

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

<b>List of Figures</b>	<b>i</b>
<b>List of Tables</b>	<b>ii</b>
<b>List of Abbreviations</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Goal of the Thesis . . . . .	2
1.2 Structure of the Thesis . . . . .	2
<b>2 Concepts and Terminology</b>	<b>3</b>
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
<b>3 Time Series Classification Algorithms</b>	<b>6</b>
3.1 Whole Time Series Algorithms . . . . .	6
3.1.1 Nearest Neighbor with ED . . . . .	6
3.1.2 Nearest Neighbor with DTW . . . . .	7
3.1.3 Nearest Neighbor with MSM Distance . . . . .	7
3.1.4 Proximity Forest . . . . .	8
3.2 Phase Dependent intervals Algorithms . . . . .	9
3.2.1 Time Series Forest . . . . .	9
3.3 Phase Independent intervals Algorithms . . . . .	10
3.3.1 Learned Shapelets . . . . .	11
3.3.2 Shapelet Transform . . . . .	11
3.4 Dictionary Based Algorithms . . . . .	12
3.4.1 WEASEL . . . . .	12
3.4.2 BOSS . . . . .	14
3.5 Deep Learning Algorithms . . . . .	15
3.5.1 Inception Time . . . . .	17
<b>4 Appendix</b>	<b>19</b>
<b>5 Declaration of Authorship</b>	<b>27</b>

## **List of Figures**

## **List of Tables**



## List of Abbreviations

BOSS	Bag of SFA Symbols
CNN	Convolutional Neural Networks
DFT	Discrete Fourier Transformation
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
MCB	Multiple Coefficient Binning
MPL	Minimum Prediction Length
MSM	Move-Split-Merge
PF	Proximity Forest
SAX	Symbolic Aggregate Approximation
SFA	Symbolic Fourier 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 [5], 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 [5], 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, 5]. The existence of the time component makes time series an abundant form of data that covers various domains like; medicine, finance, engineering and biology [42]. 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  $T_i = [t_1, t_2, \dots, t_L]$ .

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  $T_i = [t_1, t_2, \dots, t_L]$ . With  $t_j$  as a real valued number.

While the other form is when multiple measurements are captured by the observations. According to [43], 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  $T_i = [t_1, t_2, \dots, t_L]$ . With  $t_j$  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) [43].

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 [11]. 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 [16].

Time series classification has been studied with different objectives, some papers focused on attaining the highest accuracy of classification as the main goal [32, 31, 10, 42, 64, 21], while other papers focused on attaining lower time complexity [55, 5, 71, 53, 63].

In this master thesis, we are more interested in assessing the results in terms of accuracy than time complexity. We define accuracy like [65]; 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 [78]. 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 [50, 74, 79, 49]. 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 [48, 52, 38]. This is why eTSC has been applied in various domains like early medical diagnosis [26, 22], avoiding issues in network traffic flow [8], human activity recognition [80, 27] and early prediction of stock crisis [23].

We follow the definition of earliness mentioned by [65]; as the mean number of data points  $s$  after which a label is assigned.

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

As well as the objective measure, Harmonic mean (HM), mentioned by [22, 65], 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  $\beta$  can be used to give higher importance to one of the aspects over the other, but we use equal weights for both.

### 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 [32]. 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 [42]. Whole time series algorithms are best suited for problems where the unique features can exist anywhere along the whole time series[5].

One of the simplest forms of whole time series is 1-NN with Euclidean Distance [18], yet it can suprisingly attain high accuracy compared to other distance measures [77]. 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 [77, 32, 42]. 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) [5]. 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 [32, 41, 75]. 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 [31] and Derivative Dynamic Time Warping (DDTW) [33, 24] which uses derivatives of sequences as features rather than raw data to avoid singularities. On the other hand, Edit Distance with Real Penalty (ERP) [13], Time Warp Edit (TWE) [46], Longest Common Subsequence (LCSS) [15] and Move-Split-Merge (MSM) [68] are all alternatives for distance measures, yet multiple experiments have considered DTW to be relatively unbeatable [5, 1, 10]. To the extend of our knowledge, the most powerful whole time series classifiers are Elastic Ensemble (EE) [41] and Proximity Forest (PF) [45].

##### 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 [7, 31, 32]. 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 [29, 17, 3].

### 3.1.2 Nearest Neighbor with DTW

Dynamic Time Warping was a very strong baseline for time series classification for a long time [1, 5]. 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 [71].

To calculate the distance between two time series instances  $T_1 = [t1_1, t1_2, \dots, t1_n]$  and  $T_2 = [t2_1, t2_2, \dots, t2_n]$ ; a distance matrix  $M(T_1, T_2)$ , of size  $n \times n$ , 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:

- $p_1 = (t1_1, t2_1)$
- $p_p = (t1_n, t2_n)$
- 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  $n$  [63, 53]. Consequently it is usual to use a constraint with the path; to prevent comparison of points outside a certain window [71]; like the famous Sakoe-Chiba Band [58], Itakura Parallelogram [30] and Ratanamahatana-Keogh Band [55]. 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 [72].

### 3.1.3 Nearest Neighbor with MSM Distance

Move-Split-Merge is distance measure that was first introduced in [68]. 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 [41].

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 [5].

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 [68, 71].

### 3.1.4 Proximity Forest

Proximity forest was developed by [45]. It was introduced as an addition to scalable time series classification, offering a more scalable and accurate classifier than EE [71]. 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 [41, 5]. This goes back to the leave-one-out-cross-validation (LOOCV) used to optimize the parameters for each used metric [67].

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 [45].

Capitalizing on the previous research that has been put in developing specialized time series distance measures and inspired by the existing EE [21, 20]. 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 [45]. Proximity forests can scale sublinearly with training data set size, but quadratically with the length of the time series [67].

Proximity forests are based on a similar concept as Random Forests [12], 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 [45].

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 [45], 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 [21, 19]. 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 [6], 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 [16].

According to [5], 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 [16]. There are  $n(n-1)/2$  possible intervals, for a time series of length  $n$  [5]. 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 [59, 56, 16]. This introduces one of the main challenges for phase dependent algorithms, that is which intervals to consider for the feature extraction step. Which [56] proposed a solution for by only considering intervals with lengths equal to powers of two [5].

#### 3.2.1 Time Series Forest

Time Series Forest (TSF) is an algorithm that was introduced in 2013 by [16]. 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 [5] 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 [16]. 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 [63, 5]. They were first introduced in [81] 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 [10, 81, 51]. 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 [5, 11].

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 [42, 61]. 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) [54] and Learned Shapelets (LS) [25]. FS applied discretization through Symbolic Aggregate Approximation (SAX) to reduce the length of time

series, while LS tried to learn the shapelets [67]. Later on the idea of transforming time series data to an alternative space was adopted in [29], 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 [25] 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 [42, 9]. The found shapelets need not to be a subsequence of one of the training examples [5, 63].

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.[9]. 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 [5, 25].

### 3.3.2 Shapelet Transform

The first Shapelet Transform (ST) was introduced in [29]. While the original algorithm embeded shapelets discovery in decision trees and assessed candidates through enumeration and the use Information Gain (IG) at each node, ST proposed a different way that saved repeating the brute force multiple times [9]. ST 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 [5, 41]. Also [29] 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 [9].

ST 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 [29]. ST was then extended again by [11], the intuition behind it was that the previously

used assessment technique couldn't handle multi-class problems [11]. 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 [5]. For the classification, a combination of tree based, kernel based and probabilistic classifiers were used in an ensemble on the transformed data set [67, 42]. Each of the classifiers was given a weight based on its training accuracy and the final classification used weighted voting [11]. Although ST has proved to be a competent accurate classifier, it suffers from high training-time complexity [67].

### 3.4 Dictionary Based Algorithms

#### 3.4.1 WEASEL

After the success that BOSS achieved, it became one of the baseline classifiers of TSC. Other research has either included it as a building component of their ensembles [42, 4], or have used it, as a baseline for comparison of their algorithm's performance [21, 67, 45]. Later on, the same team that developed BOSS introduced a newer algorithm, Word ExtrAction for time SEries cLassification (WEASEL) [63]. A dictionary based classifier which is very similar to BOSS and can be thought of as an extension to it, with more focus on scalability [47]. WEASEL motivated their work with the absence of scalable classifiers, at the time, that can deal with big data sets; as the existing were either not scalable enough or not accurate enough. Despite being a non-ensemble classifier, WEASEL can compete with powerful ensembles in terms of accuracy without the need for long prediction times, but this comes with the cost of resource expensive training [47].

Like the other dictionary based approaches, WEASEL uses a sliding window over the data instances to extract substructures. These substructures are then discretized per window to extract features and finally a machine learning classifier is learned over the transformed features. But WEASEL differentiates itself from the other algorithms in the way it constructs and filters the features. We will discuss the new approaches that WEASEL uses for extracting discriminative features, building a model to deal with multiple occurrences of words and variable windows' lengths and selecting features to reduce runtime and exclude irrelevant features.

In the beginning, WEASEL uses multiple sliding windows of different lengths; to extract substructures. While keeping track of their order; in order to use them later on as bi-gram features. This replaces the process of building multiple models then choosing the best one, with building only one model which can learn from the concatenated high-dimensional feature vector. Then each of the substructures is normalized and a DFT is applied. Instead of filtering out the higher Fourier coefficients, WEASEL applies an ANalysis Of VAriance (ANOVA) F-test to keep real and imaginary Fourier values that best separate instances from different classes. Then the kept coefficients are discretized into words, based on alphabet of size  $c$ , using binning boundaries. Instead of using equi-depth binning, WEASEL applies an information gain based binning technique; which further more helps separating instances of classes. In the end, a histogram is built using all the windows lengths and the extracted features, uni-grams and bi-grams. Irrelevant features are then filtered out using a Chi-squared test. The final highly discriminative feature vector produced is then used to learn a logistic regression classifier.

WEASEL, like BOSS, transforms extracted windows from a time series into words using an alphabet of size  $c$ . They identify two main drawbacks for SFA, and introduce a new supervised symbolic representation technique which is based on SFA but overcomes the drawbacks. The first drawback of SFA is that it acts like a low pass filter and excludes the high frequency components from the Fourier transformation, which might discard important features in some scenarios. The other drawback, is that SFA defines the boundaries of bins during quantization independent of the class labels; which might cause SFA words of equal frequencies to appear unnecessarily in multiple classes. In order to overcome the drawbacks of SFA, WEASEL follows two steps; discriminative approximation using ANOVA F-test and discriminative quantization using information gain.

As mentioned earlier, approximation involves representing a time series of length  $n$  with a shorter, yet informative, representation of length  $l$  by applying a Fourier transformation. WEASEL aims at keeping the best class separating  $l$ , real and imaginary, Fourier coefficients by applying a one-way ANOVA F-test. The test verifies the hypothesis that the means of two or more distributions/groups differ from each other [44]. This can be tested by comparing two variances; the variance between groups and the variance within the groups. Using the notations in [63], *mean square between* ( $MS_B$ ) and *mean square within* ( $MS_W$ ) respectively. The F-values is then calculated as :

$$F = \frac{MS_B}{MS_W}$$

**Definition 12**

Then  $l$  coefficients with the highest F-values are kept; as these represent features which have big differences between the classes (high  $MS_B$ ) and small differences within the same class (low  $MS_W$ ).

After that, discretization is carried out to set the split thresholds for each of the extracted Fourier value. Discretization involves a binning process, where each Fourier value is divided into a number of bins. Each bin is represented by an upper and a lower boundary, and assigned a letter from an alphabet  $c$ . Previous quantization techniques used equi-depth or equi-width binning, but these techniques are solely based on values and ignore the distribution of classes. WEASEL applies a binning technique based on IG; assuring that for each partition the majority of the values would belong to the same class.

The result of the feature extraction phase is a high dimensional feature vector with a dimensionality of  $\mathcal{O}(\min(Nn^2, c^l))$ , where  $c$  is an alphabet,  $l$  is the length of a word,  $N$  is the total number of instances and  $n$  is the length of the time series. Since WEASEL uses bi-grams and  $\mathcal{O}(n)$  window lengths, the dimensionality of the feature space increases to  $\mathcal{O}(\min(Nn^2, c^{2l} \cdot n))$ . This enormous feature space is then reduced using a Chi-squared ( $\chi^2$ ) test. Chi-squared tests is a statistical test used to determine if there is a significant difference between the recognised frequencies and the expected frequencies of a features within a group. This implies that features which have high ( $\chi^2$ ) values are statistically frequent in certain classes. By comparing the ( $\chi^2$ ) values of features to a threshold, all features with scores lower than the threshold can be excluded from the feature space. This usually reduced the feature space by 30% - 70% and helped train the logistic regression classifier in a timely manner.

### 3.4.2 BOSS

Bag of SFA Symbols (BOSS) is a dictionary based algorithm that was introduced by [61] in 2015. Like the previous dictionary based classifiers, Bag Of Patterns (BOP) and Symbolic Aggregate approXimation in Vector Space Model (SAXVSM), BOSS also used a windowing technique to extract patterns from the data and transform them into words, but it had other significant differences to them [5]. BOSS was concerned with the issue of dealing with noisy data which, at that time, received little attention; due to the common practice of either using raw data directly or handling noise in a preprocessing stage. The goal was to introduce an algorithm faster than it's rivals of the same group, robust for existence of noise in the data and competitive with existing TSCA.

The BOSS model is divided into several steps. In the beginning, it passes a window of size  $w$  over the time series to extract substructures. Each of the extracted windows is then normalized to obtain amplitude invariance, while obtaining offset invariance by subtracting the mean value for each window is data set specific and can be decided upon based on a parameter. Then the substructures are quantized using SFA transformations which transforms the data into unordered words; this helps further reduce the noise in the data by making it phase shift invariant and makes it possible to apply string matching algorithms on the data. Since some data sets might happen to have long constant signals, this will lead to SFA transformations of their windows to produce the same words multiple times and cause higher weights to be assigned for them. BOSS applies a numerosity reduction technique adopted from [40, 39] that ignores multiple sequential occurrences of the same word. Finally time series instances can be compared for differences, using the noise reduced substructures using a customized distance measure inspired by Euclidean distance. We will, briefly, discuss the main components that are used for each of the previously mentioned steps.

To extract the substructures from a given time series instance  $T = [t_1, t_2, \dots, t_n]$ , a windowing function is used to split it into fixed sized windows  $S_{i:w} = (t_i, \dots, t_{i+w-1})$ , each of them is of a size  $w$ . The total number of windows that can be created for a time series of length  $n$  is  $n-w+1$ , and each consecutive windows overlap on  $w-1$  points. In order to achieve offset and amplitude invariance, each window is z-normalized by subtracting the mean from it's original values and then dividing the difference by the standard deviation.

**Definition 13**  $windows(T, w) = \{S_{1:w}, S_{2:w}, \dots, S_{n-w+1:w}\}$

After running the windowing function, the real values of the time points inside the windows are transformed into words using Symbolic Fourier Approximation (SFA) [62]. SFA is an alternative way of representing time series data, that instead of using real values, uses a sequence of symbols which are referred to as SFA words based on a predefined alphabet of specific length. SFA accomplishes two things; low pass filtering through removal of noisy components from the data and string representation which allows for the use of string matching algorithms.

For SFA to achieve it's goals, two main operations have to be carried out; approximation and quantization. Approximation is the process of representing a specific signal of length  $n$  using another signal of lower length  $l$ . This is achieved using Discrete Fourier Transformation (DFT); a transformation technique which is applied to a signal represented by a sequence of equally spaced values into a sequence of coefficients of complex sinusoid waves ordered by their frequencies [37]. The higher coefficients in



a DFT refer to data with rapid changes, which can be considered as noise in the signal and thus ignored. So considering only the first  $l \ll n$  coefficients acts as a low pass filter for noise producing a smoother signal.

Quantization also helps reduce noise by splitting the frequency domain into equi-depth bins, then maps each of the Fourier real and imaginary coefficients into a bin. BOSS utilizes Multiple Coefficient Binning (MCB), an adaptive technique that minimizes the loss of information which is a side effect of quantization. After the approximation step, a matrix is built from the Fourier transformations of the training data set using only the first  $\frac{l}{2}$  coefficients, including both the real and imaginary values for each coefficient. Then using an alphabet  $\Sigma$  of size  $c$ , MCB creates for each column of the matrix  $c+1$  break points using equi-depth binning. During classification, to acquire the SFA word for a Fourier transformed time series, a lookup is done on the precomputed MCB bins and a word is assigned if the value falls within its bin's boundaries.

After the transformation of instances to SFA words, BOSS uses a customized distance measure referred to as BOSS distance; to measure the similarity between the transformed instances. BOSS distance is a variation of Euclidean distance, which compares instances based on the histograms formed from their transformed forms. The appearance of the same SFA words in both instances, is considered to be a notion of similarity. While the absence of SFA words might be caused by one of two reasons; the absence of some substructures from either instances, or due to the presence of noise which disfigures the time series. Instances under BOSS distance are compared based on shared SFA words only, thus excluding words of frequencies equal to 0. BOSS distance as noted by [61] is:

**Definition 14** *Given two time series instances  $T_1$  and  $T_2$  and their corresponding BOSS histograms  $B_1 : \Sigma^l \rightarrow \mathbb{N}$  and  $B_2 : \Sigma^l \rightarrow \mathbb{N}$ , where  $l$  is the word length and  $\Sigma$  is an alphabet of size  $c$ . Their BOSS distance is:*

$$D(T_1, T_2) = \text{dist}(B_1, B_2) \\ \text{dist}(B_1, B_2) = \sum_{a \in B_1; B_1(a) > 0} [B_1(a) - B_2(a)]^2$$

Each single BOSS classifier utilizes 1-NN approach along with its BOSS model. The reason behind using 1-NN is that it is a simple technique that doesn't add to the parameters of the model, but rather proved to be a robust one. When classifying a query instance, BOSS searches for the nearest neighbor in a sample of candidates by choosing the closest instance based on the BOSS distance.

Finally, BOSS Ensemble is introduced as an ensemble of multiple BOSS classifiers. While a fixed window length is used over time series instances for BOSS classifier, BOSS Ensemble considers representing each time series instance by ensembling multiple BOSS classifiers, each of a different window length. When the training data is fitted, BOSS Ensemble acquires a group of scores for each of the different window lengths. To classify a query instance, BOSS Ensemble acquires the best accuracy score from the score sets returned during training. Then considers all window lengths that obtained accuracies within a factor of the best accuracy score. Each of the considered windows predicts a class label for the query instance using 1-NN. Majority voting is applied and the most dominant class label is assigned.

### 3.5 Deep Learning Algorithms

Deep learning is a long established group of machine learning algorithms which has proved competence in many areas [36] and which have encouraged the introduction of

deep learning algorithms for time series classification [76]. Deep learning is appealing for investigating time series data; due to the role that the time dimension play as a structure for the data and because deep learning algorithms can scale linearly with training size [67].

The general framework of deep learning neural networks is described by [19] as an architecture of  $L$  layers, each one of them is considered an input domain representation. Each of the layers consists of small computing units referred to as neurons. Neurons compute elements for the output of each layer. A layer  $l_i$  applies an activation function to its input then passes the output to the next layer  $l_{i+1}$ . The behavior of activation functions in each layer is controlled by a set of parameters  $\theta_i$  and are referred to by the term weights. The weights are assigned to the links between the inputs and outputs of the network's layers. A neural network carries out a series of computations to predict the class label for a given input  $x$ , these calculation are noted as:

**Definition 15**  $f_L(\theta_L, x) = f_{L-1}(\theta_{L-1}, f_{L-2}(\theta_{L-2}, \dots, f_1(\theta_1, x)))$

Where  $f_i$  represents the activation function applied at layer  $l_i$

In the training phase of a neural network, the weights are randomly initialized. Then a forward pass of the input  $x$  is done through the model and an output vector is computed, in which every component of the vector represent a class probability. Using a cost function, the model's prediction loss is calculated and weights of the network are updated, in a backward pass process, using a gradient descent. By continuous iteration of forward passes and weight updates by backward passes, the neural network learns the best weights to minimize the cost function. For the prediction of unseen instances, or inference phase, the input data is passed through the network in a forward pass. Then using the class probabilities of the output vector, the class with the highest probability is assigned.

Many of the deep learning research on time series focused on the use of variants of Convolutional Neural Networks (CNN). [20] CNN model is based on the idea of learning convolutional filters that can accurately represent the data. Using these filters, CNNs can learn the hierarchical structure of the data while incorporating translation invariance[35]. Multi-Scale Neural Networks (MCNN) [14] and LeNet [35] were among the first experiments of using CNN in time series classification. MCNN was composed of three stages; transformation stage, local convolution stage and full convolution stage. The transformation stage aimed at applying a low pass filter to exclude noise and to capture different temporal patterns. Local convolution is a stage where different scale features are extracted and local features of the time series are learned. In the final stage, the full convolution stage, all the features are concatenated and the final prediction is made. On the other hand, LeNet consists of 2 convolutional layers, after each one a max pooling is carried out to do sub-sampling. For the first layer, 5 filters are applied and the max pooling size is 2. While for the second there are 20 filters and max pooling size is 4. In [76] MultiLayer Perceptrons (MLP), Fully Convolutional Networks (FCN) and the Residual Networks (ResNet) were tested on 44 datasets from the UCR time series archive, where FCN was not significantly different from state of the art and MCNN was not significantly different from COTE and BOSS [63].

A recent more comprehensive review of deep learning algorithms was done by [19]. In which an experiment between nine deep learning algorithms was done. These included the classical MCNN and LeNet, in addition to the three algorithms from the previously mentioned experiment; MLP, FCV and ResNet. The other algorithms included were Encoder [66], Multi Channel Deep Convolutional Neural Network (MCD-CNN) [83, 84], Time Convolutional Neural Network [82] and Time Warping Invariant

Echo State Network (TWIESN) [73]. For more details about the algorithms and their structures, please refer to [19].

Finally, the newest deep learning algorithm introduced for time series classification is InceptionTime [21], which is the newest state of the art and which is an ensemble of deep convolutional neural networks. It can attain high accuracy scores while maintaining scalability. We will discuss in more details the InceptionTime algorithm in the next section.

### 3.5.1 Inception Time

InceptionTime [21] is a recent deep learning algorithm for TSC problems, which is able to achieve high accuracy score [57]. Motivated by the increasing interest in deep learning algorithms in the TSC domain along with the need for scalable algorithms that can deal with large data sets, either in number of instances or in length of the time series, and scale to them. InceptionTime is inspired by AlexNet [34], an algorithm that was considered a breakthrough for deep learning algorithms [2], and wanted to achieve the same but for the domain of TSC.

InceptionTime is based on the idea that the combination of deep CNN with residual connections, like in ResNet, can attain higher classification performance [19]. Since CNN has proved competency with image classification, there seemed a potential opportunity to be able to use deeper CNN for time series; since time series data is mainly structured on only one dimension which is time, while images have two spacial dimensions. This opened the door for using more complex models for TSC problems which would be computationally challenging to use for images. The building blocks of an InceptionTime network are called Inception modules, these are were introduced by [69] and evolved later on to Inceptionv4 [70]. The InceptionTime classifier is an ensemble of 5 InceptionTime networks each initialized with random weights and are assigned equal weights for the final prediction. We will discuss in more details the structure of an InceptionTime network using the notation mentioned in [21].

InceptionTime's structure is similar to that of a ResNet, but instead of using three residual blocks, InceptionTime uses only two. Each of the residual blocks is composed of three Inception modules instead of the traditional fully convolutional network. Like residual networks, a linear skip connection exists between the two residual blocks of the network, passing the input from one block to be concatenated with the input of the other; this helps passing information from earlier layers of the network to deeper layers and thus mitigating the vanishing gradient issue [28]. After the residual blocks, a Global Average Pooling (GAP) layer exists where the multivariate time series output is averaged over the time dimension. The final component of the network is a conventional fully-connected softmax layer with a count of neurons similar to the count of output classes.

Inside the inception module, there are two main components; the bottleneck layer and the variable length filters. Assuming that the input is a multivariate time series data of  $M$  dimensions. The job of the bottleneck layer is to transform the data from having  $M$  dimensions, into a multivariate data set having  $m$  dimensions, where  $m \ll M$ . This is done by passing a group of sliding filters  $m$  with length 1 and a stride of size 1. Which will substantially reduce the dimensionality of the time series data and consequently will also decrease the model's complexity making it more robust for overfitting problems on data sets of small sizes. The other benefit of including the bottleneck layer, is that it allows utilizing longer filters on the data than the original ResNet using approximately the same number of parameters; due to the lower number

of dimensions that the filters will have to deal with.

The output of the bottleneck layer is then passed for a set of variable length filters, the second component of the network, of length  $l$  where  $l \in 10, 20, 40$ . In addition to the filters, a parallel MaxPooling operation is carried out followed by a bottleneck layer; to make the model robust to small data noises. The final multivariate output is then formed by concatenating the output from each filter based convolution along with the output of the MaxPooling operation. This whole process is executed for each Inception module in the network.

In the end, an InceptionTime network is able to learn the underlying hierarchical structures of a time series by stacking multiple Inception modules and learning from the different filter sizes, which had been learned during training, inside them.

The final InceptionTime classifier is an ensemble of 5 InceptionTime networks. InceptionTime networks accuracy scores showed high variances, a problem which have been discussed by [60] and was found in ResNet networks as well [20]. This happens due to the random initialization of the networks and due to the stochastic approach used for optimization. InceptionTime classifier follows an ensembling technique for neural networks to handle TSC [20] in order to leverage the high variability; adhering to the idea that combining multiple networks would yield better results than one classifier. During the classification of an instance, InceptionTime combines the logistic output of the 5 networks and assigns an equal weight to each of them. This can be denoted by the equation

**Definition 16**  $\hat{y}_{i,c} = \frac{1}{n} \sum_{j=1}^n \sigma_c(x_i, \theta_i) | \forall_c \in [1, C]$

Where  $\hat{y}_{i,c}$  is the class probability for instance  $x_i$  as belonging to the class  $c$ , which is the averaged logistic output  $\sigma_c$  over the randomly initialized models  $n$ .

## 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] Md Zahangir Alom et al. “The history began from alexnet: A comprehensive survey on deep learning approaches”. In: *arXiv preprint arXiv:1803.01164* (2018).
- [3] 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.
- [4] Anthony Bagnall et al. “Time-series classification with COTE: the collective of transformation-based ensembles”. In: *IEEE Transactions on Knowledge and Data Engineering* 27.9 (2015), pp. 2522–2535.
- [5] 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.
- [6] 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.
- [7] 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.
- [8] Laurent Bernaille et al. “Traffic classification on the fly”. In: *ACM SIGCOMM Computer Communication Review* 36.2 (2006), pp. 23–26.
- [9] Aaron Bostrom. “Shapelet Transforms for Univariate and Multivariate Time Series Classification”. PhD thesis. University of East Anglia, 2018.
- [10] Aaron Bostrom and Anthony Bagnall. “A shapelet transform for multivariate time series classification”. In: *arXiv preprint arXiv:1712.06428* (2017).
- [11] 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](https://doi.org/10.1007/978-3-662-55608-5_2).
- [12] Leo Breiman. “Random forests”. In: *Machine learning* 45.1 (2001), pp. 5–32.
- [13] 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.

- [14] Zhicheng Cui, Wenlin Chen, and Yixin Chen. “Multi-scale convolutional neural networks for time series classification”. In: *arXiv preprint arXiv:1603.06995* (2016).
- [15] 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.
- [16] Houtao Deng et al. “A time series forest for classification and feature extraction”. In: *Information Sciences* 239 (2013), pp. 142–153.
- [17] 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.
- [18] Christos Faloutsos, Mudumbai Ranganathan, and Yannis Manolopoulos. “Fast subsequence matching in time-series databases”. In: *Acm Sigmod Record* 23.2 (1994), pp. 419–429.
- [19] 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.
- [20] 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.
- [21] Hassan Ismail Fawaz et al. “Inceptiontime: Finding alexnet for time series classification”. In: *Data Mining and Knowledge Discovery* 34.6 (2020), pp. 1936–1962.
- [22] 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.
- [23] 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.
- [24] Tomasz Górecki and Maciej Łuczak. “Using derivatives in time series classification”. In: *Data Mining and Knowledge Discovery* 26.2 (2013), pp. 310–331.
- [25] 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.
- [26] 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.
- [27] Ashish Gupta et al. “A Fault-Tolerant Early Classification Approach for Human Activities using Multivariate Time Series”. In: *IEEE Transactions on Mobile Computing* (2020).

- [28] Kaiming He et al. “Deep residual learning for image recognition”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.
- [29] Jon Hills et al. “Classification of time series by shapelet transformation”. In: *Data Mining and Knowledge Discovery* 28.4 (2014), pp. 851–881.
- [30] 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.
- [31] 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.
- [32] 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.
- [33] 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.
- [34] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. “Imagenet classification with deep convolutional neural networks”. In: *Advances in neural information processing systems* 25 (2012), pp. 1097–1105.
- [35] Arthur Le Guennec, Simon Malinowski, and Romain Tavenard. “Data augmentation for time series classification using convolutional neural networks”. In: *ECML/PKDD workshop on advanced analytics and learning on temporal data*. 2016.
- [36] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. “Deep learning”. In: *nature* 521.7553 (2015), pp. 436–444.
- [37] Xin Liao, Kaide Li, and Jiaojiao Yin. “Separable data hiding in encrypted image based on compressive sensing and discrete fourier transform”. In: *Multimedia Tools and Applications* 76.20 (2017), pp. 20739–20753.
- [38] 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.
- [39] Jessica Lin, Rohan Khade, and Yuan Li. “Rotation-invariant similarity in time series using bag-of-patterns representation”. In: *Journal of Intelligent Information Systems* 39.2 (2012), pp. 287–315.
- [40] Jessica Lin et al. “Experiencing SAX: a novel symbolic representation of time series”. In: *Data Mining and knowledge discovery* 15.2 (2007), pp. 107–144.
- [41] 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.



- [42] 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).
- [43] Markus Löning et al. “sktime: A unified interface for machine learning with time series”. In: *arXiv preprint arXiv:1909.07872* (2019).
- [44] Richard Lowry. “Concepts and applications of inferential statistics”. In: (2014).
- [45] 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.
- [46] 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.
- [47] Matthew Middlehurst, William Vickers, and Anthony Bagnall. “Scalable dictionary classifiers for time series classification”. In: *International Conference on Intelligent Data Engineering and Automated Learning*. Springer, 2019, pp. 11–19.
- [48] 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.
- [49] 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.
- [50] Usue Mori et al. “Early classification of time series using multi-objective optimization techniques”. In: *Information Sciences* 492 (2019), pp. 204–218.
- [51] 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.
- [52] Nathan Parrish et al. “Classifying with confidence from incomplete information”. In: *The Journal of Machine Learning Research* 14.1 (2013), pp. 3561–3589.
- [53] 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.
- [54] 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.

- [55] 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.
- [56] 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.
- [57] Alejandro Pasos Ruiz et al. “The great multivariate time series classification bake off: a review and experimental evaluation of recent algorithmic advances”. In: *Data Mining and Knowledge Discovery* (2020), pp. 1–49.
- [58] 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.
- [59] Tiago Santos and Roman Kern. “A Literature Survey of Early Time Series Classification and Deep Learning.” In: *Sami@ iknow*. 2016.
- [60] Simone Scardapane and Dianhui Wang. “Randomness in neural networks: an overview”. In: *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 7.2 (2017), e1200.
- [61] 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.
- [62] Patrick Schäfer and Mikael Höggqvist. “SFA: a symbolic fourier approximation and index for similarity search in high dimensional datasets”. In: *Proceedings of the 15th international conference on extending database technology*. 2012, pp. 516–527.
- [63] 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.
- [64] Patrick Schäfer and Ulf Leser. “Multivariate time series classification with WEASEL+ MUSE”. In: *arXiv preprint arXiv:1711.11343* (2017).
- [65] Patrick Schäfer and Ulf Leser. “TEASER: early and accurate time series classification”. In: *Data Mining and Knowledge Discovery* 34.5 (2020), pp. 1336–1362.
- [66] Joan Serrà, Santiago Pascual, and Alexandros Karatzoglou. “Towards a Universal Neural Network Encoder for Time Series.” In: *CCIA*. 2018, pp. 120–129.
- [67] 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.

- [68] 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.
- [69] Christian Szegedy et al. “Going deeper with convolutions”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 1–9.
- [70] Christian Szegedy et al. “Inception-v4, inception-resnet and the impact of residual connections on learning”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 31. 1. 2017.
- [71] 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.
- [72] Chang Wei Tan et al. “Time series classification for varying length series”. In: *arXiv preprint arXiv:1910.04341* (2019).
- [73] Pattreeya Tanisaro and Gunther Heidemann. “Time series classification using time warping invariant echo state networks”. In: *2016 15th IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE. 2016, pp. 831–836.
- [74] 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.
- [75] 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.
- [76] Zhiguang Wang, Weizhong Yan, and Tim Oates. “Time series classification from scratch with deep neural networks: A strong baseline”. In: *2017 International joint conference on neural networks (IJCNN)*. IEEE. 2017, pp. 1578–1585.
- [77] Zhengzheng Xing, Jian Pei, and Eamonn Keogh. “A brief survey on sequence classification”. In: *ACM Sigkdd Explorations Newsletter* 12.1 (2010), pp. 40–48.
- [78] 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.
- [79] Zhengzheng Xing, Jian Pei, and S Yu Philip. “Early classification on time series”. In: *Knowledge and information systems* 31.1 (2012), pp. 105–127.
- [80] Omolbanin Yazdanbakhsh and Scott Dick. “Multivariate Time Series Classification using Dilated Convolutional Neural Network”. In: *arXiv preprint arXiv:1905.01697* (2019).
- [81] 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.

- [82] Bendong Zhao et al. “Convolutional neural networks for time series classification”. In: *Journal of Systems Engineering and Electronics* 28.1 (2017), pp. 162–169.
- [83] Yi Zheng et al. “Time series classification using multi-channels deep convolutional neural networks”. In: *International conference on web-age information management*. Springer. 2014, pp. 298–310.
- [84] Yi Zheng et al. “Exploiting multi-channels deep convolutional neural networks for multivariate time series classification”. In: *Frontiers of Computer Science* 10.1 (2016), pp. 96–112.

## 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