - \widetilde{\text{inc}}_\ell$ the local deviation of the increment value of \widehat{T} on piece P_ℓ from the true increment of \widetilde{T} , we have that

$$\widehat{t}_j - t_{i_j} = \sum_{\ell=1}^j d_\ell =: e_{i_j}, \quad j = 0, \dots, n.$$

Recall from Sect. 4.2 that we have controlled the variance of the increment values in each cluster to be bounded by τol_s^2 . As a consequence, the increment deviations d_ℓ have bounded variance τol_s^2 , and mean zero as they correspond to deviations from their respective cluster center. It is therefore reasonable to model the “global increment errors” e_{i_j} as a random process with fixed values $e_{i_0} = e_{i_n} = 0$, expectation $E(e_{i_j}) = 0$, and variance

$$\text{Var}(e_{i_j}) = \tau\text{ol}_s^2 \cdot \frac{j(n-j)}{n}, \quad j = 0, \dots, n.$$

In the case that the d_ℓ are i.i.d. normally distributed, such a process is known as a *Brownian bridge*. See also Fig. 7 for an illustration.

Note that so far we have only considered the variance of the global increment errors e_{i_j} at the left and right endpoints of each piece P_j , but we are actually interested in analyzing the error of the reconstruction \widehat{T} on the fine time grid. To this end, we now consider a “worst-case” realization of e_{i_j} which stays s standard deviations away from its zero mean. That is, we consider a realization

$$e_{i_j} = s \cdot \tau\text{ol}_s \cdot \sqrt{\frac{j(n-j)}{n}}, \quad j = 0, \dots, n.$$

By piecewise linear interpolation of these errors from the coarse time grid i_0, i_1, \dots, i_n to the fine time grid $i = 0, 1, \dots, N$ (in accordance with the linear stitching procedure used in ABBA), we find that

$$e_i \leq \sqrt{\frac{n}{N}} \cdot s \cdot \tau\text{ol}_s \cdot \sqrt{\frac{i(N-i)}{N}}, \quad i = 0, \dots, N,$$

using that the interpolated quadratic function on the right-hand side is concave. We can now bound the squared Euclidean norm of this fine-grid “worst-case” realization as

$$\sum_{i=0}^N e_i^2 \leq \frac{n \cdot s^2 \cdot \tau\text{ol}_s^2}{N^2} \cdot \sum_{i=0}^N i(N-i) = \frac{n \cdot s^2 \cdot \tau\text{ol}_s^2}{N^2} \cdot \frac{N^3 - N}{6} \leq n \cdot s^2 \cdot \tau\text{ol}_s^2 \cdot \frac{N}{6}.$$

This is a probabilistic bound on squared Euclidean error caused by a “worst-case” realization of the Brownian bridge, and thereby a probabilistic bound on the error incurred from the digitization procedure. Equating this bound with the bound (3) on the accuracy of the compression, we find that we should choose

$$\tau\text{ol}_s = \frac{\tau\text{ol}}{s} \sqrt{\frac{6(N-n)}{Nn}},$$

with the user-specified tolerance τol . We have experimentally determined that $s = 0.2$ typically gives a good balance between the compression accuracy and the number of clusters determined using this criterion.

Fig. 7 Example of the ABBA reconstruction error forming a Brownian bridge. The blue line is the actual error, the grey lines are 50 other realizations of the random walk, and the red bounds indicate one standard deviation above and below the zero mean

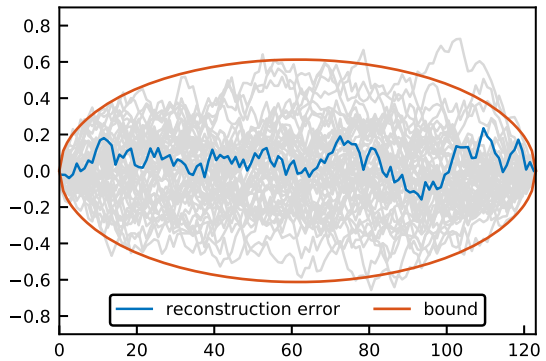
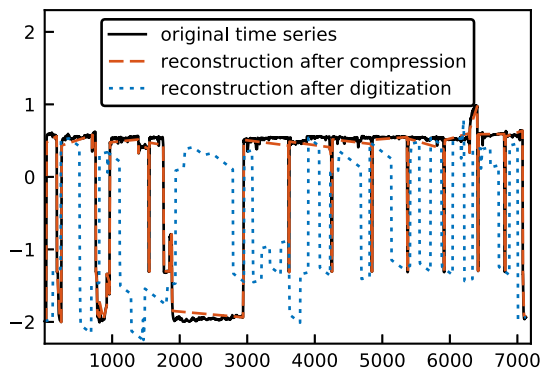


Fig. 8 ABBA representation of a time series from a heat exchanger in an ethylene cracker. With $\text{tol} = 0.1$ and $\text{scl} = 0$, the time series is reduced from 7128 points to 123 tuples using 14 symbols



Example: We now illustrate the above analysis on a challenging real-world example. Consider a time series T ($N = 7127$) consisting of temperature readings off a heat exchanger in an ethylene cracker. We use $\text{tol} = 0.1$ to compress this time series, resulting in a polygonal chain \tilde{T} with $n = 123$ pieces and an approximation error of $\text{euclid}(T, \tilde{T}) = 5.3 \leq \sqrt{N - n} \cdot \text{tol} \approx 8.4$. See Fig. 8 for a plot of the original time series T and its reconstruction \tilde{T} after compression.

We then run the ABBA digitization procedure with scaling parameter $\text{scl} = 0$, resulting in a symbolic representation S of length n using $k = 14$ symbols. In Fig. 7 we show the “global increment errors” e_{ij} of the reconstruction \hat{T} on each piece P_j , that is, the increment deviation of \hat{T} from T at the endpoints of P_j , $j = 1, \dots, n$. Note how this error is pinned at zero at $j = 0$ and $j = n$, and how it resembles a random walk in between.

The reconstruction \hat{T} on the fine time grid is also shown in Fig. 8. The reconstruction error measured in the time warping distance is $\text{dtw}(\tilde{T}, \hat{T}) = 9.5$ and the overall error is $\text{dtw}(T, \hat{T}) = 10.8$, both of which are approximately of the same order as $\sqrt{N - n} \cdot \text{tol} \approx 8.4$. Note that the ABBA reconstruction \hat{T} visually deviates a lot from T due to the rather high tolerance we have chosen for illustration, but nevertheless, the characteristic up-and-down behavior of T is well represented in \hat{T} , despite the high compression rate of $123/7128 \approx 1.7\%$.

6 Discussion and performance comparison

A Python implementation of ABBA, along with codes to reproduce the figures and performance comparisons in this paper, can be found at <https://github.com/nla-group/ABBA>

When the scaling parameter is $sc1 = 0$ or $sc1 = \infty$, our implementation calls an adaptation of the univariate k -means algorithm from the R package `Ckmeans.1d.dp` written in C++. We use SWIG, the open-source “Simplified Wrapper and Interface Generator”, to call C++ functions from Python. If $sc1 \in (0, \infty)$, we use the k -means algorithm from the Python `sklearn` library (Pedregosa et al. 2011).

ABBA uses the lengths and increments of a polygonal chain on each segment to construct its symbolic time series representation. Symbolic Polynomial (Grabocka et al. 2014) with $d = 1$ and 1d-SAX (Malinowski et al. 2013), on the other hand, use linear regression to fit a polynomial to a window of fixed pre-specified length. As we discussed in Sect. 2, Symbolic Polynomial provides no dimensional reduction and was specifically designed for time series classification problems. Most other SAX variants increase the length of the symbolic representation by enhancing the string with additional characters to capture shapes and trends. It is not clear whether these representations outperform SAX with a reduced width parameter to compensate for the increased string length. A comparison of this would be interesting but is independent of ABBA’s performance and out of the scope of this paper. SMTS (Baydogan and Runger 2015) and aSAX (Pham et al. 2010) use machine learning techniques to discretize their representation. SMTS is primarily designed for multivariate time series and provides no dimensional reduction. EN-SAX (Barnaghi et al. 2012) and aSAX suffer from a loss of the trend information in their compression step.

For these reasons, we focus on profiling the reconstructions errors of the ABBA, SAX (Lin et al. 2007), and 1d-SAX (Malinowski et al. 2013) algorithms, as these are most closely related and easily comparable. Note that none of the representations were primarily designed as compression algorithms. ABBA was designed to be adaptive in both the segment length and alphabet cardinality, whereas SAX and 1d-SAX have many other benefits such as being hashable (Chiu et al. 2003), indexable (Shieh and Keogh 2008), and permitting lower bounding distance measures. Our test set consists of all time series in the UCR Time Series Classification Archive (Dau et al. 2018) with a length of at least 100 data points. There are 128, 978 such time series from a variety of applications. Although the archive is primarily intended for benchmarking time series classification algorithms, our primary focus in this paper is on the approximation performance of the symbolic representations. Our experiment consists of converting each time series $T = [t_0, t_1, \dots, t_N]$ into its symbolic representation $S = [s_1, \dots, s_n]$, and then measuring the distance between the reconstruction $\hat{T} = [\hat{t}_0, \hat{t}_1, \dots, \hat{t}_N]$ and T in the (differenced) Euclidean and DTW norms, respectively.

Recall from Sect. 2 that both SAX and 1d-SAX require a choice for the fixed segment length. In order to provide a fair comparison, we first run the ABBA compression with an initial tolerance $\tau_{01} = 0.05$. This returns n , the number of required pieces to approximate T to this tolerance. If n turns out to be larger than $N/5$, we successively increase the tolerance by 0.05 and rerun until a compression rate of at least 20% is achieved. If a time series cannot be compressed to at least 20% even at the rather

Table 2 Tolerance used for the compression and the number of time series to which it was applied

Tolerance τ_{01}	0.05	0.10	0.15	0.20	0.25	0.30	0.35	0.40	0.45	0.50
nr of time series	75417	9247	7786	5855	2972	2236	1910	1670	2146	2650

crude tolerance of $\tau_{01} = 0.5$, we consider it as too noisy and exclude it from the test. We also exclude all time series which, after ABBA compression, result in fewer than nine pieces: this is necessary because we want to use $k = 9$ symbols for all compared methods. Table 2 shows how many of the 111, 889 remaining time series were compressed at what tolerance. The table gives evidence that most of these time series can be compressed reasonably well while maintaining a rather high accuracy. The average compression rate is 10.3 %.

After the number of pieces n has been specified for a given time series T , we determine the fixed segment length $len = \lfloor (N + 1)/n \rfloor$ to be used in the SAX and 1d-SAX algorithms. We then apply SAX and 1d-SAX to the first $n \cdot len$ points of T . This guarantees that all three algorithms (SAX, 1d-SAX, and ABBA) produce a symbolic representation of with n pieces. If $N + 1$ is not divisible by n , SAX and 1d-SAX are applied to slightly shorter time series than ABBA. The number of symbols used for the digitization is $k = 9$ for all three methods. In the case of 1d-SAX this means that three symbols are used for the mean value, and three symbols are used for the slope on each piece. Each algorithm produces a symbolic representation of length n using an alphabet of cardinality $k = 9$. SAX and 1d-SAX requires the value of w and k for the reconstruction, whereas ABBA requires the $2k$ numbers representing the lengths and increments of each cluster. In total, ABBA requires more storage to represent a time series using a string of length n and alphabet of cardinality k , but is able to represent the whole time series more accurately without truncation.

To visualize the results of our comparison we use performance profiles (Dolan and Moré 2002). Performance profiles allow to compare the relative performance of multiple algorithms over a large set of test problems. Each algorithm is represented by a non-decreasing curve in a θ - p graph. The θ -axis represents a tolerance $\theta \geq 1$ and the p -axis corresponds to a fraction $p \in [0, 1]$. If a curve passes through a point (θ, p) it means that the corresponding algorithm performed within a factor θ of the best observed performance on $100 \cdot p$ % of the test problems. For $\theta = 1$ one can read off on what fraction of all test problems each algorithm was the best performer, while as $\theta \rightarrow \infty$ all curves approach the value $p \rightarrow 1$ (unless an algorithm has failed on a fraction of the test problems, which is not the case here).

In Figs. 9a–10d we present eight performance profiles for the ABBA scaling parameters $s_{c1} = 0$ and $s_{c1} = 1$, respectively, and with four different distance measures: Euclidean and DTW distances and their differenced counterparts, respectively. Figure 9a shows the performance profile for $s_{c1} = 0$, with the distance between T and \hat{T} measured in the Euclidean norm. As expected, SAX consistently outperforms ABBA because the Euclidean distance is very sensitive to horizontal shifts in the time direction, which ABBA has completely ignored due to the $s_{c1} = 0$ parameter. However,

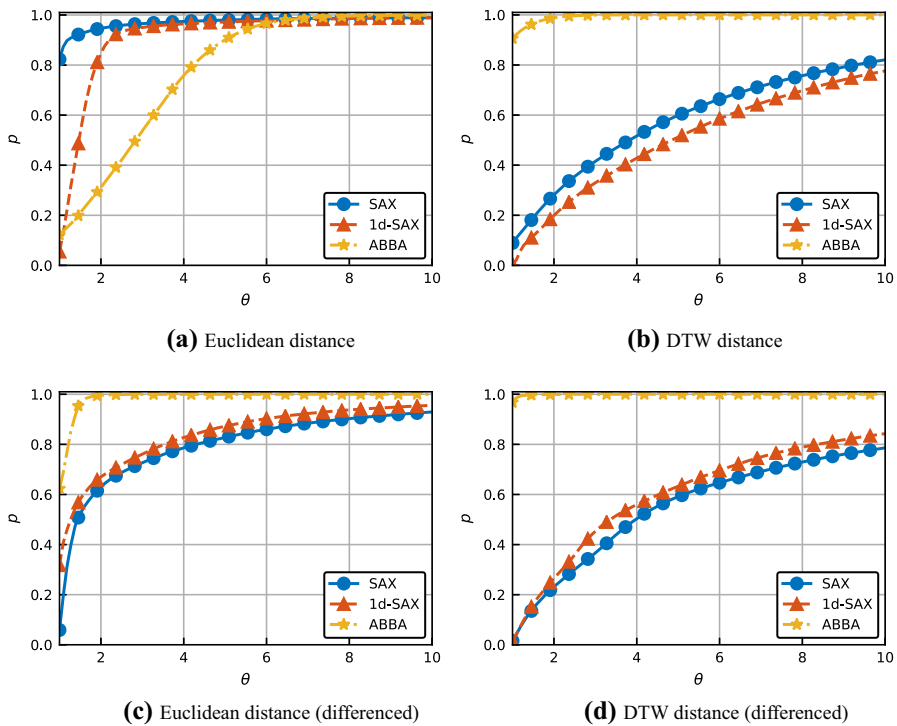


Fig. 9 Performance profiles for the reconstruction errors of SAX, 1d-SAX, and ABBA with scaling parameter $scl = 0$. **a, b** compare ABBA ($scl = 0$) with SAX and 1d-SAX using Euclidean and Dynamic Time Warping distance, respectively. **c, d** compare ABBA ($scl = 0$) with SAX and 1d-SAX using Euclidean and Dynamic Time Warping distance of the differenced time series, respectively

it is somewhat surprising that SAX also outperforms 1d-SAX. When $k = 9$, 1d-SAX allocates 3 symbols for the slopes and 3 symbols for the averages. To represent all possible combinations of slopes and averages, 9 unique symbols are needed. It appears that the use of the slope information in 1d-SAX is detrimental to the approximation accuracy and, if the number of symbols is kept constant, they should better be used to represent the averages rather than the slopes.

The performance changes when we use the DTW distance, thereby allowing for shifts in time. In this case, ABBA outperforms SAX and 1d-SAX significantly; see Fig. 9b. This is because ABBA has been tailored to preserve the up-and-down shape of the time series, at the cost of allowing for small errors in the lengths of the pieces which are easily corrected by time warping. Again SAX with $k = 9$ symbols performs better than 1d-SAX with $k = 9$ symbols.

The performance gain of ABBA becomes even more pronounced when we difference the data before computing the Euclidean and DTW distances; see Fig. 9c, d, respectively. Moreover, we observe that when differencing is used, 1d-SAX performs slightly better than SAX. Computing the Euclidean and DTW distances of the differenced data amounts to comparing the gradients of the time series, rather than their values. The gradient information is better captured when we allocate some symbols

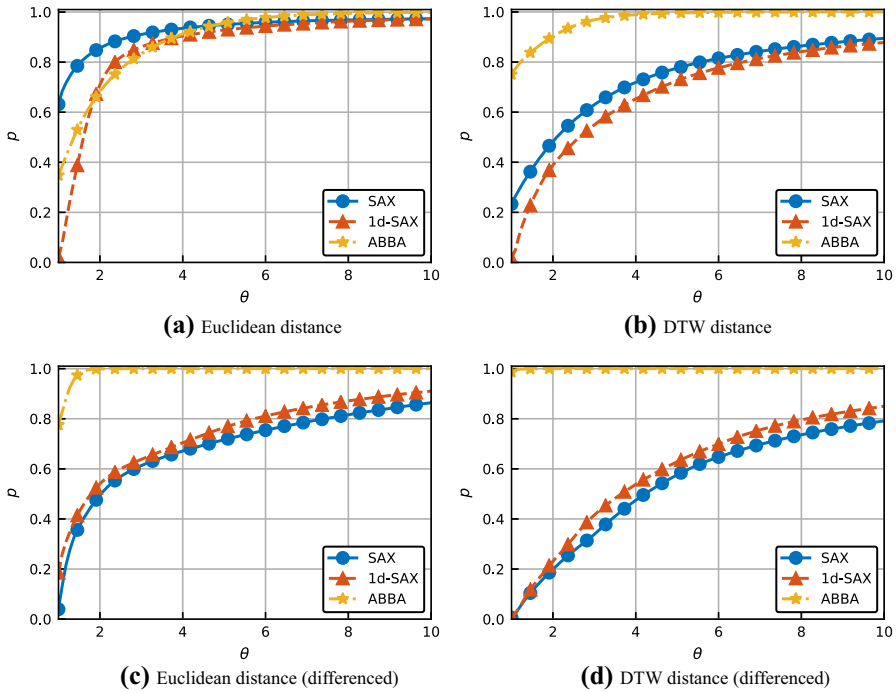


Fig. 10 Performance profiles for the reconstruction errors of SAX, 1d-SAX, and ABBA with scaling parameter $s_{cl} = 1$. **a, b** compare ABBA ($s_{cl} = 1$) with SAX and 1d-SAX using the Euclidean and Dynamic Time Warping distance, respectively. **c, d** compare ABBA ($s_{cl} = 1$) with SAX and 1d-SAX using the Euclidean and Dynamic Time Warping distance of the differenced time series, respectively

for the slope information (as in 1d-SAX) rather than allocating all symbols for the averages (as in SAX). This explains the slight advantage of 1d-SAX over SAX with differenced distance measures.

In the next four tests we set $s_{cl} = 1$, so the ABBA clustering procedure considers both the increments and lengths equally. Figure 10a, b show the resulting performance profiles using the Euclidean and DTW distance measures, respectively. As expected, ABBA becomes more competitive even for the Euclidean distance measure. Computationally, however, this comes at the cost of not being able to use a fast optimal 1d-clustering algorithm. Finally, Fig. 10c, d show the performance profiles for the Euclidean and DTW distance measures on the differenced data, respectively. As in the case $s_{cl} = 0$, differencing helps to improve the performance of ABBA in comparison to SAX and 1d-SAX even further.¹

¹ Visual comparisons of the three algorithms on the first time series in each dataset of the UCR Time Series Classification Archive, as well as codes and CSV files with all numerical results used to produce the performance profiles, can be downloaded from https://github.com/nla-group/ABBA/tree/master/paper/performance_profiles.

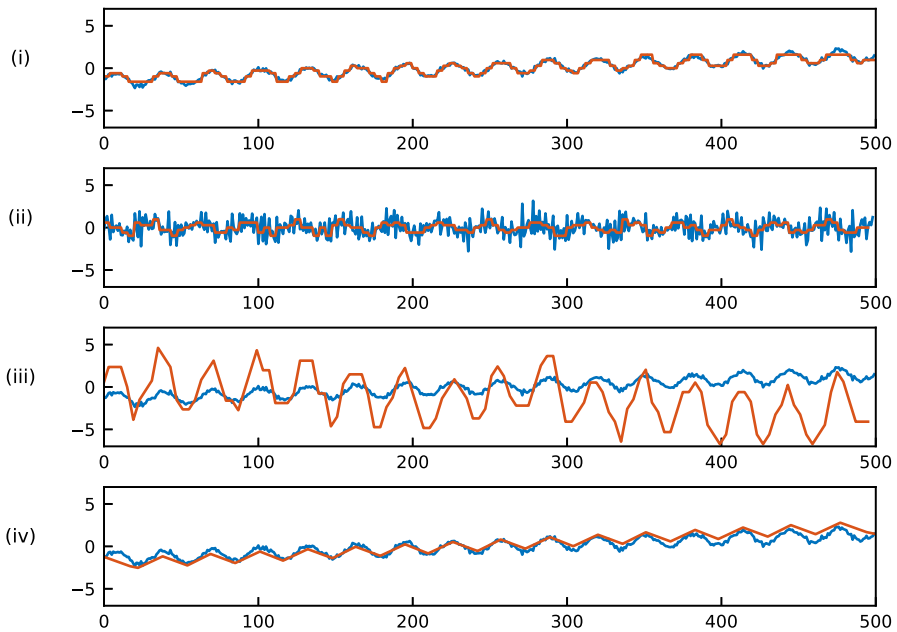


Fig. 11 Comparison of SAX and ABBA on a noisy sine wave with a gradual linear trend. (i) The original time series is shown in blue and the SAX representation is shown in orange. (ii) The differenced version of the original time series is shown in blue and the its SAX representation is shown in orange. (iii) The original time series is given in blue and the cumulative sum of the SAX representation from (ii) is shown in orange. (iv) The original time series is shown in blue and its ABBA representation is shown in orange

7 Further discussion and applications

Section 6 demonstrated that ABBA provides high compression rates while guaranteeing that the time series reconstruction is still close to the original. The high compression is a consequence of the stitching procedure during the compression stage. Section 5 showed how errors are accumulated piece by piece in the stitching process. We believe that this property prevents ABBA from admitting lower bounding distance measures as are available for SAX. SAX's lower bounding measure and indexability make it suitable for applications where multiple time series have to be compared (like time series classification). ABBA, on the other hand, appears best suited for applications where information has to be extracted from a single time series, such as anomaly detection, motif discovery, and trend prediction. As the output of ABBA is simply a string sequence, it can be combined with existing algorithms that previously used, e.g., a SAX representation. Below we discuss various aspects and applications of ABBA. *In-built differencing.* Working with the increments (instead of slopes) allows ABBA to capture linear trends in time series without preprocessing. In Fig. 11 we consider the simple test problem of a sine wave with a gradual linear trend in the presence of noise. After normalization, SAX is able to accurately represent the time series as shown in Fig. 11(i). If we used the symbolic representation for trend prediction, however, the SAX representation would be unsuitable for continuing the linear trend as new symbols

would need to be introduced. Of course, this problem could be overcome by removing the linear trend through differencing the time series. A SAX representation of the differenced time series is shown in Fig. 11(ii). Unfortunately, differencing the noisy time series amplifies the noise. Figure 11(iii) compares the original time series against the reconstructed time series from the SAX representation of the differenced data. As we can see, the increased noise level renders the SAX representation extremely inaccurate. ABBA, on the other hand, does not require any differencing as it works with increments by default. As a consequence, the ABBA reconstruction shown in Fig. 11(iv) stays very close to the original time series, capturing both the gradual linear trend as well as the characteristic up-and-down behavior.

Anomaly detection refers to the problem of finding points or intervals in time series which display surprising or unexpected behavior. Recent literature reviews of existing anomaly detection algorithms are given in Gupta et al. (2013), Atluri et al. (2018). The ABBA representation can be used for anomaly detection in a variety of ways. Trend anomalies can be detected in the digitization procedure via k -means clustering of the lengths and increments. The alphabet is ordered such that 'a' is the most frequent symbol followed by 'b' and so forth. If the k th cluster contains very few elements relative to the other clusters, then this might be considered a trend anomaly.

TARZAN. Keogh et al. (2002) is a popular anomaly detection algorithm with linear time and space complexity (Pelkonen et al. 2015). The algorithm requires two time series, a reference time series R containing normal behavior and the test time series X . Both time series are converted to a symbolic representation and stored in a suffix tree (McCreight 1976). An anomaly score is computed by comparing the frequency of a substring in X to an expected frequency computed from R . SAX can be used for the discretization process in TARZAN and has been shown to outperform other symbolic representations with no dimensional reduction (Lin et al. 2007).

If both symbolic representations are short and X contains a symbol that does not appear in R , then the TARZAN score can suffer through lack of perspective. For example, suppose the expected frequency of the substring 'abc' is 4.2 and 'abc' appears 3 times in X , then the anomaly score is $3 - 4.2 = -1.2$. Suppose the symbol 'd' does not appear in R but 'ada' appears in X . The expected frequency of the substring 'ada' is 0 and 'ada' appears only once, so the anomaly score is $0 - 1 = -1$. This implies that 'abc' is more of an anomaly than 'ada'. This issue can be overcome by dividing the anomaly score by the largest of the expected/actual frequency.

In Figs. 12, 13 we consider a simple experiment comparing SAX, 1d-SAX, and ABBA as discretization procedures for TARZAN with the modified anomaly score.² The reference time series R is a simple sine wave where each period spans 25 time samples. The time series X has a full wave replaced by a flat line of 22 time points. The SAX and 1d-SAX representations use a window length $w = 5$ and $k = 9$ symbols, whereas ABBA uses a tolerance tuned to give a symbolic representation of equal length and k is bounded by 9. The time series R and X and their symbolic reconstructions are shown in Fig. 12. If the length of the anomaly does not align with the window length w , then SAX and 1d-SAX tend to represent the sine wave following the anomaly as a

² A Python implementation of TARZAN which supports the use of SAX, 1d-SAX, and ABBA can be downloaded from <https://github.com/nla-group/TARZAN>.

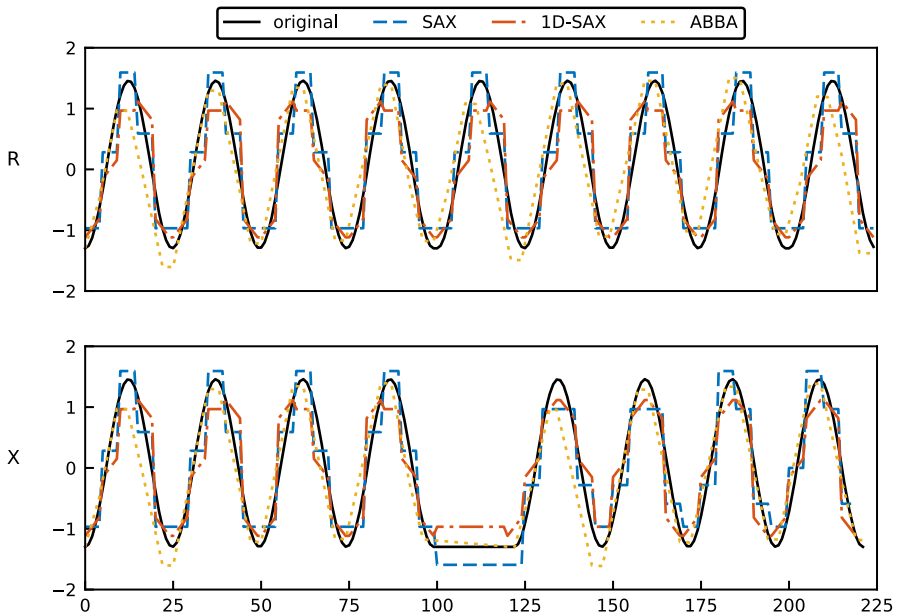


Fig. 12 A visual comparison of the symbolic representations of two time series. Here, R is the reference time series, a simple sine wave, while X is the test time series, a sine wave with a flat region that is three time points shorter than one wave period

different substring. The adapted TARZAN score is required as certain symbols appear in X that do not appear in R . Figure 13 shows the resulting TARZAN anomaly scores. Both SAX and 1d-SAX suffer from the fixed window length, returning high anomaly scores throughout time following the anomaly, whereas TARZAN using ABBA is able to recover almost immediately after the anomaly due to the adaptive segment lengths.

VizTree We finally mention the possibility of representing an ABBA output as a VizTree, a time series pattern discovery and visualization tool based on suffix trees (Lin et al. 2004a, b, 2005). The authors use SAX to discretize the time series before building a suffix tree. Each branch of the suffix tree represents a substring and the thickness of that branch represents the frequency of the substring in the symbolic representation. In principle, SAX pairs well with the visualization as the Gaussian breakpoints should ensure that each symbol appears equally likely. In practice, this is often not the case. One could use ABBA's discretization process instead of SAX by relating the thickness of each line to the frequency of the symbols determined in the clustering procedure. A poor choice of the window length w in the piecewise aggregate approximation in SAX could lead to missing motifs if the distance between is not near a multiple of w . Furthermore, SAX might fail to detect motifs if time warping has occurred, whilst VizTree via ABBA should be able to better capture time-warped motifs as the segment lengths are chosen adaptively. A further exploration of this application will be the subject of future work.

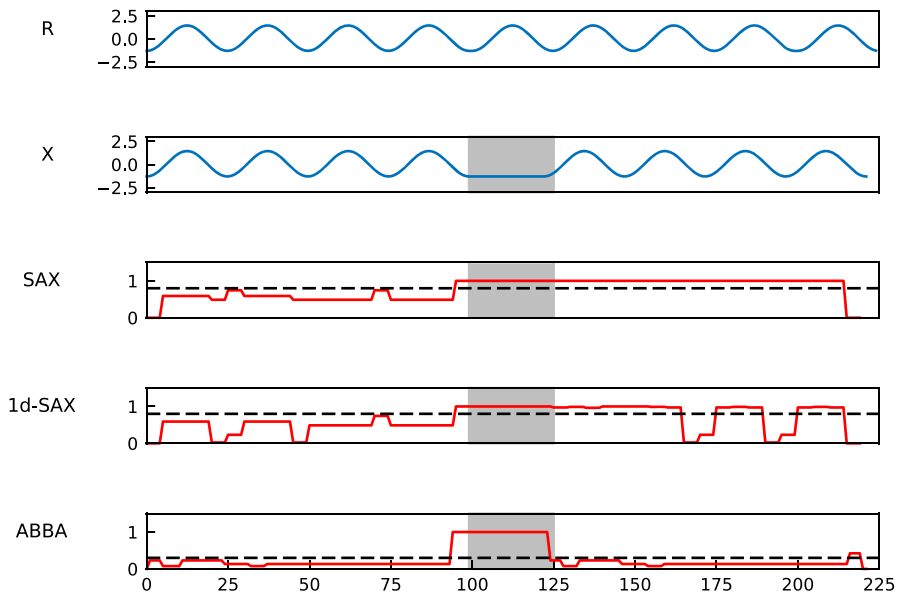


Fig. 13 A comparison of the TARZAN anomaly detection algorithm using the SAX, 1d-SAX, and ABBA representations, respectively. The first time series R is the reference, while the second time series X is to be tested. The final three plots show the adapted TARZAN anomaly scores for the SAX, 1d-SAX, and ABBA representations, respectively. The black dashed lines indicate tolerances that could be used to define the anomalies

8 Conclusions and future work

We introduced ABBA, an adaptive symbolic time series representation which aims to preserve the essential shape of a time series. We have shown that the ABBA representation has favorable approximation properties compared to other popular representations, in particular, when the dynamic time warping distance is used. Furthermore, we demonstrated the use of ABBA in some important data mining applications, including trend prediction and anomaly detection. Future research will be devoted to an online streaming version of ABBA with the necessary adaptations of the the Brownian bridge-based error analysis, as well as a more in-depth study of VizTree visualizations. Our recent work Elsworth and Güttel (2020) explores ABBA's potential for time series forecasting.

Acknowledgements This work was supported by the UK's Engineering and Physical Sciences Research Council (EPSRC) under the grant EP/N509565/1. We thank Sabisu (now AspenTech) and the EPSRC for providing SE with a CASE PhD studentship. SG acknowledges support from the Alan Turing Institute under the EPSRC Grant EP/N510129/1. We thank Timothy D. Butters for his help with C++ and SWIG, and are grateful to Eamonn Keogh and all other contributors to the UCR Time Series Classification Archive. We also thank the three anonymous referees and the editor for their helpful comments which significantly improved the paper.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence,

and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Abanda A, Mori U, Lozano JA (2019) A review on distance based time series classification. *Data Min Knowl Discov* 33(2):378–412. <https://doi.org/10.1007/s10618-018-0596-4>
- Aghabozorgi S, Shirkhorshidi AS, Wah TY (2015) Time-series clustering—a decade review. *Inf Syst* 53:16–38. <https://doi.org/10.1016/j.is.2015.04.007>
- Arthur D, Vassilvitskii S (2006) How slow is the k -means method? In: *Symposium on computational geometry*, ACM, New York, vol 6, pp 1–10
- Atluri G, Karpatne A, Kumar V (2018) Spatio-temporal data mining: a survey of problems and methods. *ACM Comput Surv (CSUR)* 51(4):83
- Barnaghi PM, Bakar AA, Othman ZA (2012) Enhanced symbolic aggregate approximation method for financial time series data representation. In: *6th International conference on new trends in information science, service science and data mining (ISSDM2012)*, IEEE, pp 790–795
- Baydogan MG, Runger G (2015) Learning a symbolic representation for multivariate time series classification. *Data Min Knowl Discov* 29(2):400–422
- Benyahmed Y, Bakar AA, Hamdan AR, Abdullah SMS (2015) A time-weighted average-based PAA representation for time series symbolization. *Int J Adv Soft Comput Appl* 7(3):1–15
- Berndt DJ, Clifford J (1994) Using dynamic time warping to find patterns in time series. *KDD Workshop* 10:359–370
- Bettaiah V, Ranganath HS (2014) An analysis of time series representation methods: data mining applications perspective. In: *Proceedings of the 2014 ACM Southeast Regional Conference*, ACM, pp 16:1–16:6. <https://doi.org/10.1145/2638404.2638475>
- Bondu A, Boullé M, Cornuéjols A (2016) Symbolic representation of time series: a hierarchical co-clustering formalization. In: *International workshop on advanced analysis and learning on temporal data*, Springer, pp 3–16
- Boullé M (2006) MODL: a Bayes optimal discretization method for continuous attributes. *Mach Learn* 65(1):131–165. <https://doi.org/10.1007/s10994-006-8364-x>
- Chiu B, Keogh E, Lonardi S (2003) Probabilistic discovery of time series motifs. In: *Proceedings of the ninth ACM SIGKDD international conference on knowledge discovery and data mining*, ACM, pp 493–498
- Dau HA, Keogh E, Kamgar K, Yeh CCM, Zhu Y, Gharghabi S, Ratanamahatana CA, Yanping, Hu B, Begum N, Bagnall A, Mueen A, Batista G (2018) The UCR time series classification archive. https://www.cs.ucr.edu/~eamonn/time_series_data_2018/. Accessed 14 Mar 2019
- Dolan ED, Moré JJ (2002) Benchmarking optimization software with performance profiles. *Math Program* 91(2, Ser. A):201–213. <https://doi.org/10.1007/s101070100263>
- Elsworth S, Güttel S (2020) Time series forecasting using LSTM networks: a symbolic approach. *Manchester Institute for Mathematical Sciences, The University of Manchester, UK*. [arXiv:2003.05672](https://arxiv.org/abs/2003.05672)
- Esmael B, Arnaout A, Fruhwirth RK, Thonhauser G (2012) Multivariate time series classification by combining trend-based and value-based approximations. In: *International conference on computational science and its applications*, Springer, pp 392–403
- Fu T (2011) A review on time series data mining. *Eng Appl Artif Intell* 24(1):164–181. <https://doi.org/10.1016/j.engappai.2010.09.007>
- Ganz F, Barnaghi P, Carrez F (2013) Information abstraction for heterogeneous real world internet data. *IEEE Sensors J* 13:3793–3805. <https://doi.org/10.1109/JSEN.2013.2271562>

- Grabocka J, Wistuba M, Schmidt-Thieme L (2014) Scalable classification of repetitive time series through frequencies of local polynomials. *IEEE Trans Knowl Data Eng* 27(6):1683–1695
- Grønlund A, Larsen KG, Mathiasen A, Nielsen JS (2017) Fast exact k -means, k -medians and Bregman divergence clustering in 1D. [arXiv:1701.07204](https://arxiv.org/abs/1701.07204)
- Gupta M, Gao J, Aggarwal CC, Han J (2013) Outlier detection for temporal data: a survey. *IEEE Trans Knowl Data Eng* 26(9):2250–2267
- Keogh E, Pazzani MJ (2001) Derivative dynamic time warping. In: Proceedings of the 2001 SIAM international conference on data mining, pp 1–11
- Keogh E, Kasetty S (2003) On the need for time series data mining benchmarks: a survey and empirical demonstration. *Data Min Knowl Discov* 7(4):349–371. <https://doi.org/10.1023/A:1024988512476386>. <https://doi.org/10.1007/s10115-004-0154-9>
- Keogh E, Ratanamahatana CA (2005) Exact indexing of dynamic time warping. *Knowl Inf Syst* 7(3):358–386. <https://doi.org/10.1007/s10115-004-0154-9>
- Keogh E, Chu S, Hart D, Pazzani M (2001) An online algorithm for segmenting time series. In: Proceedings of the 2001 IEEE international conference on data mining, pp 289–296. <https://doi.org/10.1109/ICDM.2001.989531>
- Keogh E, Lonardi S, Chiu B (2002) Finding surprising patterns in a time series database in linear time and space. In: Proceedings of the eighth ACM SIGKDD international conference on knowledge discovery and data mining, ACM, pp 550–556
- Kim JY (2000) Detection of change in persistence of a linear time series. *J Econom* 95(1):97–116. [https://doi.org/10.1016/S0304-4076\(99\)00031-7](https://doi.org/10.1016/S0304-4076(99)00031-7)
- Li G, Zhang L, Yang L (2012) TSX: a novel symbolic representation for financial time series. In: PRICAI 2012: trends in artificial intelligence, Springer, pp 262–273
- Lin J, Keogh E, Lonardi S, Lankford JP, Nystrom DM (2004a) Visually mining and monitoring massive time series. In: Proceedings of the tenth ACM SIGKDD international conference on knowledge discovery and data mining, ACM, pp 460–469
- Lin J, Keogh E, Lonardi S, Lankford JP, Nystrom DM (2004b) Viztree: a tool for visually mining and monitoring massive time series databases. In: Proceedings of the 30th international conference on very large data bases, vol 30 (VLDB '04), VLDB Endowment, pp 1269–1272
- Lin J, Keogh E, Lonardi S (2005) Visualizing and discovering non-trivial patterns in large time series databases. *Inf Vis* 4(2):61–82
- Lin J, Keogh E, Wei L, Lonardi S (2007) Experiencing SAX: a novel symbolic representation of time series. *Data Min Knowl Discov* 15(2):107–144. <https://doi.org/10.1007/s10618-007-0064-z>
- Lkhagva B, Suzuki Y, Kawagoe K (2006) New time series data representation ESAX for financial applications. In: 22nd international conference on data engineering workshops (ICDEW'06), pp x115–x115. <https://doi.org/10.1109/ICDEW.2006.99>
- Luo G, Yi K, Cheng SW, Li Z, Fan W, He C, Mu Y (2015) Piecewise linear approximation of streaming time series data with max-error guarantees. In: 2015 IEEE 31st international conference on data engineering, pp 173–184. <https://doi.org/10.1109/ICDE.2015.7113282>
- Malinowski S, Guyet T, Quiniou R, Tavenard R (2013) 1d-SAX: a novel symbolic representation for time series. In: Advances in intelligent data analysis XII, Springer, pp 273–284
- McCreight EM (1976) A space-economical suffix tree construction algorithm. *J ACM (JACM)* 23(2):262–272
- Mörchen F, Ultsch A (2006) Finding persisting states for knowledge discovery in time series. In: From data and information analysis to knowledge engineering, Springer, pp 278–285
- Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E (2011) Scikit-learn: machine learning in Python. *J Mach Learn Res* 12:2825–2830
- Pelkonen T, Franklin S, Teller J, Cavallaro P, Huang Q, Meza J, Veeraraghavan K (2015) Gorilla: a fast, scalable, in-memory time series database. *Proc VLDB Endow* 8(12):1816–1827
- Pham ND, Le QL, Dang TK (2010) HOT aSAX: a novel adaptive symbolic representation for time series discords discovery. In: Proceedings of the second international conference on intelligent information and database systems: part I, Springer, pp 113–121
- Shieh J, Keogh E (2008) iSAX: indexing and mining terabyte sized time series. In: Proceedings of the 14th ACM sigkdd international conference on knowledge discovery and data mining, ACM, pp 623–631
- Wang H, Song M (2011) Ckmeans.1d.dp: optimal k -means clustering in one dimension by dynamic programming. *R J* 3(2):29–33

Zhang K, Li Y, Chai Y, Huang L (2018) Trend-based symbolic aggregate approximation for time series representation. In: 2018 Chinese control and decision conference, IEEE, pp 2234–2240. <https://doi.org/10.1109/CCDC.2018.8407498>

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.