

OTTO-VON-GUERICKE-UNIVERSITY MAGDEBURG
Faculty of Computer Science



MASTER THESIS

A comparison of Time Series Classification Algorithms based on their ability to learn on diminishing time series

AUTHOR:
ISMAIL, WAHBA

MATRICULATION NUMBER:
217526

EXAMINER AND SUPERVISOR:
PROF. DR. MYRA SPILIOPOULOU
KNOWLEDGE MANAGEMENT AND DISCOVERY LAB
INSTITUTE OF TECHNICAL AND BUSINESS INFORMATION SYSTEMS
OTTO-VON-GUERICKE-UNIVERSITY MAGDEBURG

2ND SUPERVISOR:
NAME
INSTITUTE
UNIVERSITY

day.month.year

Wahba, Ismail:

A comparison of Time Series Classification Algorithms based on their ability to learn
on diminishing time series

Master Thesis, Otto-von-Guericke-University Magdeburg, year.

Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Curabitur eget porta erat. Morbi consectetur est vel gravida pretium. Suspendisse ut dui eu ante cursus gravida non sed sem. Nullam sapien tellus, commodo id velit id, eleifend volutpat quam. Phasellus mauris velit, dapibus finibus elementum vel, pulvinar non tellus. Nunc pellentesque pretium diam, quis maximus dolor faucibus id. Nunc convallis sodales ante, ut ullamcorper est egestas vitae. Nam sit amet enim ultrices, ultrices elit pulvinar, volutpat risus.

Acknowledgement

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Curabitur eget porta erat. Morbi consectetur est vel gravida pretium. Suspendisse ut dui eu ante cursus gravida non sed sem. Nullam sapien tellus, commodo id velit id, eleifend volutpat quam. Phasellus mauris velit, dapibus finibus elementum vel, pulvinar non tellus. Nunc pellentesque pretium diam, quis maximus dolor faucibus id. Nunc convallis sodales ante, ut ullamcorper est egestas vitae. Nam sit amet enim ultrices, ultrices elit pulvinar, volutpat risus.

Table of Contents

List of Figures	i
List of Tables	ii
List of Abbreviations	iii
1 Introduction	1
1.1 Goal of the Thesis	2
1.2 Structure of the Thesis	2
2 Concepts and Terminology	3
2.1 Time Series Data	3
2.1.1 Nature of Time Series Data	3
2.2 Time Series Classification	4
2.3 Early Time Series Classification	4
2.4 notes	5
3 Time Series Classification Algorithms	8
3.1 Whole Time Series Algorithms	8
3.1.1 Nearest Neighbor with ED	8
3.1.2 Nearest Neighbor with DTW	9
3.1.3 Nearest Neighbor with MSM Distance	9
3.1.4 Proximity Forest	10
3.2 Phase Dependent intervals Algorithms	11
3.2.1 Time Series Forest	11
3.3 Phase Independent intervals Algorithms	12
3.3.1 Learned Shapelets	13
3.3.2 Shapelet Transform	13
3.4 Dictionary Based Algorithms	14
3.4.1 WEASEL	14
3.4.2 BOSS	14
3.5 Deep Learning Algorithms	14
3.5.1 Inception Time	14
4 Appendix	15
5 Declaration of Authorship	21

List of Figures

List of Tables

List of Abbreviations

DTW	Dynamic Time Warping
EE	Elastic Ensemble
ERP	Edit Distance with Real Penalty
eTSCA	Early Time Series Classification Algorithms
FS	Fast Shapelets
HM	Harmonic Mean
IG	Information Gain
i.i.d	Independent and Identically Distributed
KNN	K-Nearest Neighbor Algorithm
KNNED	K-Nearest Neighbor using Euclidean Distance
LCSS	Longest Common Subsequence
LOOCV	Leave One Out Cross Validation
LS	Learned Shapelets
MPL	Minimum Prediction Length
MSM	Move-Split-Merge
PF	Proximity Forest
SAX	Symbolic Aggregate Approximation
ST	Shapelet Transform
TSC	Time Series Classification
TSCA	Time Series Classification Algorithms
TSF	Time Series Forest
TWE	Time Warp Edit
WDTW	Weighted Dynamic Time Warping

1 Introduction

Time Series Classification is a field of machine learning that has grabbed the attention of many researchers in the last decade. Time series data exists, by nature, in numerous real scenarios; medical examination records of patients{reference}, signal processing{reference}, weather forecasting{reference} and astronomy{reference} are some of them.

Classification of time series data has been tackled with different objectives; the first is concerned with the accuracy of classification as well as space and time complexity. This objective is referred to simply as Time Series Classification (TSC). While the second objective adds the factor of earliness as a primary goal and is referred to as Early Time Series Classification (eTSC).

Numerous algorithms have been introduced to tackle the problem of Time Series Classification. According to the [3], these algorithms can be divided, based on their technique, into six groups.

Whole time series algorithms{reference} compare two time series, usually by employing an elastic distance measure between all data points of both time series. Phase dependent interval algorithms{reference} operate by extracting informative features from intervals of time series, they are more suitable for long and noisy data than whole time series algorithms. Phase independent interval algorithms{reference} are used when a class can be identified using a single or multiple patterns regardless of when they occur during the time series. Dictionary based algorithms{reference} consider the number of repetitions of patterns as a factor of classification and not just simple occurrence of one. Ensembles{reference} combine the power of different algorithms, either of different or same core technique, then make the final classification decision based on voting. In addition to the previous algorithms, there are also deep learning time series algorithms which build classifiers using generative as well as discriminative models.

On the other hand, Early Time Series Classification algorithms are designed to deal with less data in order to achieve earliness of prediction, but of course this comes with a price of accuracy. Many of the ideas applied in TSC have also been applied in eTSC; including 1-NN with Minimum Prediction Length (MPL){reference}, Phase independent intervals{reference}, generative classifiers{reference} and ensembles{reference}.

Both, Time Series Classification Algorithms (TSCA) and Early Time Series Classification Algorithms (eTSCA), have introduced well performing algorithms in terms of their respective performance measures. Their algorithms have been tested on publicly available archives{reference}; to benchmark their performance on a diverse set of datasets with different characteristics.

According to [3], based on the "No free lunch theorem", no specific algorithm has proven to prevail over all others. This means that different problems with different datasets would require a choice between the algorithms based on how they perform on them, specially for non-public or non-experimented datasets. In this thesis, we tackle this idea; by offering a framework that runs state-of-the-art algorithms on the provided dataset and provides analyses about the performance of each algorithm.

Also due to their different objectives, TSCA and eTSCA have been dealt with as two different families. Which leaves studying the relationship between both algorithm families an open area for research. We study the relationship between TSCA and eTSCA, by extending TSCA to deal with earliness as a main objective and compare how they perform in an early time series classification problem context.

Goals

1.1 Goal of the Thesis

This master thesis had two main goals. The first goal was to create a testbed for comparing different algorithms on a non-public dataset. While the second one was to study the relationship between the two families of algorithms; TSCAs and eTSCAs.

The first goal was motivated by {reference to the great bake-off}, one of the most comprehensive review papers in the time series field. With it's release, Bagnall et. al has set the foundation methodology for accurately benchmarking the performance of TSCAs for the ,at that time, currently existing and for algorithms that will be developed in the future. In their experiment, they have used 85 datasets publicly available from UCR and UEA, the biggest two data archives. Our goal was to offer a testbed, which can be used on private datasets. It runs state of the art algorithms, then provides analysis about their classification performance. The provided analysis can help, based on empirical evidence, choose the best fitting algorithm in accordance with the problem at hand.

As for the second goal, we extended the study of relationship between TSCAs and eTSCAs. Both families offer a wide variety of algorithms, but have different objectives and thus have different approaches in their learning processes. TSCAs focus primarily on the accuracy of the classification. In order to achieve this goal, full utilization of the whole time series data is done to achieve the highest possible accurate results. While eTSCAs objective tries to maximize both accuracy and earliness together, which is hard to attain because of the contradicting nature between both{reference}. This is why eTSCAs try to learn with as least possible data points as possible while maintaining classification accuracy. This study investigated the ability of TSCAs to perform in a simulated early classification context. TSCAs were trained on shortened training data, while keeping record of models' accuracy measure in comparison to a baseline utilizing complete training data points.

1.2 Structure of the Thesis

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Curabitur eget porta erat. Morbi consectetur est vel gravida pretium. Suspendisse ut dui eu ante cursus gravida non sed sem. Nullam sapien tellus, commodo id velit id, eleifend volutpat quam. Phasellus mauris velit, dapibus finibus elementum vel, pulvinar non tellus. Nunc pelentesque pretium diam, quis maximus dolor faucibus id. Nunc convallis sodales ante, ut ullamcorper est egestas vitae. Nam sit amet enim ultrices, ultrices elit pulvinar, volutpat risus.

2 Concepts and Terminology

This Chapter discusses the definitions and background of the topics mentioned in this thesis. We discuss the nature of time series data, the two problems of time series classification (TSC) and early time series classification (eTSC) then present the different techniques encompassed by them.

2.1 Time Series Data

A time series is a finite sequence of ordered observations, either based on time or another aspect [1, 3]. The existence of the time component makes time series an abundant form of data that covers various domains like; medicine, finance, engineering and biology [32]. A time series dataset is a collection of time series instances.

Definition 1 A set T of n time series instances, $T = \{T_1, T_2, \dots, T_n\}$.

Each of the time series instances T_i consists of a sequence of observations.

Definition 2 A time series T_i , of length L is represented as $T_i = [t_{i1}, t_{i2}, \dots, t_{iL}]$.

Time series data can come in different forms. It is important to comprehend what different forms the data can take and what implicit assumptions they convey; to be able to choose the suitable algorithms and tools to deal with it.

The first form is when the observations of instances capture a singular value, this is referred to as univariate time series.

Definition 3 Univariate time series T_i , of length L is represented as $T_i = [t_{i1}, t_{i2}, \dots, t_{iL}]$. With t_{ij} as a real valued number.

While the other form is when multiple measurements are captured by the observations. According to [33], it is essential to differentiate between the two ways multiple time series can be generated; panel data and multivariate time series data.

If more than one variable is being observed during a single experiment, with each variable representing a different measurement; this is called multivariate time series.

Definition 4 Multivariate time series T_i , of length L is represented as $T_i = [t_{i1}, t_{i2}, \dots, t_{iL}]$. With t_{ij} having M dimensions, each is a univariate time series.

While panel data is when the same kind of measurements is collected from independent instances; like different patients or diverse industrial processes.

For panel data, it is possible to assume that the different instances are i.i.d, but this assumption doesn't hold for observation of a single instance. The same goes for multivariate time series, individual univariate observations are assumed to be statistically dependant.

2.1.1 Nature of Time Series Data

Having discussed the dependency assumptions in time the different forms of time series data. It is this dependency that makes time series data challenging for conventional machine learning algorithms, which are used for tabular and cross-sectional data. Tabular and cross-sectional data assume observations to be independent and identically distributed (i.i.d) [33].

If we were to tabularize time series data; convert it into a tabular form by considering each observation as an individual feature. Then it would be possible to apply conventional machine learning algorithms, under the implicit modelling assumption that observations are not ordered. This means that if the order of the features was changed, still the model result will not change. This assumption can work for some problems, but it doesn't have to work for all problems.

2.2 Time Series Classification

Time series classification is a subtype of the general classification problem, which considers the unique property of dependency between adjacent features of instances [9]. The main goal of time series classification is to learn a function f , which given a training dataset $T = \{T_1, T_2, \dots, T_n\}$ of time series instances along with their corresponding class labels $Y = \{y_1, y_2, \dots, y_n\}$ where $y_i \in \{1, 2, \dots, C\}$, can predict class labels for unseen instances [13].

Time series classification has been studied with different objectives, some papers focused on attaining the highest accuracy of classification as the main goal [28, 27, 8, 32, 49, 18], while other papers focused on attaining lower time complexity [43, 3, 53, 41, 48].

In this master thesis, we are more interested in assessing the results in terms of accuracy than time complexity. We define accuracy like [50]; as the percentage of correctly classified instances for a given dataset D , either being a training or testing dataset.

Definition 5 *Accuracy* = $\frac{\text{number of correct classifications}}{|D|}$

2.3 Early Time Series Classification

On another side, early time series classification is also a classification problem which considers the temporal nature of data, but with a slightly different objective and used for different scenarios other than time series classification.

eTSC's main objective is to learn a model which can classify unseen instances as early as possible, while maintaining a competitive accuracy compared to a model that uses full length data or to a user defined threshold [58]. Which is a very challenging objective; due to the, naturally, contradicting nature of earliness and accuracy. In general, the more data is made available for the model to learn the better accuracy it can attain [38, 55, 59, 37]. This is why many eTSC researches consider it as a problem of optimizing multiple objectives.

eTSC is needed in situations in which waiting for more data to arrive can be costly or when making late decisions can cause unfavorable results [36, 40, 30]. This is why eTSC has been applied in various domains like early medical diagnosis [23, 19], avoiding issues in network traffic flow [6], human activity recognition [60, 24] and early prediction of stock crisis [20].

We follow the definition of earliness mentioned by [50]; as the mean number of data points s after which a label is assigned.

Definition 6 *Earliness* = $\frac{\sum_{T_i \in D} \frac{s}{len(T_i)}}{|D|}$

As well as the objective measure, Harmonic mean (HM), mentioned by [19, 50], which includes both accuracy and earliness. For the problem we have, HM is a weighted average between accuracy and earliness.

Definition 7 $F_\beta = (1 + \beta^2) \frac{\text{accuracy}(1 - \text{earliness})}{\beta^2(1 - \text{accuracy}) + \text{earliness}}$

The value of β can be used to give higher importance to one of the aspects over the other, but we use equal weights for both.

2.4 notes

General notes about time series classification problem:

1. Types of Data

(a) Static Data

Data which describe characteristics of the studied instances or properties that won't change with time. These could be the date of birth of a patient, or the species of an animal.

(b) Dynamic Data

i. Several variables from one or more objects observed in a series of time

A. On one object - i time series data

B. On multiple objects - i Panel data

ii. Data Balance

A. Balanced: Observation carried out thoroughly on all objects in a series of time

B. Unbalanced: When several variables of each object can't be full observed in the same timeframe

2. Types of Models

(a) Univariate Model

We Predict an object with some characteristics will be in which group based on a categorical variable

(b) Multivariate Model

Same as univariate but on multiple categorical variables simultaneously

3. Different Time Series Classification Techniques

(a) Similarity-based techniques

(b) Interval-based techniques

(c) Shapelet-based techniques

(d) Dictionary-based techniques

(e) Combination of transformations

1. A time series forest for classification and feature extraction:

Deng, Houtao, et al. "A time series forest for classification and feature extraction." Information Sciences 239 (2013): 142-153.

The paper introduces a new Tree ensemble classifier for time series data called the Time Series Forest (TSF).

It tries to overcome the shortcomings of time-series instance based classifiers; like 1-NN with Euclidean distance and Nearest Neighbor with Dynamic Time Warping (NNDTW) because they provide few insights on the temporal features which are important for distinguishing different time series classes.

It uses simple summary statistics features (mean, std, slope) but outperforms the others.

It introduces a new technique for choosing the best split called Entrance gain (Entropy & distance) which is better than and cheaper than previous techniques.

It has lower computational complexity of $O(M)$ instead of $O(M^2)$ from the previous methods.

Can be extended by using more complex features like wavelets.

It assumes that input time series are of the same length, so it can be extended by using techniques that align time series with different lengths like Dynamic Time Warping (DTW)

2. A Shapelet Transform for Time Series Classification

Lines, Jason, et al. "A shapelet transform for time series classification." Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining. 2012.

The paper tries to improve using shapelets for time series classification

Shapelets are subsequences of a time series that are considered representatives

Shapelets are easily interpretable, compact and classify new instances fastly, allow for the detection of phase-independent shape-based similarity of subsequences.

It proposes a shapelet transformation. Which is separating the process of finding shapelets from the classification step (this is how it is done in the original technique).

This allows for using any classifier.

Shapelet transform, tries to reduce complexity of original algorithm by choosing top K candidate shapelets instead of keeping all of them.

Then the candidate shapelets are used to transform data instances into a number of features, that can be used with any classifier.

It also proposes a new shapelet evaluation method to use with multi class problems (Compare F-Statistic with Information gain)

Can be extended by doing clustering for the extracted shapelets and not using top K , because there were a lot of similar shapelets.

3. Early Prediction on Time Series: A Nearest Neighbor Approach

Xing, Zhengzheng, Jian Pei, and S. Yu Philip. "Early prediction on time series: a nearest neighbor approach." Twenty-First International Joint Conference on Artificial Intelligence. 2009.

The paper introduces a new concept called Minimum Prediction Length (MPL), which allows for Early Classification of Time Series (ECTS) using 1-Nearest Neighbor

ECTS should be able to make earlier predictions using shorter time series than normal 1-NN on Time series data using full-length time series. While retaining accuracy.

It compares to 1-NN with Euclidean distance as distance measure, because it has proved itself to be one of the best techniques in Time Series clustering.

It keeps using shorter subsequences as long as they give the same accuracy of the full time series.

Using 1NN and 1RNN (reverse nearest neighbor) they try to identify the most confident minimum prediction length (MPL) for early prediction

Can be extended for streaming data

4. Faster and More Accurate Classification of Time Series by Exploiting a Novel Dynamic Time Warping Averaging Algorithm

Petitjean, François, et al. "Faster and more accurate classification of time series by exploiting a novel dynamic time warping averaging algorithm." Knowledge and Information Systems 47.1 (2016): 1-26.

The paper tries to extend 1NN with Dynamic Time Warping (DTW)

It uses the Nearest Centroid Classifier (NCC), an algorithm that generalizes Nearest Neighbor by introducing representative prototype (center of mass) for each class. This allows for way cheaper and faster classification $O(1)$ instead of $O(N)$.

NCC in some scenarios offer higher accuracy than 1NN. So NCC is preferred in cases of similar or higher accuracy due to its less resource requirements

The problem with creating centroid using the traditional DTW is that the resulting average can be a value that is not even a representative of any existing instance

Instead the paper uses DTW Barycenter Averaging (DBA), one of the best averaging algorithm for time series. It defines an average sequence and iteratively refines it following an expectation maximization scheme

5. TS-CHIEF: a scalable and accurate forest algorithm for time series classification Shifaz, Ahmed, et al. "Ts-chief: A scalable and accurate forest algorithm for time series classification." *Data Mining and Knowledge Discovery* (2020): 1-34.

A very recent technique

The new state of the art in time series classification

An ensemble classifier that competes with the previous HIVE-COTE and FLAT-COTE ensemble algorithms, but defeats them in time

It starts by using Proximity Forest, dictionary-based and interval-based algorithms to build an ensemble of classification trees. The splits of these trees are a set of time series references, an object would go down the path of the most similar reference.

At each node candidate splits are created, then the best split is selected using weighted Gini index

For classification, After a time series instance is passed down to the leaf nodes of the trees. The final classification of instance is made using a majority vote of K trees.

TS-CHIEF has an overall almost linear complexity in respect to training size

Can be extended for multivariate time series data and variable length datasets

3 Time Series Classification Algorithms

This chapter will introduce different types of TSCA. There are multiple ways to divide TSCAs

3.1 Whole Time Series Algorithms

Whole time series similarity algorithms, also called distance-based algorithms, are methods that compare pairs of time series instances. An unlabeled time series instance is given the class label of the nearest instance in a training data set [28]. There are two main techniques for carrying out the comparison; either by comparing vector representations of the time series instances, or by combining a defined distance function with a classifier, KNN being the most common one [32]. Whole time series algorithms are best suited for problems where the unique features can exist anywhere along the whole time series[3].

One of the simplest forms of whole time series is 1-NN with Euclidean Distance [15], yet it can suprisingly attain high accuracy compared to other distance measures [57]. But Nearest Neighbor with Euclidean Distance (KNNED) is an easy to beat baseline, due to it's sensitivity for distortion and inability to handle time series of unequal lengths [57, 28, 32]. This lead many of the researchers to focus on defining more advanced and elastic distance measures that can compensate for misalignment between time series [1]. The standard and most common baseline classifier utilizing elastic distance measures is 1-NN with Dynamic Time Warping (DTW) [3]. In contrast to the idea that more powerful machine learning algorithms will be able to defeat the simple KNN and an elastic measure, DTW has proved to be a very tough opponent to other algorithms and other elastic distance measures as well [28, 31, 56]. But there were also other distance metrics that have been introduced in literature, these include extensions of DTW on one hand like; Weighted Dynamic Time Warping (WDTW) which penalizes large warpings based on distance [27] and Derivative Dynamic Time Warping (DDTW) [29, 21] which uses derivatives of sequences as features rather than raw data to avoid singularities. On the other hand, Edit Distance with Real Penalty (ERP) [11], Time Warp Edit (TWE) [35], Longest Common Subsequence (LCSS) [12] and Move-Split-Merge (MSM) [52] are all alternatives for distance measures, yet multiple experiments have considered DTW to be relatively unbeatable [3, 1, 8]. To the extend of our knowledge, the most powerful whole time series classifiers are Elastic Ensemble (EE) [31] and Proximity Forest (PF) [34].

3.1.1 Nearest Neighbor with ED

The Euclidean distance is a remarkably simple technique to calculate the distance between time series instances. Given two instances $T_1 = [t_{11}, t_{12}, \dots, t_{1n}]$ and $T_2 = [t_{21}, t_{22}, \dots, t_{2n}]$, the euclidean distance between them can be determined as:

Definition 8 $ED(T_1, T_2) = \sqrt{\sum_{i=1}^n (t_{1i} - t_{2i})^2}$

Euclidean distnace has been preferred to other classifiers due to it's space and time efficiency, but it suffers from two main shortcomings [5, 27, 28]. The first one is that it cannot handle comparisons between time series of different lengths. While the second one is it's sensitivity to minor discrepancies between time series; it would calculate large distance values for small shiftings or misalignments. Although other metrics

have been introduced to overcome the drawbacks of euclidean distance, experimental proof showed that this is only the case for small datasets, but for larger datasets the accuracy of other elastic measures converge with euclidean distance [25, 14, 2].

3.1.2 Nearest Neighbor with DTW

Dynamic Time Warping was a very strong baseline for time series classification for a long time [1, 3]. It was first introduced as a technique to recognize spoken words that can deal with misalignments between time series that Euclidean Distance couldn't handle [53].

To calculate the distance between two time series instances $T_1 = [t1_1, t1_2, \dots, t1_m]$ and $T_2 = [t2_1, t2_2, \dots, t2_m]$; a distance matrix $M(T_1, T_2)$, of size $m \times m$, is calculated for T_1 and T_2 . With $M_{i,j}(t1_i, t2_j)$ representing the distance between $t1_i \in T_1$ and $t2_j \in T_2$. The goal of DTW is to find an optimal path that minimizes the cumulative distance between points of T_1 and T_2 .

A candidate path $P = [p_1, p_2, \dots, p_p]$ is to be found by traversing M . For a path to be valid it must conform to some conditions:

- $p_1 = (t1_1, t2_1)$
- $p_p = (t1_m, t2_m)$
- for all $i < m$:
 - $0 \leq t1_{i+1} - t1_i \leq 1$
 - $0 \leq t2_{i+1} - t2_i \leq 1$

Finding an optimal path under DTW can be computationally expensive with complexity of $O(n^2)$ for a time series of length n [48, 41]. Consequently it is usual to use a constraint with the path; to prevent comparison of points outside a certain window [53]; like the famous Sakoe-Chiba Band [45], Itakura Parallelogram [26] and Ratanamahatana-Keogh Band [43]. Typically DTW can make use of it's warping window to handle distortion in time series, but still it is vulnerable to cases where the difference in length between instances length is larger than the warping window [54].

3.1.3 Nearest Neighbor with MSM Distance

Move-Split-Merge is distance measure that was first introduced in [52]. The main purpose of introducing MSM was to combine certain characteristics within one distance measure. These are; robustness to misalignments between time series instances, being an actual metric unlike other distance measures like DTW and LCSS, assure translation invariance and achieve quadratic run-time [31].

The way MSM works is pretty much like other edit distance methods; it determines the similarity between two instances through the usage of a set of operations to transform one of them to the other. These operations, as the name indicates, are; move, split and merge [3].

Move is the substitution of one single time point of a series with another. Split divides a single time point of a series into two consecutive time points holding the same value as the original time point. Merge can be seen as the opposite of Split, it combines two consecutive time points holding the same value into one time point.

Each of the previously mentioned operations is associated with a cost. The cost of a move is equal to the absolute difference between the old and the new value of the time point. The costs of split and merge are equal and they are set to a constant to satisfy the symmetry property of metricity [52, 53].

3.1.4 Proximity Forest

Proximity forest was developed by [34]. It was introduced as an addition to scalable time series classification, offering a more scalable and accurate classifier than EE [53]. On one side, EE was an accurate classifier being one the state of the art algorithms and the best among distance based algorithms, as it combines 11 NN-algorithms each using a different elastic measure. But on the other hand, EE’s training process was very slow as it scales quadratically with the training size of the data set [31, 3]. This goes back to the leave-one-out-cross-validation (LOOCV) used to optimize the parameters for each used metric [51].

Proximity Forest wanted to achieve two main goals. The first was to offer an adaptable algorithm that can scale with huge data sets consisting of millions of time series instances. Beating EE, by orders of magnitude, and other state of the art algorithms in terms of training and testing run time complexity. While the other goal was to develop a competitive algorithm on the UCR data sets archive without the need to sacrifice accuracy for scalability as is the case with BOSS-VS [34].

Capitalizing on the previous research that has been put in developing specialized time series distance measures and inspired by the existing EE [18, 17]. Proximity forests combine the the eleven elastic distances from EE along with a tree-based algorithms to form an ensemble of decision trees. The reason behind using tree-based algorithms lies in the divide-and-conquer strategy that they adopt, which makes them scalable for large data sets. Also a stochastic process is used for the selection of distance measures and their hyper-parameters, which usually hinders the performance of other algorithms, like KNN, that need to learn the hyper-parameters of the utilized distance measure for each data set before using it [34]. Proximity forests can scale sublinearly with training data set size, but quadratically with the length of the time series [51].

Proximity forests are based on a similar concept as Random Forests [10], another tree-based ensemble, which learns only on a subset of the available features for building tree nodes. This process insinuates in a factor of variability between the trees that form the ensemble but each with a low bias. The collective classification accuracy of the ensemble then tends to provide better results than any if it’s single classifiers [34].

The building unit of a proximity forest is called the proximity tree. A proximity tree and a decision tree are similar on all aspects, but they differ in the tests they apply in internal nodes. A conventional decision tree builds it’s nodes using attributes. When an instance is being tested, it is compared to the value of the attribute and then follows the branch to which it conforms.

Unlike conventional decision trees, that use feature values for their nodes, proximity trees build their nodes using randomly selected exemplars. When an instance to be tested, an elastic distance measure is calculated and then it follows the branch of the nearest exemplar.

An internal node in the tree holds two attributes; *measure* and *branches*. As noted in [34], a measure is function $object \times object \rightarrow \mathbb{R}$. Proximity Forest uses the same 11 distance measures used by EE; Euclidean distance (ED) Dynamic time warping using the full window (DTW); Dynamic time warping with a restricted warping window (DTW-R); Weighted dynamic time warping (WDTW); Derivative dynamic time warping using the full window (DDTW); Derivative dynamic time warping with a restricted warping window (DDTW-R); Weighted derivative dynamic time warping (WDDTW); Longest common subsequence (LCSS); Edit distance with real penalty (ERP); Time

warp edit distance (TWE); and, Move-Split-Merge (MSM). Proximity Forest saves a lot of the computational cost by replacing parameter searches with random sampling [18, 16]. While branches is a vector of the possible branches to follow, each branch holding two attributes; *exemplar* and *subtree*. *exemplar* is a time series instance to which a query instance is compared, and *subtree* refers to the tree an instance should follow in case it is closest to a specific exemplar.

If all time series in a specific node share the same class, then a leaf node is created and the value of the class label is assigned to the *class* attribute of this node. During classification, if a query instance is to reach this node, it is directly labeled with the value of its *class* attribute.

When a query time series is to be classified, it starts at the root node of a proximity tree. The distance between the query and each of the randomly selected exemplars is calculated, using the randomly selected distance measure at the node. Then the query travels down the branch of the nearest exemplar. This process is repeated, passing through the internal nodes of the tree till the query reaches a leaf node, where it is assigned the class label of that node. This whole process is then repeated for all the trees constituting the forest. The final classification of the forest is made by majority voting between its trees.

3.2 Phase Dependent intervals Algorithms

Phase dependent algorithms is a group of algorithms that extract temporal features from intervals of time series. These temporal features help with the interpretability of the model; as they give insights about the temporal characteristics of the data [4], unlike whole time series algorithms that base their decisions solely on the similarities between instances. Another advantage of phase dependent algorithms is that they can also handle distortions and misalignments of time series data [13].

According to [3], phase dependent algorithms are best used for problems where discriminatory information from intervals exist, this would be the case with long timer series instances and which might include areas of noise that can easily deceive classifiers. Like the case with the SmallKitchenAppliances dataset, in which the usage of three classes; a kettle, a microwave and toaster is recorded every 2 minutes for one day. Not only the pattern of usage is discriminatory in such case, but also the time of usage

Typically using interval features requires a two phase process; first by extracting the temporal features and then training a classifier using the extracted features [13]. There are $n(n-1)/2$ possible intervals, for a time series of length n [3]. There is also a wide variety of features, also called literals, to extract for each interval. These cover simple statistical measures as well as local and global temporal features [46, 44, 13]. This introduces one of the main challenges for phase dependent algorithms, that is which intervals to consider for the feature extraction step. Which [44] proposed a solution for by only considering intervals with lengths equal to powers of two [3].

3.2.1 Time Series Forest

Time Series Forest (TSF) is an algorithm that was introduced in 2013 by [13]. They motivated their model with two main criteria; contributing to interpretable time series classification through the use of simple statistical temporal features and reaching this goal by creating an efficient and effective classifier.

TSF considers three types of interval features; mean, standard deviation and slope. If

we were to consider an interval with starting point t_1 and with ending point t_2 . Let v_i be the value at a specific point t_i . Then the three features can be denoted as:

Definition 9

$$mean(t_1, t_2) = \frac{\sum_{i=t_1}^{t_2} v_i}{t_2 - t_1 + 1}$$

Definition 10

$$std(t_1, t_2) = \begin{cases} \frac{\sum_{i=t_1}^{t_2} (v_i - mean(t_1, t_2))^2}{t_2 - t_1} & \text{if } t_1 < t_2 \\ 0 & \text{if } t_1 = t_2 \end{cases}$$

Definition 11

$$slope(t_1, t_2) = \begin{cases} m & \text{if } t_1 < t_2 \\ 0 & \text{if } t_1 = t_2 \end{cases}$$

Where m denotes the slope of the least squares regression line for the training dataset. For building the trees, TSF introduced a new splitting criteria at the tree nodes, which they called the Entrance. A combination of Entropy and distance; to break the ties between features of equal entropy gain by preferring splits that have the furthest distance to the nearest instance. They also use a specific number of evaluation points rather than checking all split points for the highest information gain. In their experiment [3] found these two criteria to have negative effect on accuracy.

As mentioned earlier the feature space for creating interval features is huge, TSF adopts the same random sampling technique that Random Forests use reducing the feature space from $O(M^2)$ to only $O(M)$, by considering only $O(\sqrt{M})$ random interval sizes and $O(\sqrt{M})$ random starting points at each tree node [13]. The final classification of a testing instance is done using majority voting of all time series trees created.

3.3 Phase Independent intervals Algorithms

Phase independent shapelets, or just shapelets as less formally known, are subseries which are ultimately distinctive of classes regardless of their place on the time series [48, 3]. They were first introduced in [61] as an alternative for KNN approaches; to overcome their shortcomings.

Shapelets reduce the space and time complexity needed by KNN, because they are formed from subsequences which are shorter than the original time series. Needing only one shapelet at classification time, they form a compressed format of the classification problem [8, 61, 39]. While KNN classify based on comparison to other instances, shapelets provide insight about the unique features of classes and thus more interpretable results of how the classification was carried out. Finally, shapelets are best suited for problems where a certain pattern can differentiate instances which is harder to detect when comparing whole series [3, 9].

The original shapelet algorithm enumerated all possible shapelets and embedded the best ones, based on information gain assessment, in a decision tree. Together with a calculated distance threshold, the shapelets and the threshold are used together as splitting criteria [32, 47]. There have been many attempts to speed up the process of shapelets discovery, by determining good shapelets in a faster manner. Two of them are; Fast Shapelets (FS) [42] and Learned Shapelets (LS) [22]. FS applied discretization through Symbolic Aggregate Approximation (SAX) to reduce the length of time

series, while LS tried to learn the shapelets [51]. Later on the idea of transforming time series data to an alternative space was adopted in [25], the transformed data consists of distances to the best k shapelets, then classification is done using an ensemble of eight classifiers.

3.3.1 Learned Shapelets

Learned Shapelets (LS) was proposed in [22] as a new prospective for approaching time series shapelets. Instead of searching for shapelets through enumeration of all candidates, LS learns K near-to-optimal shapelets that can linearly separate instances through a stochastic gradient objective [32, 7]. The found shapelets need not to be a subsequence of one of the training examples [3, 48].

LS follows a two steps technique. In the beginning LS looks for a set of shapelets from the training data set using two parameters; L controls the length of shapelets searched, while R controls the scaling of subsequences. Then these shapelets are clustered using a K-Means algorithm and instances are represented in a new K -dimensional format where the values of the features represent the minimum distance between the instance and one of the shapelets.

For the second step, using the new features representation, LS can start learning class probabilities for instances by considering a logistic regression model for each of the classes and optimizing a regularized logistic loss function. The regularized loss function updates the shapelets and the weights of features. This process keeps iteratively going untill either the model converges or the maximum number of iterations is reached.[7]. In summary, the main objective of the algorithm is to learn collectively the optimal shapelets and the weights linear hyper-plane that minimizes the objective function [3, 22].

3.3.2 Shapelet Transform

The first Shapelet Transform (TS) was introduced in [25]. While the original algorithm embeded shapelets discovery in decision trees and assessed candidates through enumeration and the use Information Gain (IG) at each node, TS proposed a different way that saved repeating the brute force multiple times [7]. TS segregated the shapelets discovery process from the classifier. This segregation opened the door for choosing classifiers freely and considering more accurate classifiers than decision trees [3, 31]. Also [25] experimented with other shapelet assessment metrics like Kruskal-Wallis, F-stat and Mood's median to find out that F-stat attained higher accuracies than the other three and than IG [7].

TS follows a three step procedure. In the beginning, a data transformation phase is carried out by utilizing a single-scan algorithm and extracting K best shapelets from the training data set where K represents a cutoff threshold for the maximum number of shapelets to extract without affecting the quality of the shapelets extracted. Then a reduction process is done by clustering the shapelets together untill they reach a user defined number. Finally, The clustered shapelets are then used to transform the original dataset, by representing instances in terms of their disatances to each one of the extracted shapelets. They experimented with different classifiers other than decision trees, these are; C4.5 tree, 1-NN, naive Bayes, Bayesian network, Random Forest, Rotation Forest, and support vector machine, for which decision trees proved to be the worst among all, while support vector machine proved to be the best [25]. TS was then extended again by [9], the intuition behind it was that the previously

used assessment technique couldn't handle multi-class problems [9]. Instead of assessing shapelets that discriminate between all classes, they accommodated a one-vs-all technique so that shapelets are assessed on their ability to separate one class to all other classes. They also introduced a balancing technique to represent each of the classes with the same number of shapelets [3]. For the classification, a combination of tree based, kernel based and probabilistic classifiers were used in an ensemble on the transformed data set [51, 32]. Each of the classifiers was given a weight based on its training accuracy and the final classification used weighted voting [9]. Although ST has proved to be a competent accurate classifier, it suffers from high training-time complexity [51].

3.4 Dictionary Based Algorithms

3.4.1 WEASEL

3.4.2 BOSS

3.5 Deep Learning Algorithms

3.5.1 Inception Time

4 Appendix

References

- [1] Amaia Abanda, Usue Mori, and Jose A Lozano. “A review on distance based time series classification”. In: *Data Mining and Knowledge Discovery* 33.2 (2019), pp. 378–412.
- [2] Anthony Bagnall et al. “Transformation based ensembles for time series classification”. In: *Proceedings of the 2012 SIAM international conference on data mining*. SIAM. 2012, pp. 307–318.
- [3] Anthony Bagnall et al. “The great time series classification bake off: a review and experimental evaluation of recent algorithmic advances”. In: *Data Mining and Knowledge Discovery* 31.3 (2017), pp. 606–660.
- [4] Mustafa Gokce Baydogan and George Runger. “Time series representation and similarity based on local autopatterns”. In: *Data Mining and Knowledge Discovery* 30.2 (2016), pp. 476–509.
- [5] Mustafa Gokce Baydogan, George Runger, and Eugene Tuv. “A bag-of-features framework to classify time series”. In: *IEEE transactions on pattern analysis and machine intelligence* 35.11 (2013), pp. 2796–2802.
- [6] Laurent Bernaille et al. “Traffic classification on the fly”. In: *ACM SIGCOMM Computer Communication Review* 36.2 (2006), pp. 23–26.
- [7] Aaron Bostrom. “Shapelet Transforms for Univariate and Multivariate Time Series Classification”. PhD thesis. University of East Anglia, 2018.
- [8] Aaron Bostrom and Anthony Bagnall. “A shapelet transform for multivariate time series classification”. In: *arXiv preprint arXiv:1712.06428* (2017).
- [9] Aaron Bostrom and Anthony Bagnall. “Binary Shapelet Transform for Multiclass Time Series Classification”. In: *Transactions on Large-Scale Data- and Knowledge-Centered Systems XXXII: Special Issue on Big Data Analytics and Knowledge Discovery*. Ed. by Abdelkader Hameurlain et al. Berlin, Heidelberg: Springer Berlin Heidelberg, 2017, pp. 24–46. ISBN: 978-3-662-55608-5. DOI: 10.1007/978-3-662-55608-5_2. URL: https://doi.org/10.1007/978-3-662-55608-5_2.
- [10] Leo Breiman. “Random forests”. In: *Machine learning* 45.1 (2001), pp. 5–32.
- [11] Lei Chen and Raymond Ng. “On the marriage of lp-norms and edit distance”. In: *Proceedings of the Thirtieth international conference on Very large data bases-Volume 30*. 2004, pp. 792–803.
- [12] Gautam Das, Dimitrios Gunopulos, and Heikki Mannila. “Finding similar time series”. In: *European Symposium on Principles of Data Mining and Knowledge Discovery*. Springer. 1997, pp. 88–100.
- [13] Houtao Deng et al. “A time series forest for classification and feature extraction”. In: *Information Sciences* 239 (2013), pp. 142–153.

- [14] Hui Ding et al. “Querying and mining of time series data: experimental comparison of representations and distance measures”. In: *Proceedings of the VLDB Endowment* 1.2 (2008), pp. 1542–1552.
- [15] Christos Faloutsos, Mudumbai Ranganathan, and Yannis Manolopoulos. “Fast subsequence matching in time-series databases”. In: *Acm Sigmod Record* 23.2 (1994), pp. 419–429.
- [16] Hassan Ismail Fawaz et al. “Deep learning for time series classification: a review”. In: *Data Mining and Knowledge Discovery* 33.4 (2019), pp. 917–963.
- [17] Hassan Ismail Fawaz et al. “Deep neural network ensembles for time series classification”. In: *2019 International Joint Conference on Neural Networks (IJCNN)*. IEEE. 2019, pp. 1–6.
- [18] Hassan Ismail Fawaz et al. “Inceptiontime: Finding alexnet for time series classification”. In: *Data Mining and Knowledge Discovery* 34.6 (2020), pp. 1936–1962.
- [19] Mohamed F Ghalwash and Zoran Obradovic. “Early classification of multivariate temporal observations by extraction of interpretable shapelets”. In: *BMC bioinformatics* 13.1 (2012), p. 195.
- [20] Mohamed F Ghalwash, Vladan Radosavljevic, and Zoran Obradovic. “Utilizing temporal patterns for estimating uncertainty in interpretable early decision making”. In: *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2014, pp. 402–411.
- [21] Tomasz Górecki and Maciej Luczak. “Using derivatives in time series classification”. In: *Data Mining and Knowledge Discovery* 26.2 (2013), pp. 310–331.
- [22] Josif Grabocka et al. “Learning time-series shapelets”. In: *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2014, pp. 392–401.
- [23] M Pamela Griffin and J Randall Moorman. “Toward the early diagnosis of neonatal sepsis and sepsis-like illness using novel heart rate analysis”. In: *Pediatrics* 107.1 (2001), pp. 97–104.
- [24] Ashish Gupta et al. “A Fault-Tolerant Early Classification Approach for Human Activities using Multivariate Time Series”. In: *IEEE Transactions on Mobile Computing* (2020).
- [25] Jon Hills et al. “Classification of time series by shapelet transformation”. In: *Data Mining and Knowledge Discovery* 28.4 (2014), pp. 851–881.
- [26] Fumitada Itakura. “Minimum prediction residual principle applied to speech recognition”. In: *IEEE Transactions on acoustics, speech, and signal processing* 23.1 (1975), pp. 67–72.
- [27] Young-Seon Jeong, Myong K Jeong, and Olufemi A Omitaomu. “Weighted dynamic time warping for time series classification”. In: *Pattern recognition* 44.9 (2011), pp. 2231–2240.

- [28] Rohit J Kate. “Using dynamic time warping distances as features for improved time series classification”. In: *Data Mining and Knowledge Discovery* 30.2 (2016), pp. 283–312.
- [29] Eamonn J Keogh and Michael J Pazzani. “Derivative dynamic time warping”. In: *Proceedings of the 2001 SIAM international conference on data mining*. SIAM. 2001, pp. 1–11.
- [30] Yu-Feng Lin et al. “Reliable early classification on multivariate time series with numerical and categorical attributes”. In: *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer. 2015, pp. 199–211.
- [31] Jason Lines and Anthony Bagnall. “Time series classification with ensembles of elastic distance measures”. In: *Data Mining and Knowledge Discovery* 29.3 (2015), pp. 565–592.
- [32] Jason Lines, Sarah Taylor, and Anthony Bagnall. “Time series classification with HIVE-COTE: The hierarchical vote collective of transformation-based ensembles”. In: *ACM Transactions on Knowledge Discovery from Data* 12.5 (2018).
- [33] Markus Löning et al. “sktime: A unified interface for machine learning with time series”. In: *arXiv preprint arXiv:1909.07872* (2019).
- [34] Benjamin Lucas et al. “Proximity forest: an effective and scalable distance-based classifier for time series”. In: *Data Mining and Knowledge Discovery* 33.3 (2019), pp. 607–635.
- [35] Pierre-François Marteau. “Time warp edit distance with stiffness adjustment for time series matching”. In: *IEEE transactions on pattern analysis and machine intelligence* 31.2 (2008), pp. 306–318.
- [36] Usue Mori et al. “Early classification of time series by simultaneously optimizing the accuracy and earliness”. In: *IEEE transactions on neural networks and learning systems* 29.10 (2017), pp. 4569–4578.
- [37] Usue Mori et al. “Reliable early classification of time series based on discriminating the classes over time”. In: *Data mining and knowledge discovery* 31.1 (2017), pp. 233–263.
- [38] Usue Mori et al. “Early classification of time series using multi-objective optimization techniques”. In: *Information Sciences* 492 (2019), pp. 204–218.
- [39] Abdullah Mueen, Eamonn Keogh, and Neal Young. “Logical-shapelets: an expressive primitive for time series classification”. In: *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2011, pp. 1154–1162.
- [40] Nathan Parrish et al. “Classifying with confidence from incomplete information”. In: *The Journal of Machine Learning Research* 14.1 (2013), pp. 3561–3589.

- [41] François Petitjean et al. “Faster and more accurate classification of time series by exploiting a novel dynamic time warping averaging algorithm”. In: *Knowledge and Information Systems* 47.1 (2016), pp. 1–26.
- [42] Thanawin Rakthanmanon and Eamonn Keogh. “Fast shapelets: A scalable algorithm for discovering time series shapelets”. In: *proceedings of the 2013 SIAM International Conference on Data Mining*. SIAM. 2013, pp. 668–676.
- [43] Chotirat Ann Ratanamahatana and Eamonn Keogh. “Making time-series classification more accurate using learned constraints”. In: *Proceedings of the 2004 SIAM international conference on data mining*. SIAM. 2004, pp. 11–22.
- [44] Juan Jose Rodriguez and Carlos J Alonso. “Support vector machines of interval-based features for time series classification”. In: *International Conference on Innovative Techniques and Applications of Artificial Intelligence*. Springer. 2004, pp. 244–257.
- [45] Hiroaki Sakoe and Seibi Chiba. “Dynamic programming algorithm optimization for spoken word recognition”. In: *IEEE transactions on acoustics, speech, and signal processing* 26.1 (1978), pp. 43–49.
- [46] Tiago Santos and Roman Kern. “A Literature Survey of Early Time Series Classification and Deep Learning.” In: *Sami@ iknow*. 2016.
- [47] Patrick Schäfer. “The BOSS is concerned with time series classification in the presence of noise”. In: *Data Mining and Knowledge Discovery* 29.6 (2015), pp. 1505–1530.
- [48] Patrick Schäfer and Ulf Leser. “Fast and accurate time series classification with weasel”. In: *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*. 2017, pp. 637–646.
- [49] Patrick Schäfer and Ulf Leser. “Multivariate time series classification with WEASEL+ MUSE”. In: *arXiv preprint arXiv:1711.11343* (2017).
- [50] Patrick Schäfer and Ulf Leser. “TEASER: early and accurate time series classification”. In: *Data Mining and Knowledge Discovery* 34.5 (2020), pp. 1336–1362.
- [51] Ahmed Shifaz et al. “Ts-chief: A scalable and accurate forest algorithm for time series classification”. In: *Data Mining and Knowledge Discovery* (2020), pp. 1–34.
- [52] Alexandra Stefan, Vassilis Athitsos, and Gautam Das. “The move-split-merge metric for time series”. In: *IEEE transactions on Knowledge and Data Engineering* 25.6 (2012), pp. 1425–1438.
- [53] Chang Wei Tan, François Petitjean, and Geoffrey I Webb. “FastEE: Fast Ensembles of Elastic Distances for time series classification”. In: *Data Mining and Knowledge Discovery* 34.1 (2020), pp. 231–272.
- [54] Chang Wei Tan et al. “Time series classification for varying length series”. In: *arXiv preprint arXiv:1910.04341* (2019).

- [55] Romain Tavenard and Simon Malinowski. “Cost-aware early classification of time series”. In: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer. 2016, pp. 632–647.
- [56] Xiaoyue Wang et al. “Experimental comparison of representation methods and distance measures for time series data”. In: *Data Mining and Knowledge Discovery* 26.2 (2013), pp. 275–309.
- [57] Zhengzheng Xing, Jian Pei, and Eamonn Keogh. “A brief survey on sequence classification”. In: *ACM Sigkdd Explorations Newsletter* 12.1 (2010), pp. 40–48.
- [58] Zhengzheng Xing, Jian Pei, and S Yu Philip. “Early prediction on time series: a nearest neighbor approach”. In: *Twenty-First International Joint Conference on Artificial Intelligence*. Citeseer. 2009.
- [59] Zhengzheng Xing, Jian Pei, and S Yu Philip. “Early classification on time series”. In: *Knowledge and information systems* 31.1 (2012), pp. 105–127.
- [60] Omolbanin Yazdanbakhsh and Scott Dick. “Multivariate Time Series Classification using Dilated Convolutional Neural Network”. In: *arXiv preprint arXiv:1905.01697* (2019).
- [61] Lexiang Ye and Eamonn Keogh. “Time series shapelets: a new primitive for data mining”. In: *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2009, pp. 947–956.

5 Declaration of Authorship

I hereby declare that I have written this thesis "TITLE TITLE TITLE" without any help from others and without the use of documents and aids other than those stated above. Furthermore, I have mentioned all used sources and have cited them correctly according to the citation rules. Moreover, I confirm that the paper at hand was not submitted in this or similar form at another examination office, nor has it been published before.

Magdeburg, DATE, SIGNATURE