

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



## MASTER THESIS

---

# A comparison of Time Series Classification Algorithms based on their ability to learn in an Early Classification Context

---

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, GERMANY

2ND SUPERVISOR:  
PROF. PHD. PANAGIOTIS PAPAPETROU  
DEPARTMENT OF COMPUTER AND SYSTEMS SCIENCES  
STOCKHOLM UNIVERSITY, SWEDEN

25.05.2021

Wahba, Ismail:  
A comparison of Time Series Classification Algorithms based on their ability to learn  
in an Early Classification Context  
Master Thesis, Otto-von-Guericke-University Magdeburg, 2021.

## Abstract

Time series classification (TSC) is a field of machine learning that aims at learning predictive models that can assign class labels for time series. Whether the time series is univariate or multivariate, time series data exists in many real life applications around us. The abundance of time series data lead to an interest in extending TSC to incorporate early decision making. A specialized field in learning early classification algorithms for time series, called Early Time Series Classification (eTSC), emerged. This new field helps in scenarios where time is an important factor; like diseases diagnosis and natural disasters prediction. Inspired by the eTSC problem and the competent Time Series Classification Algorithms (TSCAs), we define a context that simulates incremental learning of time series data and study its effect on the performance of TSCAs. We review existing algorithms from both TSC and eTSC, as well as previous frameworks for comparing TSCAs. Then implement a framework that simulates our proposed context through chopping of training data. We conduct experiments on 77 data sets from the UCR and UEA data archives and make comparisons within and across different TSCAs. The results demonstrate that each TSCA can achieve competent performance in the early classification context compared to its performance on the full length data. When compared against other on full data length, the TSCAs are at their full capacity and so they all attain high performance scores on the data sets. There is no statistical difference between any of the classifiers at that stage. On the other hand, when compared on early chunks of data, their performances are affected by the data chopping and statistical differences exist between their performances. TSF ranks as the best performing classifier on early chunk data compared to the others. To evaluate our framework, we show that it is possible to learn a recommender to predict the performance of the TSCAs in the early classification context for unseen data sets and suggest the good performing ones.

## 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 . . . . .	iv
List of Abbreviations . . . . .	vii
<b>1 Introduction</b>	<b>1</b>
<b>2 Classification Algorithms</b>	<b>4</b>
2.1 Time Series Data . . . . .	4
2.1.1 Definitions . . . . .	4
2.1.2 Nature of Time Series Data . . . . .	5
2.2 Time Series Classification . . . . .	5
2.3 Time Series Classification Algorithms . . . . .	6
2.3.1 Whole Time Series Algorithms . . . . .	6
2.3.2 Phase Dependent intervals Algorithms . . . . .	10
2.3.3 Phase Independent intervals Algorithms . . . . .	12
2.3.4 Dictionary Based Algorithms . . . . .	14
2.3.5 Deep Learning Algorithms . . . . .	20
2.4 Early Time Series Classification . . . . .	23
2.5 Early Time Series Classification Algorithms . . . . .	24
2.5.1 ECTS . . . . .	24
2.5.2 Shapelet Based Algorithms . . . . .	25
2.5.3 ECDIRE . . . . .	26
2.5.4 TEASER . . . . .	26
<b>3 Related Work on Comparing Time Series Classification Algorithms</b>	<b>28</b>
3.1 The great time series classification bake off . . . . .	28
3.1.1 Inclusion Criteria . . . . .	28
3.1.2 Experimental Design . . . . .	29
3.2 Deep learning for time series classification: a review . . . . .	29
3.2.1 Inclusion Criteria . . . . .	29
3.2.2 Experimental Design . . . . .	30
3.3 The great multivariate time series classification bake off . . . . .	30
3.3.1 Inclusion Criteria . . . . .	31
3.3.2 Experimental Design . . . . .	31
3.4 Evaluation Criteria . . . . .	32

<b>4 Framework for Comparing Time Series Classification Algorithms in Early Classification Context</b>	<b>33</b>
4.1 Early Classification Context . . . . .	33
4.2 Conceptual Design . . . . .	34
4.2.1 The Testbed . . . . .	34
4.2.2 The Recommender . . . . .	36
4.3 Implementation . . . . .	36
4.3.1 Design Protocol . . . . .	36
4.3.2 Toolkits . . . . .	37
4.3.3 The Testbed Implementation . . . . .	37
4.3.4 The Recommender Implementation . . . . .	45
<b>5 Data Sets and Experimental Design</b>	<b>48</b>
5.1 UCR/UEA Archives . . . . .	48
5.1.1 UCR Archive . . . . .	48
5.1.2 UEA Archive . . . . .	48
5.1.3 Used Data Sets . . . . .	49
5.2 Experimental Design . . . . .	55
5.2.1 The Testbed Experimental Design . . . . .	55
5.2.2 The Recommender Experimental Design . . . . .	59
<b>6 Results</b>	<b>60</b>
6.1 Testbed Results . . . . .	60
6.1.1 Results Across Classifiers . . . . .	68
6.1.2 Data Set characteristics and performance in early classification context . . . . .	69
6.1.3 Results Within Classifier . . . . .	77
6.1.4 Results on Runtime Duration . . . . .	78
6.2 Recommender Results . . . . .	82
6.2.1 Are the chunk learners able to learn F-scores for the classifiers ? . . . . .	82
6.2.2 Which features contribute the most to the F-Score Prediction ? . . . . .	85
6.2.3 Are the final recommendations of the framework consistent with actual results ? . . . . .	86
<b>7 Conclusion and Future Work</b>	<b>90</b>
<b>A Appendix - Testbed Results</b>	<b>92</b>
A.1 Results Across Classifiers . . . . .	92
A.2 Data Set characteristics and performance in early classification context . . . . .	94
A.3 Results Within Classifier . . . . .	100
A.4 Results on Runtime Duration . . . . .	101
<b>B Recommender Results</b>	<b>118</b>
B.1 Are the chunk learners able to learn F-scores for the classifiers ? . . . . .	118
B.2 Which features contribute the most to the F-Score Prediction ? . . . . .	123
B.3 Are the final recommendations of the framework consistent with actual results ? . . . . .	124

# List of Figures

2.1	Mapping of Euclidean distance (lock-step measure) versus mapping of DTW distance (elastic measure) [AML19]. . . . .	7
2.2	Examples of the Move, Split, Merge operations [SAD12]. . . . .	8
2.3	Visual depiction of the root node for the “Trace” dataset (simplified to 2 classes). The top chart represents the data at the root node (one colour per class) while the data at the bottom left and right represent the data once split by the tree. The two time series in the middle left and right are the exemplars on which the tree is splitting. The scatter plot at the center represents the distance of each time series at the root node with respect to the left and right exemplars (resp. x- and y-axes) (Color figure online) [Luc+19]. . . . .	10
2.4	An illustration of two shapelets S1, S2 (leftmost plots) learned on the Coffee dataset. Series’ distances to shapelets can optimally project the series into a 2-dimensional space, called the shapelet-transformed representation [Lin+12] (rightmost plot). The middle plots show the closest matches of the shapelets on series of two classes having light-blue and black colors [Gra+14]. . . . .	12
2.5	The BOSS workflow [Sch15] . . . . .	16
2.6	A time series is (a) approximated (low pass filtered) using DFT and (b) quantised using MCB resulting in the SFA word DAAC [Sch15] . . . . .	17
2.7	WEASEL Pipeline: Feature extraction using a novel supervised symbolic representation, and a novel bag-of-patterns model [SL17a] . . . . .	19
2.8	On the left: Distribution of Fourier coefficients for a sample dataset. The high F-values on imag3 and real0 (red text at bottom) should be selected to best separate the samples from class labels ’A’ and ’B’. On the right: Zoom in on imag3. Information gain splitting is applied to find the best bins for the subsequent quantization. High information gain (IG) indicates pure (good) split points [SL17a] . . . . .	20
2.9	Inside Inception module for time series classification. For simplicity a bottleneck layer of size $m = 1$ is used [Faw+20] . . . . .	23

2.10 TEASER is given a snapshot of an energy consumption time series. After seeing the first s measurements, the first slave classifier $sc_1$ performs a prediction which the master classifier $mc_1$ rejects due to low class probabilities. After observing the i-th interval which includes a characteristic energy burst, the slave classifier $sc_i$ (correctly) predicts RECEIVER, and the master classifier $mc_i$ eventually accepts this prediction. When the prediction of RECEIVER has been consistently derived v times, it is output as final prediction [SL20] . . . . .	27
4.1 Conceptual design of the proposed framework . . . . .	35
5.1 Metadata of the used data sets . . . . .	50
5.2 Comparison between Balanced Accuracy (left) and $F_1$ (right) for CBoss . . . . .	58
5.3 Comparison between Balanced Accuracy (left) and $F_{0.5}$ (right) for CBoss . . . . .	58
6.1 Critical Difference diagrams per chunk across all classifiers . . . . .	69
6.2 Crtitical difference diagrams between versions of the same classifier trained on different chunks ( $F_\beta$ ) . . . . .	77
6.3 Crtitical difference diagrams between versions of the same classifier trained on different chunks (balanced accuracy) . . . . .	78
6.4 Train Time (CPU Time in Log Scale) for CBoss across all chunks . . . . .	79
6.5 Train Time (CPU Time in Log Scale) for PForest across all chunks . . . . .	79
6.6 Train Time (CPU Time in Log Scale) for ST across all chunks . . . . .	80
6.7 Train Time (CPU Time in Log Scale) for TSF across all chunks . . . . .	80
6.8 Train Time (CPU Time in Log Scale) for WEASEL across all chunks . . . . .	81
6.9 Histograms for distribution of RMSE per chunk over 50 runs . . . . .	83
6.10 Avg MAE for classifiers per chunk over 50 runs . . . . .	84
6.11 Avg RMSE for classifiers per chunk over 50 runs . . . . .	84
6.12 Mean Feature Importance per chunk across 50 runs . . . . .	85
6.13 Actual vs. Predicted F-score for chunk learners (experiment seed = 22) . . . . .	87
A.1 Critical Difference diagrams per chunk across all classifiers . . . . .	92
A.2 Critical Difference diagram for all classifiers and all chunks ( $F_\beta$ ) . . . . .	93
A.3 Critical Difference diagram for all classifiers and all chunks (balanced accuracy) . . . . .	93
A.4 Crtitical difference diagrams between versions of the same classifier trained on different chunks ( $F_\beta$ ) . . . . .	100
A.5 Crtitical difference diagrams between versions of the same classifier trained on different chunks (balanced accuracy) . . . . .	100
A.6 Train Time (CPU Time in Log Scale) for all 5 classifiers for 10% chunk . . . . .	108
A.7 Train Time (CPU Time in Log Scale) for all 5 classifiers for 20% chunk . . . . .	108
A.8 Train Time (CPU Time in Log Scale) for all 5 classifiers for 30% chunk . . . . .	109
A.9 Train Time (CPU Time in Log Scale) for all 5 classifiers for 100% chunk . . . . .	109
A.10 Train Time (CPU Time in Log Scale) for CBOss per series length . . . . .	110
A.11 Train Time (CPU Time in Log Scale) for PForest per series length . . . . .	110
A.12 Train Time (CPU Time in Log Scale) for ST per series length . . . . .	111
A.13 Train Time (CPU Time in Log Scale) for TSF per series length . . . . .	111
A.14 Train Time (CPU Time in Log Scale) for WEASEL per series length . . . . .	112
A.15 Train Time (CPU Time in Log Scale) for CBOss per train size . . . . .	112
A.16 Train Time (CPU Time in Log Scale) for PForest per train size . . . . .	113

A.17 Train Time (CPU Time in Log Scale) for ST per train size . . . . .	113
A.18 Train Time (CPU Time in Log Scale) for TSF per train size . . . . .	114
A.19 Train Time (CPU Time in Log Scale) for WEASEL per train size . . .	114
A.20 Train Time (CPU Time in Log Scale) for CBOss per number of dimensions	115
A.21 Train Time (CPU Time in Log Scale) for PForest per number of dimensions . . . . .	115
A.22 Train Time (CPU Time in Log Scale) for ST per number of dimensions	116
A.23 Train Time (CPU Time in Log Scale) for TSF per number of dimensions	116
A.24 Train Time (CPU Time in Log Scale) for WEASEL per number of dimensions . . . . .	117
B.1 Histograms for distribution of MAE per chunk over 50 runs . . . . .	123
B.2 Mean Feature Importance stacked by chunk . . . . .	124

# List of Tables

4.2	Input parameters required for using the framework . . . . .	38
4.4	Parameters and configuration for TSCAs . . . . .	39
4.6	Analysis collected for each experiment . . . . .	45
5.1	Summary of the 77 data sets used in the experiments . . . . .	54
6.1	$F_\beta$ scores for CBoss, PForest and ST on 77 data sets for all chunks . . . . .	62
6.1	$F_\beta$ scores for CBoss, PForest and ST on 77 data sets for all chunks (continued) . . . . .	63
6.1	$F_\beta$ scores for CBoss, PForest and ST on 77 data sets for all chunks (continued) . . . . .	64
6.2	$F_\beta$ scores for TSF and WEASEL on 77 data sets for all chunks . . . . .	65
6.2	$F_\beta$ scores for TSF and WEASEL on 77 data sets for all chunks (continued) . . . . .	66
6.2	$F_\beta$ scores for TSF and WEASEL on 77 data sets for all chunks (continued) . . . . .	67
6.3	Performance of classifiers with respect to length in comparison to baseline on 10% chunk data sets . . . . .	70
6.4	Performance of classifiers with respect to train size in comparison to baseline on 10% chunk data sets . . . . .	71
6.5	Performance of classifiers with respect to train size in comparison to baseline on 20% chunk data sets . . . . .	71
6.6	Performance of classifiers with respect to train size in comparison to baseline on 30% chunk data sets . . . . .	71
6.7	Performance of classifiers with respect to train size in comparison to baseline on 100% chunk data sets . . . . .	73
6.8	Performance of classifiers with respect to problem type in comparison to baseline on 10% chunk data sets . . . . .	73
6.9	Performance of classifiers with respect to problem type in comparison to baseline on 20% chunk data sets . . . . .	74
6.10	Performance of classifiers with respect to problem type in comparison to baseline on 30% chunk data sets . . . . .	74
6.11	Performance of classifiers with respect to problem type in comparison to baseline on 100% chunk data sets . . . . .	75
6.12	Performance of classifiers with respect to number of classes in comparison to baseline on 10% chunk data sets . . . . .	75
6.13	Performance of classifiers with respect to number of classes in comparison to baseline on 20% chunk data sets . . . . .	75

6.14 Performance of classifiers with respect to number of classes in comparison to baseline on 30% chunk data sets . . . . .	76
6.15 Performance of classifiers with respect to number of classes in comparison to baseline on 100% chunk data sets . . . . .	76
6.16 Overall performance of chunk learners in predicting actual $F_\beta$ scores over 50 runs . . . . .	82
6.17 Overall Performance of chunk learners in predicting good and bad performers over 50 runs . . . . .	86
6.18 Misestimation results per classifier per chunk learner over 50 runs . . . . .	88
 A.2 Performance of classifiers with respect to length in comparison to baseline on 20% chunk data sets . . . . .	94
A.4 Performance of classifiers with respect to length in comparison to baseline on 30% chunk data sets . . . . .	94
A.6 Performance of classifiers with respect to length in comparison to baseline on 100% chunk data sets . . . . .	95
A.8 Performance of classifiers with respect to train size in comparison to baseline on 10% chunk data sets . . . . .	95
A.10 Performance of classifiers with respect to train size in comparison to baseline on 20% chunk data sets . . . . .	95
A.12 Performance of classifiers with respect to train size in comparison to baseline on 30% chunk data sets . . . . .	96
A.14 Performance of classifiers with respect to train size in comparison to baseline on 100% chunk data sets . . . . .	96
A.15 Performance of classifiers with respect to problem type in comparison to baseline on 10% chunk data sets . . . . .	96
A.16 Performance of classifiers with respect to problem type in comparison to baseline on 20% chunk data sets . . . . .	97
A.17 Performance of classifiers with respect to problem type in comparison to baseline on 30% chunk data sets . . . . .	97
A.18 Performance of classifiers with respect to problem type in comparison to baseline on 100% chunk data sets . . . . .	98
A.19 Performance of classifiers with respect to number of classes in comparison to baseline on 10% chunk data sets . . . . .	98
A.20 Performance of classifiers with respect to number of classes in comparison to baseline on 20% chunk data sets . . . . .	98
A.21 Performance of classifiers with respect to number of classes in comparison to baseline on 30% chunk data sets . . . . .	99
A.22 Performance of classifiers with respect to number of classes in comparison to baseline on 100% chunk data sets . . . . .	99
A.23 Train Time (CPU Time) for CBoss, PForest and ST on 77 data sets for all chunks . . . . .	102
A.23 Train Time (CPU Time) for CBoss, PForest and ST on 77 data sets for all chunks (continued) . . . . .	103
A.23 Train Time (CPU Time) for CBoss, PForest and ST on 77 data sets for all chunks (continued) . . . . .	104
A.24 Train Time (CPU Time) for TSF and WEASEL on 77 data sets for all chunks . . . . .	105
A.24 Train Time (CPU Time) for TSF and WEASEL on 77 data sets for all chunks (continued) . . . . .	106

A.24 Train Time (CPU Time) for TSF and WEASEL on 77 data sets for all chunks (continued) . . . . .	107
B.1 Performance of chunk learners in predicting $F_\beta$ per chunk over 50 runs	122
B.2 Confusion Matrix of chunk learners in predicting good and bad classi- fiers per chunk over 50 runs . . . . .	128

## List of Abbreviations

ABR	Auditory Brainstem Response
BOSS	Bag of SFA Symbols
CNN	Convolutional Neural Networks
DFT	Discrete Fourier Transformation
DTW	Dynamic Time Warping
ECTS	Early Classification of Time Series
EE	Elastic Ensemble
ERP	Edit Distance with Real Penalty
eTSC	Early Time Series Classification
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	Mimimum Prediction Length
MSM	Move-Split-Merge
MTSC	Multivariate Time Series Classification
PAA	Piecewise Aggregate Approximation
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



# Chapter 1

## Introduction

Time Series Classification (TSC) is a research field which has grabbed a lot of attention in the data mining community during the last decade. This goes back to the fact that time series data exists, by nature, in numerous scenarios in our life; medical examination records of patients [GM01; GO12], consumption of electrical devices in smart homes and smart grids [Gao+14; JZ14] and astronomy [Pro+06; RK13] are only few examples of them.

The problem of time series data classification has been tackled in different ways and motivated by various objectives. When the main driving objective is to attain high classification accuracy, then this is referred to simply as Time Series Classification (TSC). TSC has also been extended to cover space and time complexity along with accuracy. On the other hand, if time plays a critical role in the problem being tackled beside accuracy; like early detection of malfunctions in industrial processes or diagnosis of a patient with a life threatening disease. In this case the classification problem is referred to as Early Time Series Classification (eTSC).

A lot of algorithms have been introduced to tackle the problem of TSC. According to [Bag+17], these algorithms can be divided, based on their technique, into six groups. Distance-based time series algorithms compare instances of time series, usually using 1-NN classifiers and an elastic distance measure between the data points of both time series. Phase dependent interval algorithms operate by extracting informative features from intervals of time series, they are more suitable for long and noisy data than distance-based time series algorithms. Phase independent interval algorithms 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 consider the number of repetitions of patterns as a factor of classification and not just simple occurrence of one. Ensembles 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 neural networks models.

On the other hand, Early Time Series Classification algorithms are designed to deal with less data in order to achieve earliness of prediction, without sacrificing accuracy as much as possible. Since accuracy and earliness are contradicting by nature, eTSC's main goal is to find a tradeoff between both objectives [Gup+20]. Many of the techniques that have been applied in TSC have also been applied for eTSC. Using 1-NN classifiers is a very popular TSC technique, which has been extended to eTSC by

combining it with Minimum Prediction Length (MPL)[XPP12]. Phase independent intervals[GO12; He+15], generative and discriminative model based classifiers[LLF14; Che+18] as well as deep learning algorithms [HLT18; Ruß+19] have all been extended to tackle eTSC problems.

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. The existence of comprehensive benchmark data archives [Dau+18; Bag+18] has helped regulate and standardize the process of comparison between classifiers. Algorithms can be benchmarked on a diverse set of datasets with different characteristics to avoid biased results and cherry picking.

Like other fields, the *No free lunch theorem* exists in the TSC domain as well [Bag+17; Bos18]. This means that there is no one specific algorithm that can generalize to all problems and prevail over all other competing algorithms. With that being the case, this leaves an open door for continuous investigation of how the performance of algorithms would change on different data sets and in different contexts.

Inspired by the of the problem of eTSC and its application in diverse areas. We define a situation in which a classifier is put into test to learn on only a part of the available data; to see if it will be able to gracefully degrade. We consider a classifier to be graceful, or fault tolerant, if it is able to maintain statistically significant results by learning on less data in comparison to its results when learning on full length data. We refer to this situation as the early classification context.

The motivation behind experimenting with such a context is driven by the medical field, where sometimes patients have to go through strenuous tests to be diagnosed. This is specially the case in audiological examinations; as patients are exposed to signals which pushes the limits of their hearing discomfort to be able to examine their problems. A classifier which can attain comparable results while only learning on a fraction of the whole time series data would shorten this process and thus relief patients from their discomfort.

## Goal of this Thesis

The goals of this master thesis can be addressed as:

- Provide a solution for adapting existing time series workflows for an early classification context. The solution should run experiments on a group of data sets from the UCR/UEA archive and measure the performance of classifiers within the early classification context. By the end of the experiments the solution should collect the results and provide an analysis of the results, the analysis should:
  - Identify if classifiers can sustain good performance while learning on fewer data points in comparison to learning on the full length
  - Identify if there is a winning classifier across the executed experiments
  - Identify the aspects of the data set which affect the performance of classifiers within the context. Such as train size, number of classes, number of dimensions, length, etc.
- Identify a suitable evaluation criteria for classifiers in the proposed early classification context, trading classification accuracy against time series length.
- Evaluate the proposed solution for adapting the time series workflows to be used for future recommendations



## Chapter 2

# Classification Algorithms

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. We will discuss TSCAs in more details as they are the main focus of our framework, but a review of the problem of early classification and a review of its algorithms is also included; as it motivates the context in which the framework operates.

## 2.1 Time Series Data

### 2.1.1 Definitions

A time series is a finite sequence of ordered observations, either based on time or another aspect [AML19; Bag+17]. The existence of the time component makes time series an abundant form of data that covers various domains like; medicine, finance, engineering and biology [LTB18].

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 collected at consequent time points.

#### Definition 2

*A time series  $T_i$ , of length  $L$  is represented as  $T_i = [t_1, t_2, \dots, t_L]$ .*

An observation  $t_i$  represents the value of the collected measurement at some point of time  $i$ .

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 of time series data is referred to as univariate time series. This form of data exists when the observations of instances capture a singular value across time.

**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.*

The other form is when multiple measurements are captured by the observations. According to [Lön+19], 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 one of them 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 independent and identically distributed (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.

### 2.1.2 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 i.i.d [Lön+19].

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 [BB17b]. 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 [Den+13].

Time series classification has been studied with different objectives, some papers focused on attaining the highest accuracy of classification as the main goal [Kat16; JJO11; BB17a; LTB18; SL17b; Faw+20], while other papers focused on attaining lower time complexity [RK04; Bag+17; TPW20; Pet+16; SL17a].

In this master thesis, we are more interested in assessing the results in terms of accuracy than time complexity. We define accuracy like [SL20]; as the percentage of

correctly classified instances for a given dataset D, either being a training or testing dataset.

#### Definition 5

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

## 2.3 Time Series Classification Algorithms

This chapter will introduce different types of TSCAs and discuss the various techniques that each type apply. There are multiple ways to divide TSCAs in order to better understand them. In this thesis we follow the grouping defined by [Bag+17].

### 2.3.1 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 [Kat16]. 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 [LTB18]. Whole time series algorithms are best suited for problems where the unique features can exist anywhere along the whole time series[Bag+17].

One of the simplest forms of whole time series is 1-NN with Euclidean Distance [FRM94], yet it can surprisingly attain high accuracy compared to other distance measures [XPK10]. But Nearest Neighbor with Euclidean Distance (KNNED) is an easy to beat baseline, due to its sensitivity for distortion and inability to handle time series of unequal lengths [XPK10; Kat16; LTB18]. This lead many of the researchers to focus on defining more advanced and elastic distance measures that can compensate for misalignment between time series [AML19]. The standard and most common baseline classifier utilizing elastic distance measures is 1-NN with Dynamic Time Warping (DTW) [Bag+17]. 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 [Kat16; LB15; Wan+13]. 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 [JJO11] and Derivative Dynamic Time Warping (DDTW) [KP01; GL13] which uses derivatives of sequences as features rather than raw data to avoid singularities. On the other hand, Edit Distance with Real Penalty (ERP) [CN04], Time Warp Edit (TWE) [Mar08], Longest Common Subsequence (LCSS) [DGM97] and Move-Split-Merge (MSM) [SAD12] are all alternatives for distance measures, yet multiple experiments have considered DTW to be relatively unbeatable [Bag+17; AML19; BB17a]. To the extent of our knowledge, the most powerful whole time series classifiers is Proximity Forest (PF) [Luc+19].

#### 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  $T_2 =$

$[t_2, t_2, \dots, t_2]$ , the euclidean distance between them can be determined as:

$$ED(T_1, T_2) = \sqrt{\sum_{i=1}^n (t_{1i} - t_{2i})^2} \quad (2.1)$$

Euclidean distance has been preferred to other classifiers due to its space and time efficiency, but it suffers from two main shortcomings [BRT13; JJO11; Kat16]. The first one is that it cannot handle comparisons between time series of different lengths. While the second one is its sensitivity to minor discrepancies between time series; it would calculate large distance values for small shifting 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 [Hil+14; Din+08; Bag+12].

### Nearest Neighbor with DTW

Dynamic Time Warping was a very strong baseline for time series classification for a long time [AML19; Bag+17]. 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 [TPW20]. Figure 2.1 shows a comparison between ED and DTW while computing the distance between two time series instances.

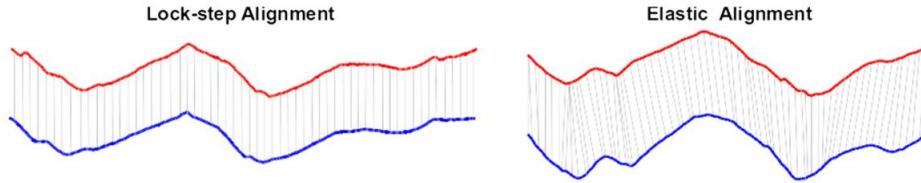


Figure 2.1: Mapping of Euclidean distance (lock-step measure) versus mapping of DTW distance (elastic measure) [AML19].

To calculate the distance between two time series instances  $T_1 = [t_{11}, t_{12}, \dots, t_{1n}]$  and  $T_2 = [t_{21}, t_{22}, \dots, t_{2n}]$ ; a distance matrix  $M(T_1, T_2)$ , of size  $n \times n$ , is calculated for  $T_1$  and  $T_2$ . With  $M_{i,j}(t_{1i}, t_{2j})$  representing the distance between  $t_{1i} \in T_1$  and  $t_{2j} \in T_2$ . 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 three conditions:

- $p_1 = (t_{11}, t_{21})$
- $p_p = (t_{1n}, t_{2n})$
- for all  $i < m$  :
  - $0 \leq t_{1i+1} - t_{1i} \leq 1$
  - $0 \leq t_{2i+1} - t_{2i} \leq 1$

Finding an optimal path under DTW can be computationally expensive with complexity of  $O(n^2)$  for a time series of length  $n$  [SL17a; Pet+16]. Consequently it is usual to use a constraint with the path; to prevent comparison of points outside a certain window [TPW20]; like the famous Sakoe-Chiba Band [SC78], Itakura Parallelogram [Ita75] and Ratanamahatana-Keogh Band [RK04]. Typically DTW can make use of its 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 [Tan+19].

### Nearest Neighbor with MSM Distance

Move-Split-Merge is distance measure that was first introduced in [SAD12]. 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 [LB15].

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 [Bag+17].

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.

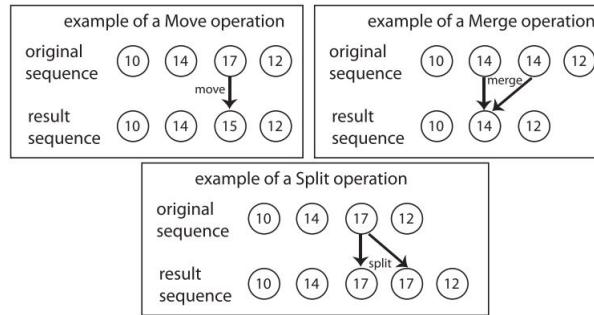


Figure 2.2: Examples of the Move, Split, Merge operations [SAD12].

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 [SAD12; TPW20].

### Proximity Forest

Proximity forest was developed by [Luc+19]. It was introduced as an addition to scalable time series classification, offering a more scalable and accurate classifier than EE [TPW20]. On one side, EE was an accurate classifier being the state of the

art algorithm 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 [LB15; Bag+17]. This goes back to the leave-one-out-cross-validation (LOOCV) used to optimize the parameters for each used metric [Shi+20].

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 [Luc+19].

Capitalizing on the previous research that has been put in developing specialized time series distance measures and inspired by the existing EE [Faw+20; Faw+19b]. Proximity forests combine 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 [Luc+19]. Proximity forests can scale sublinearly with training data set size, but quadratically with the length of the time series [Shi+20].

Proximity forests are based on a similar concept as Random Forests [Bre01], 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 its single classifiers [Luc+19].

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 its 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 [Luc+19], 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 [Faw+20; Faw+19a]. 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 (subtree) that the instance should follow based on its closest 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.

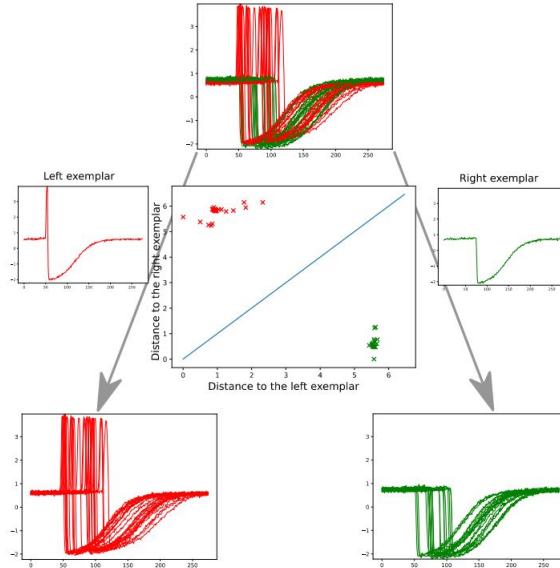


Figure 2.3: Visual depiction of the root node for the “Trace” dataset (simplified to 2 classes). The top chart represents the data at the root node (one colour per class) while the data at the bottom left and right represent the data once split by the tree. The two time series in the middle left and right are the exemplars on which the tree is splitting. The scatter plot at the center represents the distance of each time series at the root node with respect to the left and right exemplars (resp. x- and y-axes) (Color figure online) [Luc+19].

### 2.3.2 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 [BR16],

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 [Den+13].

According to [Bag+17], phase dependent algorithms are best used for problems where discriminatory information from intervals exist, this would be the case with long time 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 [Den+13]. There are  $n(n-1)/2$  possible intervals, for a time series of length  $n$  [Bag+17]. 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 [SK16; RA04; Den+13]. This wide space of features represents one of the main challenges that phase dependent algorithms face. [RA04] proposed a solution for this by only considering intervals with lengths equal to powers of two [Bag+17] in the features extraction step.

### Time Series Forest

Time Series Forest (TSF) is an algorithm that was introduced in 2013 by [Den+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:

$$\text{mean}(t_1, t_2) = \frac{\sum_{i=t_1}^{t_2} v_i}{t_2 - t_1 + 1} \quad (2.2)$$

$$\text{std}(t_1, t_2) = \begin{cases} \frac{\sum_{i=t_1}^{t_2} (v_i - \text{mean}(t_1, t_2))^2}{t_2 - t_1} & \text{if } t_1 < t_2 \\ 0 & \text{if } t_1 = t_2 \end{cases} \quad (2.3)$$

$$\text{slope}(t_1, t_2) = \begin{cases} m & \text{if } t_1 < t_2 \\ 0 & \text{if } t_1 = t_2 \end{cases} \quad (2.4)$$

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 Entrace. 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 [Bag+17] 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 [Den+13]. The final

classification of a testing instance is done using majority voting of all time series trees created.

### 2.3.3 Phase Independent intervals Algorithms

Phase independent intervals, or just shapelets as less formally known, are subseries which are ultimately distinctive of classes regardless of their place on the time series [SL17a; Bag+17]. They were first introduced in [YK09] 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 [BB17a; YK09; MKY11]. 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 [Bag+17; BB17b]. Figure 2.4 shows how shapelets can separate instances by calculating distances and provide interpretable explanation by projecting the data in a shapelet transform space.

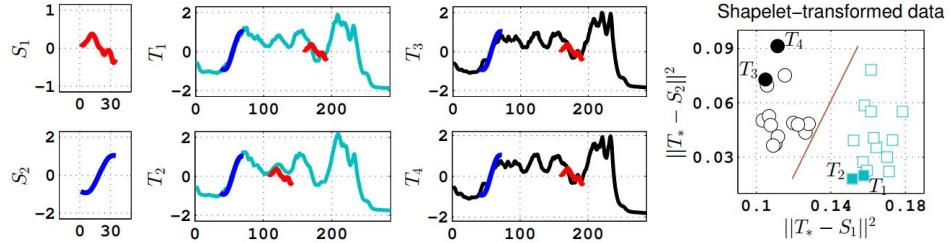


Figure 2.4: An illustration of two shapelets  $S_1, S_2$  (leftmost plots) learned on the Coffee dataset. Series' distances to shapelets can optimally project the series into a 2-dimensional space, called the shapelet-transformed representation [Lin+12] (rightmost plot). The middle plots show the closest matches of the shapelets on series of two classes having light-blue and black colors [Gra+14].

The original shapelet algorithm enumerated all possible shapelets and embeded 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 [LTB18; Sch15]. 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) [RK13] and Learned Shapelets (LS) [Gra+14]. FS applied discretization through Symbolic Aggregate Approximation (SAX) to reduce the length of time series, while LS tried to learn the shapelets [Shi+20]. Later on the idea of tranforming time series data to an alternative space was adopted in [Hil+14], the transformed data consists of distances to the best k shapelets, then classification is done using an ensemble of eight classifiers.

### Learned Shapelets

Learned Shapelets (LS) was proposed in [Gra+14] 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 [LTB18; Bos18]. The found shapelets need not to be a subsequence of one of the training examples [Bag+17; SL17a].

LS follows a two steps technique. In the begining 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.[Bos18]. In summary, the main objective of the aglorithm is to learn collectively the optimal shapelets and the weights linear hyper-plane that minimizes the objective function [Bag+17; Gra+14].

### Shapelet Transform

The first Shapelet Transform (ST) was introduced in [Hil+14]. While the original algorithm embeded shapelets discovery in decision trees and assesed 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 [Bos18]. ST seggregated the shapelets discovery process from the classifier. This segregation opened the door for choosing classifiers freely and considering more accurate classifiers than decision trees [Bag+17; LB15]. Also [Hil+14] 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 [Bos18].

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 [Hil+14].

ST was then extended again by [BB17b], the intuition behind it was that the previously used assessment technique couldn't hanlde multi-class problems [BB17b]. Instead of assessing shapelets that discriminate between all classes, they accomodated 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 [Bag+17]. For the classification, a

combination of tree based, kernel based and probabilistic classifiers were used in an ensemble on the transformed data set [Shi+20; LTB18]. Each of the classifiers was given a weight based on it's training accuracy and the final classification used weighted voting [BB17b]. Although ST has proved to be a competent accurate classifier, it suffers from high training-time complexity [Shi+20].

ST is one of multiple classifier in the bigger ensemble HIVE-COTE [LTB18]. It also is among the slowest components of the ensemble; due to the exhaustive shapelet enumeration which reaches a complexity of  $O(n^2 \cdot l^4)$ , where  $n$  is the size of the data set and  $l$  is the length of the time series [Shi+20]. During a modification that was introduced in [Bag+20] for HIVE-COTE, some enhancements were applied to ST. Some of these enhancements were new, others were already suggested before. The first enhancement was to add a time contract; so that it limits the exhaustive enumeration process to happen only during the contracted time. The value of the time contract is configured using a new added parameter. The second enhancement was to train a Rotation Forest classifier on the transformed data instead of the mixed ensemble.

In the beginning of the training process, ST explores the number of shapelets that exist in the data set during the contract time, then approximates the number of shapelet samples it can acquire from each time series. Based on the results, it modifies the time of the model using a reinforcement learning technique. For the shapelets sampling process a random search is executed. If the time contract is long enough, then all shapelets are considered. Then, an evaluation process starts, where all the acquired samples are compared and filtered; to exclude redundant shapelets. The filtered samples are then evaluated using information gain (IG); to keep the best distinctive shapelets of the classes. For multi-class problems, a one vs all evaluation is carried out as per the suggestion by [BB17b]. After that the weak shapelets are excluded and the remaining shapelets are merged into a pool and the transformation process starts. Finally, the rotation forest classifier is trained. In our framework, we use a contractable version of ST instead of the original. The implemented version allows for specifying a contrat time and uses a random forest classifier.

### 2.3.4 Dictionary Based Algorithms

As previously discussed, shapelet algorithms are very good at identifying unique sub-series which can be related to class labels either by their presence or absence. This technique although very helpful for some problems, but it can deceived if the class labels are defined by the relative frequency of some patterns and not just the a single match [Lar+18; Bag+17]. Dictionary based algorithms is a family of algorithms that uses the frequency of words as the foundation of finding discriminatory features of classes [MVB19]. They can address this type of problems by tracking the frequency counts of repeating patterns, which makes same resilient to noisy data [Shi+20]. To achieve this, dictionary based algorithms transform time series data into words; to reduce their dimensionality and approximate their values. Then based on the distribution of the words they measure similarity between instances [LKL12; Sch15].

In general, the process of dictionary based classifiers involves using a sliding window, of some length  $w$ , over the data. Then for each window, an approximation process is carried out to reduce its length to from  $w$  to a shorter length  $l$  [Bag+20]. A quantization process then follows; to discretize the real values into words by assigning letters to each value. When all the windows are finished, the number of occurences of words is calcualted and each time series is represented by a histogram [Shi+20]. In the end, similarity between instances is calculated based on similarity between histograms.

Different dictionary based algorithms apply different approximation and discretization techniques. For example two of the older algorithms; Bag of Patterns (BoP) [LKL12] and Symbolic Aggregate Approximation-Vector Space Model (SAX-VSM) [SM13] use an approximation method called Symbolic Aggregate Approximation (SAX) [Lin+07]. SAX extends Piecewise Aggregate Approximation (PAA), a method that divides a time series into equal consecutive segments and concatenates their means [Shi+20]. Then to transform the series to words, it creates break points using quantiles from the normal distribution. In addition to using SAX, SAX-VSM merges it with the vector space model from the Information Retrieval domain. The major difference between BoP and SAX-VSM, is that SAX-VSM calculates its frequencies on the class level rather than on the series level. Then applies weighting using inverse document frequency [Bag+17].

On the other hand, the newer algorithms; BOSS [Sch15], BOSS-VS [Sch16] and WEASEL [SL17a] use another method which is called Symbolic Fourier Approximation (SFA) which approximates data using coefficients of Discrete Fourier Transformation (DFT). Then apply an adaptive technique called Multiple Coefficient Binning (MCB) to quantize the data. These two techniques; SFA and MCB, will be discussed in more details in the next parts.

## BOSS

Bag of SFA Symbols (BOSS) is a dictionary based algorithm that was introduced by [Sch15] in 2015. Like the previous dictionary based classifiers, BoP and SAX-VSM, 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 [Bag+17]. 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 its 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 [Lin+07; LKL12] 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

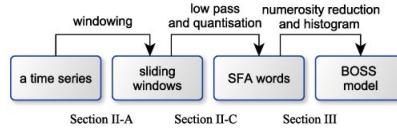


Figure 2.5: The BOSS workflow [Sch15]

subtracting the mean from its original values and then dividing the difference by the standard deviation.

$$\text{windows}(T, w) = \{S_{1:w}, S_{2:w}, \dots, S_{n-w+1:w}\} \quad (2.5)$$

After running the windowing function, the real values of the time points inside the windows are transformed into words using SFA [SH12]. 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 its 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 [LLY17]. 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 it 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 [Sch15] is:

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:

$$D(T_1, T_2) = dist(B_1, B_2) \quad (2.6)$$

where

$$dist(B_1, B_2) = \sum_{a \in B_1; B_1(a) > 0} [B_1(a) - B_2(a)]^2 \quad (2.7)$$

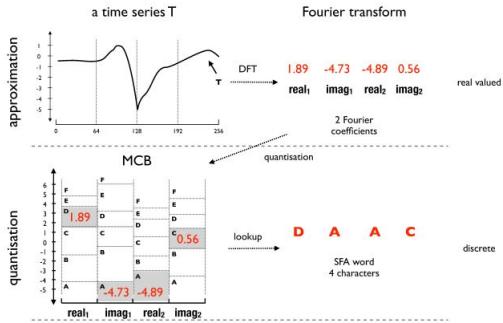


Figure 2.6: A time series is (a) approximated (low pass filtered) using DFT and (b) quantised using MCB resulting in the SFA word DAAC [Sch15]

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.

In [MVB19], some enhancements were introduced to BOSS. These included randomising the parameter settings, subsampling training data for each classifier, adding a voting scheme to assign weights to classifiers based on the probability distribution from cross validation on training dataset, fixing the number of classifiers and adding a time contract.

The new enhanced version of BOSS is named the Contractable BOSS or CBoss. In the contractable BOSS, the ensemble is assigned to build as many individual BOSS

classifiers as possible during a specified duration of time (the contracting time) despite the properties of the data set. In cases of small data sets, this might require huge memory resources; as the number of possible classifiers to create would increase. This was overcome by introducing a parameter to limit the number of maximum classifiers to create. On the other side, for huge data sets there is a chance of building problems to happen and the loss of progress. Which was overcome by adding checkpoints to periodically save the current state of the classifier.

Multiple experiments were carried out using different variations of the enhancements. None of the enhancements alone proved competence to the original BOSS performance, but a combination of them together made it possible to achieve significant results. More details about the experiments can be found in [MVB19]. In our framework, we use CBoss as a replacement for the original BOSS; as recommended by [MVB19] and already replaced in the new improved version of the ensemble HIVE-COTE in [Bag+20].

## 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 [LTB18; Bag+15], or have used it, as a baseline for comparison of their algorithm's performance [Faw+20; Shi+20; Luc+19]. Later on, the same team that developed BOSS introduced a newer algorithm, Word ExtrAction for time SEries cLassification (WEASEL) [SL17a]. 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 [MVB19]. 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 [MVB19].

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, unigrams 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.

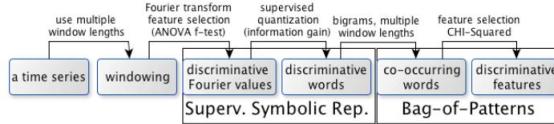


Figure 2.7: WEASEL Pipeline: Feature extraction using a novel supervised symbolic representation, and a novel bag-of-patterns model [SL17a]

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 over come 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 [Low14]. This can be tested by comparing two variances; the variance between groups and the variance within the groups. Using the notations in [SL17a], *mean square between* ( $MS_B$ ) and *mean square within* ( $MS_W$ ) respectively. The F-values is then calculated as :

$$F = \frac{MS_B}{MS_W} \quad (2.8)$$

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

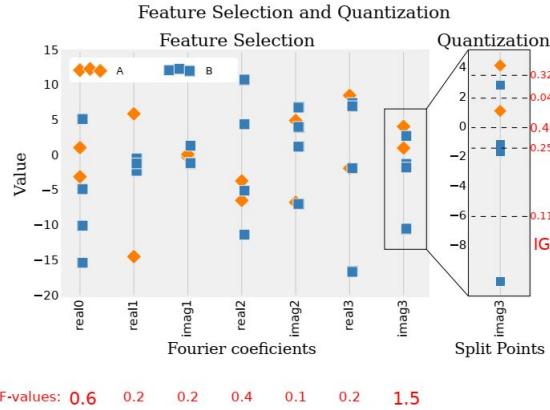


Figure 2.8: On the left: Distribution of Fourier coefficients for a sample dataset. The high F-values on imag3 and real0 (red text at bottom) should be selected to best separate the samples from class labels 'A' and 'B'. On the right: Zoom in on imag3. Information gain splitting is applied to find the best bins for the subsequent quantization. High information gain (IG) indicates pure (good) split points [SL17a]

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 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.

### 2.3.5 Deep Learning Algorithms

Deep learning is a long established group of machine learning algorithms which has proved competence in many areas [LBH15] and which have encouraged the introduction of deep learning algorithms for time series classification [WYO17]. 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 [Shi+20].

The general framework of deep learning neural networks is described by [Faw+19a] 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:

$$f_L(\theta_L, x) = f_{L-1}(\theta_{L-1}, f_{L-2}(\theta_{L-2}, \dots, f_1(\theta_1, x))) \quad (2.9)$$

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). [Faw+19b] 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[LMT16]. Multi-Scale Neural Networks (MCNN) [CCC16] and LeNet [LMT16] 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 [WYO17] 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 [SL17a].

A recent more comprehensive review of deep learning algorithms was done by [Faw+19a]. 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, FCN and ResNet. The other algorithms included were Encoder [SPK18], Multi Channel Deep Convolutional Neural Network (MCDCCNN) [Zhe+14; Zhe+16], Time Convolutional Neural Network [Zha+17] and Time Warping Invariant Echo State Network (TWIESN) [TH16]. For more details about the algorithms and their structures, please refer to [Faw+19a].

Finally, the newest deep learning algorithm introduced for time series classification is InceptionTime [Faw+20], 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.

### InceptionTime

InceptionTime [Faw+20] is a recent deep learning algorithm for TSC problems, which is able to achieve high accuracy score [Rui+20]. 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 [KSH12], an algorithm that was considered a breakthrough for deep learning algorithms [Alo+18], 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 [Faw+19a]. 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 were introduced by [Sze+15] and evolved later on to Inceptionv4 [Sze+17]. The InceptionTime classifier is an ensemble of 5 InceptionTime networks that are 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 [Faw+20].

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 [He+16]. 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 in 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, it is then followed out 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.

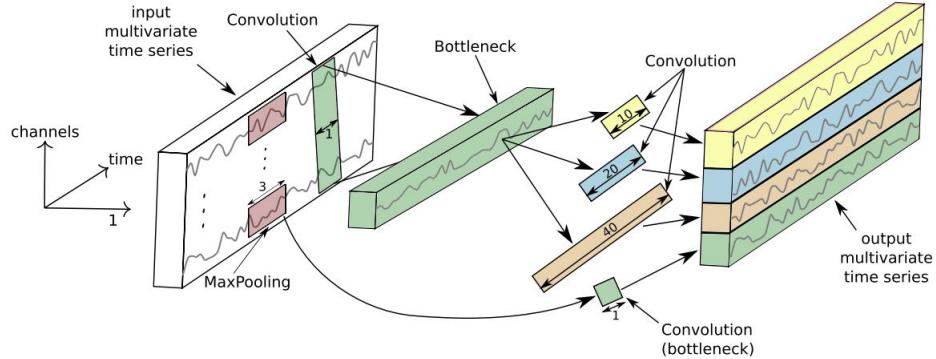


Figure 2.9: Inside Inception module for time series classification. For simplicity a bottleneck layer of size  $m = 1$  is used [Faw+20]

InceptionTime networks accuracy scores showed high variances, a problem which have been discussed by [SW17] and was found in ResNet networks as well [Faw+19b]. 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 [Faw+19b] 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

$$\hat{y}_{i,c} = \frac{1}{n} \sum_{j=1}^n \sigma_c(x_i, \theta_j) | \forall c \in [1, C] \quad (2.10)$$

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  $n$ .

## 2.4 Early Time Series Classification

On another side, Early Time Series Classification (eTSC) is also a classification problem which considers the temporal nature of data, but with a focus on slightly different objectives than TSC.

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 compared to a user defined threshold[XPP09]. 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 [Mor+19; TM16; XPP12; Mor+17b]. This is why eTSC is considered as a problem of optimizing multiple objectives.

We have already defined accuracy for TSC in section 2.2. For earliness, we follow the definition mentioned by [SL20]; as the mean number of data points  $s$  after which a label is assigned.

**Definition 6**

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

There are multiple ways to compare algorithms for eTSC based on the two objectives accuracy and earliness. For instance, one could fix the value of earliness and compare accuracies of algorithms at a defined early point of time, fix a certain accuracy value to be achieved and compare how early classifiers can reach it, or combine both accuracy and earliness in one score [SL20]. The  $F_\beta - \text{measure}$  is an evaluation measure which combines both accuracy and earliness in one equation. It is the calculation of the weighted average between earliness and accuracy. The value of  $\beta$  is the weighting variable; it can be used to give higher importance to either of the aspects over the other.  $F_\beta - \text{measure}$  definition is:

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

A special case of the  $F_\beta - \text{measure}$  is the Harmonic mean (HM) or the  $F_1$  score. The HM assigns equal weights for earliness and accuracy by assigning the value of  $\beta$  to 1. It is a popular choice of evaluation for eTSC problems and has been previously used by [GO12] and [SL20]. For the scope of our framework, we will also consider the weighted average  $F_\beta - \text{measure}$  to evaluate our classifiers.

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 [Mor+17a; Par+13; Lin+15]. This is why eTSC has been applied in various domains like early medical diagnosis [GM01; GO12], avoiding issues in network traffic flow [Ber+06], human activity recognition [YD19; Gup+20] and early prediction of stock crisis [GRO14].

## 2.5 Early Time Series Classification Algorithms

### 2.5.1 ECTS

The idea of eTSC evolved over time. In the beginning, similar concepts to eTSC existed where the main goal was to classify prefixes of data, but these techniques lacked the reasoning of how to provide reliable classification results [Lv+19; SK16]. One of the earliest papers that formally defined eTSC as a tradeoff between earliness and accuracy was the paper by [XPP09], proposing a technique which is named after the problem itself; Early Time Series Classification (ECTS). The main idea was to examine the stability of Nearest Neighbor approaches, to keep correctly classifying data instances, in the space of full length series and the subspaces created from subsequences of the training data set [Lv+19; Mor+17a]. ECTS utilizes the concept of reverse nearest neighbors (RNN) to learn for each training instance the Minimum Prediction Length (MPL). MPL represents the earliest point in time for which a training instance can be used to classify a query instance without loss of accuracy. The definition mentioned by [XPP09] defines MPL as:

**Definition 8** *In a training data set  $T$  with full-length  $L$ , for a time series  $t \in T$ , the minimum prediction length (MPL for short) of  $t$ ,  $MPL(t) = k$  if for any  $l$  ( $k \leq l \leq L$ ), (1)  $RNN^l(t) = RNN^L(t) \neq \emptyset$  and (2)  $RNN^{k-1}(t) \neq RNN^L(t)$ . Specifically, if  $RNN^L(t) = \emptyset$ , then  $MPL(t) = L$ .*

During classification MPL is calculated for each instance (in the same study a relaxed variant is introduced using clusters of instances). As for the classification, as the data

points of the testing instance start to arrive, the classification method tries to find the nearest instance which has reached its MPL. If there was no suh instance found, then more data points are needed till a trusted classification can be done. ECTS was criticized [He+15; GRO14; Xin+11] because it is based on nearest neighbor approaches; which only provide classification results based on distance without deriving patterns from the data. This means that data users cannot attain insights about the problem using the results.

### 2.5.2 Shapelet Based Algorithms

A group of methods [Xin+11; GO12; GRO14; He+15] focused on the interpretability of the results obtained by eTSCAs, motivated by domains where not only correct classification is important, but also getting insights and understanding what contributed to the final result. One of these domains would be the medical field, where undersrstanding the reason for a patient's disease is very important than just linking the patient to previously seen patients [Xin+11]. Shapelets provided a good solution for these meth-ods. Since shapelets, as previously mentioned, are defined as unique subsequences which can uniquely identify classes. Early Distinctive Shapelet Classification (EDSC) was a framework consisting of 2 steps; feature extraction and feature seleton. In the feature extraction step, all local shapelets are extracted and a distance threshold is learned. EDSC developed a technique which they called the best-matching-distance (BMD), which consideres a shapelet to be distinctive, by comparing it to all instances in the training data set and considering target classes distribution. If the majority of instances closer to the shapelet belong to a target class, it is considered to be distinctive of this class. In the feature selection step, EDSC selects a smaller set of features based on three main criteria; earliness, frequency and distinctivness. EDSC was later extended in two different directions.

In [GO12] EDSC was extended in the aspect of data dimensionality; as the original EDSC only convered univariate data sets. The extension covered; adapting the feature extraction to discover multivariate shapelets along with calculating distance thresholds for each dimension and using a distance threshold based on information gain, adapting the feature selection step, by modifying the equations calculating the notions of earliness, frequency and distinctivness to multiple dimensions. In [GRO14] EDSC was extended in the aspect of data uncertainty. The main idea was that shapelets use a distance threshold when being compared to instances, if the distance between the shapelet and the instance is less than the threshold, the instance is given the label of the shapelet. On the other hand, since the threshold represents a range, a shapelet is more certain of the class label given to an instance which has a distance closer to 0 than the class label given to an instance falling on the limits of the threshold. Motivated by the interest of uncertainty in the medical field, the goal was to provide a notion of uncertainty along with the interpretability of shapelets. The proposed Modified EDSC with Uncertainty estimates (MEDSC-U) calculated uncertainty as  $(1 - \text{confidence})$ , where confidence is a numeric value represting two aspects; confidence that the distance between the shapelet and the instance is less than the threshold and confidence that the shapelet can correctly classify the instance. This confidence notion could then be generalized for classes represented by multiple shapelets. The issue of confidence for eTSC had been discussed in other papers like [Par+13] and [Lv+19].

### 2.5.3 ECDIRE

The paper by [Mor+17b] proposed a new framework; Early Classification of Time Series based on Discriminating Classes Over Time (ECDIRE). There were two goals that ECDIRE wanted to achieve [Mor+17a]. The main goal of the framework was to track the accuracy evolution for a group of probabilistic classifiers across time; to identify safe time points after which predictions can be made. This would allow for avoiding unnecessary checking at each time point. The other goal was to offer an outliers discarding mechanism using a reliability condition. ECDIRE consisted of three training steps. In the first step the framework analyzes the data set and tries to identify the safe points at which classes can be differentiated from each other. This is done by defining time points as percentages of the full length and training, at each time point, a probabilistic classifier. The early points are then defined as those which maintain an accuracy of a desired percentage from that obtained on the complete data [SK16]. In the second step, a reliability threshold is determined using the class probabilities of the assigned class and the second highest class; this helps avoid uncertain classifications or very soon time points. In the last step, the final probabilistic classifiers are trained using the early time points identified in the first step. As a backup plan, if the last early time point did not coincide with the full length of the time series, ECDIRE acts like other methods and builds classifiers on each time point.

### 2.5.4 TEASER

One of the latest classifiers introduced is the Two-tier Early and Accurate Series clasifiER (TEASER) [SL20]. TEASER's main focus was challenging the logic of previous classifiers on determining the point of time at which a classification can be assigned to a time series instance. Previous methods would try to learn from the training data set a fixed prefix needed to be seen before it can start classification; like MPL in ECTS and EDSC, or to learn a safe point like in ECDIRE. These techniques would require the classifiers to always wait till they reach these fixed points of time before they can assign a label. Whereas TEASER was designed to give a classification once confident of its decision, and not to depend on any fixed time points nor assume that all data instances were recorded at the same starting point. TEASER models the problem of eTSC with a two tier solution. In the first tier, a slave classifier is regularly checking the incoming data and assigning class probabilities. These probabilities are then assessed in the second tier by a master classifier, which either accepts the classification from the slave, or rejects it and decides to wait for more data to arrive. The classifiers are paired through out the whole process, such that for each slave classifier there is a master classifier that reads its output and assesses its results. TEASER starts off by defining the number of classifier pairs that will be needed, this is calculated by dividing the length of the longest data instance in the training data set by the window length which is a user defined parameter. Which is noted as  $S = [n_{max}/w]$ , where  $n_{max}$  is the length of the longest instance and  $w$  is the window length. During training, slave classifiers are trained on subsequences  $s_i$ , or snapshots, of the time series where the length of  $s_i = i.w$ . These snapshots are z-normalized based on the values existing only in them and not on values from the full length data as recommended by [Mor+17a].

The slave classifiers would then calculate three features; the class label with the highest probability  $c(s_i)$ , the vector of class probabilities for all available classes  $P(s_i)$ ,

and the difference between the class probability of the highest two classes denoted as  $\Delta d$ . These three features are then passed to the paired master classifier, which trains a model over them using a one-class Support Vector Machine (oc-SVM). A parameter that plays an important role in this process is  $v$ , it is the number of consecutive times the same class label is given to an instance. Once the same label is given  $v$  times for consecutive snapshots, the master classifier is confident of this class label. Figure 2.10 shows an example of how TEASER components handle a misleading data instance from the traces of microwave and a digital receiver.

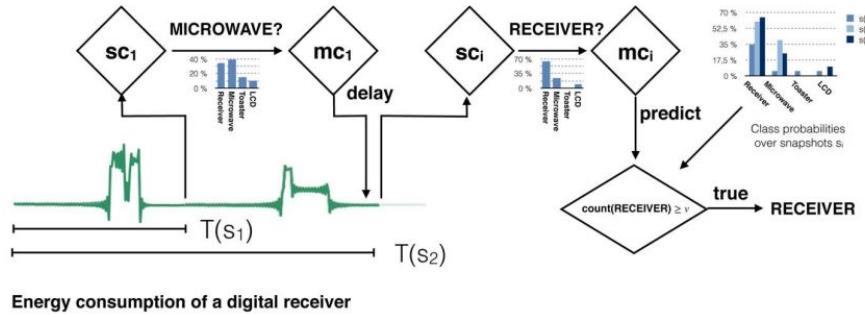


Figure 2.10: TEASER is given a snapshot of an energy consumption time series. After seeing the first  $s$  measurements, the first slave classifier  $sc_1$  performs a prediction which the master classifier  $mc_1$  rejects due to low class probabilities. After observing the  $i$ -th interval which includes a characteristic energy burst, the slave classifier  $sc_i$  (correctly) predicts RECEIVER, and the master classifier  $mc_i$  eventually accepts this prediction. When the prediction of RECEIVER has been consistently derived  $v$  times, it is output as final prediction [SL20]

# Chapter 3

## Related Work on Comparing Time Series Classification Algorithms

The main contribution of this thesis is to provide a framework that compares TSCAs in an early classification context. The idea of creating frameworks to standardize the process of comparison between algorithms is not new to the time series data domain. It was actually needed due to the increasing interest in time series data and the vast amount of newly introduced algorithms. In this chapter we will briefly discuss 3 other frameworks that had been developed to compare TSCA on univariate and multivariate data sets.

### 3.1 The great time series classification bake off

One of the most famous studies in comparing TSCAs is [Bag+17] and which has inspired some of the work done in this thesis. At the time when this paper was published, the TSC domain had already witnessed hundreds of algorithms being proposed, each claiming superior performance. There was no clear structure, at the time, for experimenting and benchmarking algorithms' performance, which pushed for having a more concrete structure to be followed in succeeding research. The primary goal of the experiment was to assess the average accuracy of the different classifiers on a broad set of univariate data sets. The scalability and efficiency of the classifiers were considered as secondary goals.

#### 3.1.1 Inclusion Criteria

The experiment collected 18 different classifiers to be compared with 2 baseline classifiers on the UCR data archive (discussed in 5.1.1). Algorithms were selected based on three criteria; the algorithm had to be published in a high impact conference or journal, it had to be previously assessed using some data sets from the UCR archive and the code for the algorithm had to be made publicly available.

A taxonomy was introduced by [Bag+17] which grouped the algorithms together based on the techniques they exploit. Each group of algorithms share a similar way for extracting features that they use to discriminate between classes of problems. This taxonomy helps describe the various techniques applied for solving TSC problems and understand for which problem types the techniques best apply. These groups were; whole series similarity (also referred to as distance-based algorithms), phase dependent intervals, phase independent intervals (shapelets), dictionary based and ensemble algorithms. We follow the same taxonomy for our experiment as well and describe these groups in more details in section 2.3.

### 3.1.2 Experimental Design

The framework ran 100 resampling folds for each of the data sets for each classifier, totalling 8500 (85 data sets  $\times$  100 resamples) experiments for each classifier. The first run used the default train/test split provided with the data set, while the remaining resamples were stratified to preserve the class distribution in the default split. Using the same parameters as the relevant published papers, the experiment limited the number of parameter values for each classifier to a maximum of 100. Also different number of cross validation folds were applied based on the type of algorithm. LOOCV was used for algorithms utilizing distance based measures, while for the other classifiers a new model was created for each parameter set. Resulting in a total of 850,000 models per classifier.

## 3.2 Deep learning for time series classification: a review

Another paper that focused on comparison of deep learning algorithms for TSC was [Faw+19a]. Deep learning algorithms had already proved competency in many application fields including image classification, speech recognition and natural language processing [He+16; SK16; KSH12; Gua+19] and was gaining more popularity in TSC problems [Zhe+14; Zhe+16; Zha+17]. The main goals were to benchmark the performance of deep learning algorithms on TSC problems; as they were less studied than the other algorithms and to provide an open source framework for deep learning on TSC.

### 3.2.1 Inclusion Criteria

The framework included 9 discriminative model end-to-end deep learning algorithms developed to work on TSC problems. A taxonomy for the deep learning algorithms was also introduced which divided the algorithms into two main groups; generative and discriminative approaches.

Generative models involve an unsupervised step which tries to learn a quality representation of the time series, then this new representation is fed to another classifier in a learning phase. But this family of algorithms was excluded due to its complexity and incompetence compared to the other group [LGI17; Bag+17].

For the discriminative models, algorithms that needed preprocessing of features prior to learning were excluded; to avoid the bias of hand crafted features. The remaining discriminative approaches were; Multi Layer Perceptron (MLP), Fully Convolutional Neural Network (FCN), Residual Network (ResNet), Encoder, Multi-scale

Convolutional Neural Network (MCNN), Time Le-Net, Multi Channel Deep Convolutional Neural Network (MCDCCNN), Time Convolutional Neural Network (Time-CNN) and Time Warping Invariant Echo State Network (TWIESN).

These 9 algorithms were compared to each other and to 7 other classifiers. A group consisting of the best 4, out of the 18, classifiers that were included in the experiment by [Bag+17] which are; Elastic Ensemble (EE), Bag Of SFA Symbols (BOSS), Shapelet Transform (ST), Collective of Transformation-based Ensemble(COTE). Hierarchical Vote Collective of Transformation-Based Ensembles (HIVE-COTE) was also included, which is an extension of COTE using hierarchical voting and two extra classifiers. In addition to Proximity Forest, an ensemble based on the same concept of Random Forests but using class exemplars instead of feature split values. Finally, the classic 1-NN classifier utilizing Dynamic Time warping (DTW) elastic distance measure with the warping window set through cross-validation on training data set. We covered many of these classifiers on in section 2.3.

### 3.2.2 Experimental Design

The experiment covered both univariate data sets and multivariate data sets. The univariate data sets were represented by the same 85 z-normalized data sets from the UCR archive as [Bag+17]. As for the multivariate data sets, 12 data sets from the archive by Mustafa Baydogan<sup>1</sup> were used. Due to the existence of instances with unequal lengths in the multivariate archive, the linear interpolation method suggested by [RK05] was used; so that all instances inside the same data set were adjusted to the length of the longest instance. Z-normalization was not applied to any of the Baydogan archive.

In order to avoid bias from the initial weights assigned to the classifiers, for each of the data sets for each classifier, the framework used 10 runs for training and then considered the mean accuracy of all runs together. The original train/test split was used in all 10 runs, but the initial weights were randomized. The framework applied optimization for the hyperparameters of the deep learning algorithms, but not for the other classifiers. With the exception of TWIESN, the number of epochs used during optimization ranged between 100 and 5000. A model checkpoint procedure was involved. This meant that after training a model for 1000 epochs, the model which attains the least error on the validation data set is chosen for evaluation. All the models were initialized randomly using the methodology from [GB10] and were optimized using variants of stochastic gradient descent; like Adam [KB14] and AdaDelta [Zei12]. For FCN, MLP, and ResNet if the training loss had not improved for 50 consecutive epochs, then learning rate was decreased by 0.5 to a minimum of 0.0001.

## 3.3 The great multivariate time series classification bake off

A recent study was carried out by [Rui+20] focusing on comparing multivariate TSCAs. With previous research paying more attention to univariate classifiers and univariate problems, multivariate time series classification (MTSC) received less attention. Which meant that the MTSC domain is now in the same position that univariate TSC

---

<sup>1</sup><http://www.mustafabaydogan.com/multivariate-time-series-discretization-forclassification.html>

had been in some years ago and there is a need for a benchmarking framework to guide upcoming research on how to compare multivariate TSCAs to already existing ones.

MTSC is simply the classification of time series data collected, where multiple features are collected at each time point and a single label is assigned. MTSC is more challenging than univariate TSC because the discriminative features can be in the interaction between the multiple dimensions and not only based on the interdependency of consecutive values in one dimension. The framework included two techniques that algorithms follow to deal with multivariate data; either dedicated multivariate TSCAs which are, by design, able to handle multiple dimensions, or adapted univariate classifiers which can handle multivariate time series instances. The adapted univariate classifiers technique is based on two pillars; assuming independence between the dimensions of the data and fitting separate univariate classifiers to each dimension then ensemble their results. We used this technique in our experiment for handling multivariate data sets.

### 3.3.1 Inclusion Criteria

Although univariate and multivariate TSC problems differ in their complexity, but the algorithms for both can be grouped, based on the technique they apply, in the same way. These are; distance based similarity, phase dependent intervals, phase independent intervals (shapelets), dictionary based and deep learning algorithms.

The framework collected 16 classifiers from these different groups and compared them on the 30 data sets from the UEA multivariate data archive (discussed in 5.1.2). The included classifiers had to meet 2 simple criteria; code availability and runnable code. The classifiers encompassed 3 variations of the classic DTW which represented the baseline; DTW using Independent Warping ( $DTW_I$ ) which calculates DTW distance for each dimension separately then sums all distances, DTW using Dependent Warping ( $DTW_D$ ) which simultaneously calculates the DTW distance across all dimensions for each time point and DTW using Adaptive Warping ( $DTW_A$ ) which adapts to each instance by calculating both  $DTW_I$  and  $DTW_D$  and then chooses between them based on a threshold score. The multivariate classifiers were; Generalized random shapelet forest (gRFS), WEASEL+MUSE, Canonical interval forest (CIF), Random Convolutional Kernel Transform (ROCKET), Multiple Representation Sequence Learner (MrSEQL), Residual network (ResNet), InceptionTime, Time series attentional prototype network (TapNet). The adapted univariate classifiers were; the ensemble HIVE-COTE and each of its single components; Shapelet Transform (ST), Time Series Forest (TSF), Contractable Bag of Symbolic Fourier Approximation Symbols (CBOSS) and Random Interval Spectral Ensemble (RISE). Two specialized toolkits in time series machine learning were used for the experiment. One of them is python based, while the other is java based. Most of the algorithms are implemented in both, but few exist only in one. More information is available in [Rui+20].

### 3.3.2 Experimental Design

The framework ran 30 resamples for each data set for each classifier to calculate performance metrics, the first run using the default train/test split of the data set, while the remaining 29 runs using stratified sampling; to maintain the class distribution of the original split. In contrast to the previous experiments where only accuracy was measured, this experiment had other measures which are; area under the ROC,

balanced accuracy, F1, negative log likelihood, Matthew’s correlation coefficient, recall/sensitivity, precision and specificity. Still accuracy was the primary performance measure used in comparisons; due to its easier interpretation.

The classifiers were initialized using the default structure and hyperparameters in their original published papers. Apart from the classifiers that have internal tuning processes, no other external tunings were applied. Even for the window size of DTW, which was proved to have a small but significant improvement of performance for univariate data sets [RK05]. Using a side experiment, where an untuned DTW was compared to a naive implementation of DTW with window sizes between 0% and 100% on 21 of the data sets. The experiment had proved that the untuned DTW was better on 14 data sets, but with no significant difference between both. Which meant that tuning was not needed as far as the experiment is concerned.

The data sets of the UEA archive are not normalized and were presented to the classifiers without any preprocessing. This might have caused a disadvantage for WEASEL+MUSE, which doesn’t have an internal normalization process, against other classifiers like, ROCKET; gRSF; TapNet; InceptionTime; ResNet; CBOSS; STC and CIF, which do. The decision of excluding data normalization was based on another side experiment that they carried out. A comparison was done using the three DTW variants, once on normalized data and once on non-normalized data. The experiment had proved that data normalization had no significant effect on the performance of the classifiers, but overall the performance of DTW declined with normalization. The same experiment was repeated for HIVE-COTE and its components and the same result was observed. This meant that data normalization was also not considered important for the scope of the experiment.

### 3.4 Evaluation Criteria

All three frameworks follow the same methodology, introduced by [Dem06], for comparing their classifiers over multiple data sets. When using multiple resamples over data sets, the score of a classifier on a given data set is calculated as the average of its accuracy over all resamples. The methodology recommends using a Friedman ranking test to refute the null hypothesis; that there is no difference between the ranks of the classifiers. After refuting the null hypothesis and applying ranks to each classifier, a pairwise post-hoc analysis is done according to the recommendation of [BCM16]. Which applies a Wilcoxon signed rank test using a Holm’s alpha correction of value 5% to form cliques of classifiers. Each clique represents a group of classifiers where there is no significance between their pairwise comparisons. Results are described by a critical difference diagram, where the average rank of each classifier is noted and cliques are expressed by a thick line.

## Chapter 4

# Framework for Comparing Time Series Classification Algorithms in Early Classification Context

This chapter describes our proposed framework which compares TSCAs in a context that is inspired by the problem of eTSC. We begin by defining what we mean by early classification context and how our framework design simulate such context for TSCAs to operate in. Then we discuss the conceptual structure of the framework, our implementation for its different components and how to evaluate it. In this chapter we will be answering our research questions:

- How to adapt existing TSC workflows for the early TSC context ?
- How to evaluate classifiers in the early TSC context ?
- How to evaluate the proposed solution ?

### 4.1 Early Classification Context

As motivated earlier, eTSC is a field which is concerned with classification of time series data with earliness and accuracy as the main objectives. Dedicated eTSCAs focus on building models that can learn class labels of the data as early as possible while maintaining accuracy [Mor+17a]. Our framework, on the other hand, investigates the adaptation of the context in which conventional TSCA operate to learn a specific classification problem and the effect it has on the classifiers' performance.

We define our early classification context as the problem in which a conventional TSCA is trained on labeled but incomplete time series data. Then this trained model is used to classify unseen data instances of full length. This context helps investigate if TSCA are able to learn about the problem at hand without having to wait for the full data to be made available. In other words, the benefit acquired from waiting for

more data points to be made available for learning is less than the benefit of having a slightly less accurate result but at an earlier point of time.

## 4.2 Conceptual Design

To answer our first research question *How can we adapt existing TS classification workflow for the early TS scenario?* We propose a framework that simulates an early classification context by shortening the length of the data before feeding it to TSCAs to learn on. These classifiers are then evaluated by comparing their performance to a baseline represented by their performance on the full length data. In the beginning, the framework divides whole timeseries data instances into smaller subsequences using a chopping algorithm, and then starts multiple learning processes.

For the first learning process, it feeds only one subsequence for the classifier to learn on. Then in each following learning process, the framework adds the next subsequence to the training data in addition to the previous ones. These processes simulate an incremental learning scenario, where the classifier is not exposed to all the data points collected from the beginning to the end at once, but rather exposed to the data in a sequential procedure. During each learning process the classifier gets to learn more information about the data instances than the previous process. Finally, to assess how good the performance of the classifiers that were exposed to less data points is, we compare their performance to that of a baseline. The baseline for each classifier is determined by its performance when exposed to full length training data.

Our framework design can be divided into two main components; the comparison testbed and the recommender. Each of these components consists of smaller elements which are responsible for specialized tasks which we will discuss in more details. Figure 4.1 shows the conceptual design of the framework and all its constituent components.

### 4.2.1 The Testbed

The first component is the comparison testbed. It is the main engine of the framework for configuring and executing learning processes. There are multiple units that make up the testbed. These are; the testbed management console, the model preparation agent, the chopper and the experiment executor. Each of these units is responsible for a specific task and communicate information either to the other units in the testbed or to other units in the recommender.

#### The Testbed Management Console

In order for the testbed to be able to configure its components and experiments, it needs to collect information about the desired shape of experiments. The testbed management console is the interface which the user interacts with to provide all the necessary information. The console allows the user to specify details about the data sets to be used, the configuration of the chopping algorithm and the structure of the experiments to be executed in terms of models specification, training and testing processes design. The testbed management console creates the central configuration which will be followed by all the other elements.

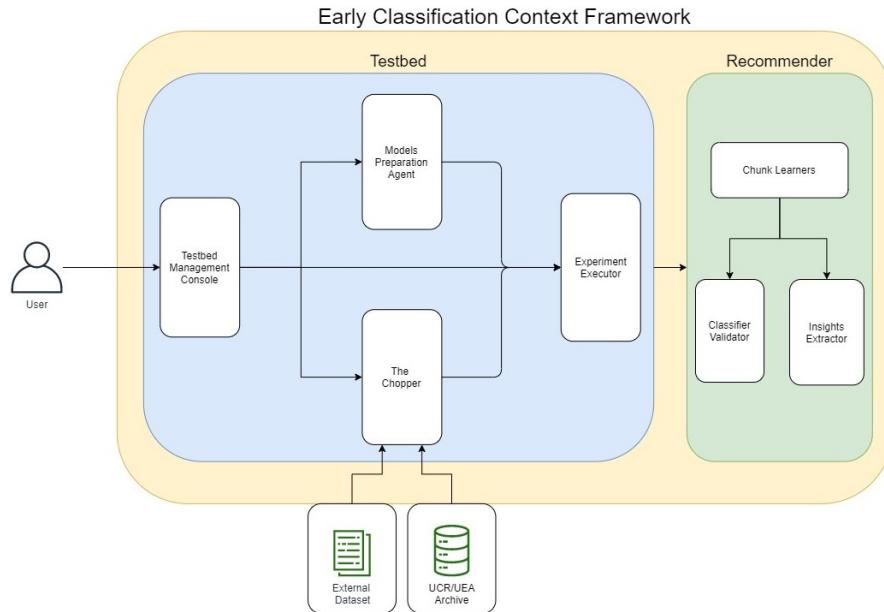


Figure 4.1: Conceptual design of the proposed framework

### The Model Preparation Agent

The model preparation agent is responsible for constructing and composing the classification models that will be used in the learning processes. Based on the configuration from the testbed management console, the model preparation agent includes or excludes models from the experiments and supply the hyperparameter space for each model. The configured models are then passed to the experiment executor.

### The Chopper

The chopper algorithm is the essence of the early classification context. It is the component which simulates the incremental learning scenario for the classifiers. Based on the configuration from the testbed, the chopper collects information about the different data sets that will be used in the learning processes. Then for each data set, it applies the chopping on its data instances based on their length. The chopped data would then be combined with the algorithm; to start the learning processes.

### The Experiment Executor

The last element in the testbed is the experiment executor. It collects the configured models from the model preparation agent along with the chopped data sets from the chopper for the actual execution of the experiments. Based on the central configuration, it configures the environment for the learning processes to provide unbiased results, as well as sets the evaluation scheme of the testing process. After the experiments are finished, it provides insights from the results for each process. It also communicates the results and analysis for the recommender to learn on.

### 4.2.2 The Recommender

The other component of the framework is the recommender. As the name suggests, the recommender is responsible for providing reliable suggestions about the most suitable classifiers to use for unseen data sets. Based on results from the testbed experiments, it learns to predict the performance of different classifiers on data sets in the early classification context. Then when provided with an unseen data set, the recommender would suggest which classifier or multiple classifiers will attain good results if used to learn on incomplete subsequences of data. The recommender consists of three elements; the chunk learners, the classifier validator and the insights extractor.

#### The Chunk Learners

The chunk learners are a group of classifiers which learn on the results from the experiments of the testbed. The word chunk in the name refers to the chunks (subsequences) that were shown for the classifiers to learn on. Each learner is responsible for learning on the performance of classifiers for a given chunk and consequently predicting the performance of these classifiers for unseen data sets. The prediction results are then passed to the classifier validator and the insights extractor.

#### The Classifier Validator

The classifier validator is a component which assesses the prediction results from the chunk learners. Its main objective is to decide if the predicted performance results for the classifiers on a data set qualify as good recommendations or not. The classifier validator should exclude classifiers from the recommendations if it evaluates their performances as weak and not competent. The result from the classifier validator is the actual final recommendations of the framework.

#### The Insights Extractor

The insights extractor is responsible for analyzing and extracting insights from the results of the chunk learners. It reveals hidden patterns and provides understandable insights about how the performance predictions were made and how recommendations were chosen.

## 4.3 Implementation

This section describes our implementation for the framework and all of its components. In the beginning we discuss our design protocol for building the framework. Then we demonstrate our technical design for each of the components, the selected toolkits for building the components and the flow of information through communication channels that exist between the framework's units.

### 4.3.1 Design Protocol

Our design protocol covers the main design aspects in the framework like; the selection criteria of the classifiers to be included and the evaluation criteria for assessing the performance of the classifiers.

There is a wide variety of TSCAs that has been introduced in the last decade. Since our goal is to compare TSCAs in the early classification context, we specify an inclusion criteria which would cover this wide spectrum of classifiers. To represent the currently available techniques; we follow the grouping criteria mentioned in [Bag+17] that breaks down classifiers into 5 main groups based on the technique applied. In addition to a 6<sup>th</sup> group which represents deep learning TSCAs. We represent each of the groups by 2 classifiers; one which is an ensemble and another which is a non-ensemble. The reason we include both ensembles and individual classifiers; is that although for each group there is a classifier which can outperform the others, yet combining classifiers in an ensemble have proven to outperform all of its individual classifiers [Faw+19a].

The other aspect of the framework is to choose an evaluation criteria for assessing the performance of classifiers. Since our framework is inspired by the problem of eTSC; we are interested in evaluating classifiers based on their performance in predicting correct class labels, as well as considering how early they make the prediction. Previous eTSC research have used the harmonic mean of accuracy and earliness as an evaluator for classifiers [SL20]; because it is allows combining both in a single score. The harmonic mean considers both earliness and accuracy to be equally important. In our framework we use  $F_\beta$ ; a generalized form of the harmonic mean which allows the assignment of different weights to earliness and accuracy using the parameter  $\beta$ . We prefer using  $F_\beta$  to offer more flexibility in adjusting weights to suit cases where either of the objectives is more important than the other.

### 4.3.2 Toolkits

We implemented our framework in python language. For the implementation of the non-deep learning algorithms we used the open source libraires *sktimeV.0.5.3* [Lön+19] and *pytsV.0.11.0* [FJ20], while for the implementation of the deep learning algorithm InceptionTime we used the library *sktimedlV.0.2.0* which provides an interface for the implementation provided in the study [Faw+20]. Finally, we used the library *TslearnV.0.5.0.5* [Tav+20] as a mediator between the different model implementations; due to its capacity to process raw data into the needed input format by any of the libraries, as well as converting data format between the different libraires.

### 4.3.3 The Testbed Implementation

#### The Testbed Management Console Implementation

The Testbed Management Console is the tool that allows the user to interact with the framework and configure the experiments and the learning processes which will be the basis of the final recommendations. At this stage, the user provides the necessary information required by the framework to formulate and start its learning experiments. These information provide the guidelines for the other components in the framework like the model preparation agent and the data chopper. Which are essential for the preparation of data sets and models for the learning processes. Table 4.2 shows a summary of each of the needed information, represented as a parameter for the framework.

We implemented the user input interface using *docopt*<sup>1</sup>. A library which offers command-line interfaces through the use of simple arguments and elements. It also allows setting default values for parameters which can be overriden if a value is passed.

---

<sup>1</sup><https://github.com/docopt/>

Parameter	Data Type	Description
Dataset	String List	The name of the data set(s)
Splits	Integer	The number of splits to apply on the data set
Hyperparam	Boolean	Learn hyperparameters for the model or skip
NumIterations	Integer	Number of hyperparameters combinations to sample
ChunksToKeep	Integer List	The indexes of the chunks to use
FromBeg	Boolean	Reveal data chunks from beginning or end
ScoringFunction	String	The strategy for evaluating the model on test data

Table 4.2: Input parameters required for using the framework

To start running the framework, one provides the necessary parameter values through a simple terminal command. Listing 4.1 shows a sample command for running our master python file *run.py* on the Computers data set using hyperparameters optimization, while applying a split into 10 chunks.

```
python run.py --etsc --dataset=Computers --cv --split=10
```

Listing 4.1: Sample command for providing user input to the framework

### The Model Preparation Agent Implementation

The model preparation agent is responsible for preparing the time series classifiers and creating their hyperparameter space for each learning process. We design our model preparation agent to adjust the preparation of models to handle univariate and multivariate data sets.

As mentioned in section 2.3, we follow the grouping criteria by [Bag+17] which arranges TSCAs into 5 main categories based on their techniques. We have included classifiers representing each of the 5 groups, in addition to a 6<sup>th</sup> group which represents the family of deep learning time series algorithms. For each group, we have selected the state of the art classifier based on the results of previous comparison frameworks discussed in chapter 3. We tried, as far as the implemented libraries allowed, to represent each group with 2 classifiers; one that is a non-ensemble and another which is an ensemble. A total of 10 classifiers are included in the framework covering all 6 groups, these can be broken down into the following:

- Distance based : 1NN using MSM distance and PForest
- Phase dependent algorithms : TSF
- Shapelets: LS and ST
- Dictionary based : WEASEL and BOSS
- Deep learning : InceptionTime
- Classical baselines: 1NN using euclidean distance and 1NN using DTW

Table 4.4 shows the configuration and hyperparameters for all the included classifiers. We used the same hyperparameter space like that of the original published papers as closely as possible.

Classifier	Parameters
MSM	$c = \{0.01, 0.1, 1, 10, 100\}$
PForest	$k = 100$ $r = 5$
TSF	$r = 500$ splitting = {entropy, gini} maxfeatures = {sqrt, log2}
LS	$\lambda_w = \{0.01, 0.1, 1\}$ $L_{min} = \{0.025, 0.075, 0.125, 0.175, 0.2\}$ $R = \{1, 2, 3\}$ $K = \{0.05, 0.15, 0.3\}$ $\eta = 0.01$ maxIter = {2000, 5000, 10000}
ST	$t = 60$ mins $n=\{500,100\}$
CBOSS	$t = 60$ mins $k=\{50, 100, 250, 500\}$ $s=250$ $p=[0.5 - 1]$
WEASEL	ANOVA = {True, False} bigrams = {True, False} binning = {equi-depth, equi-width, information-gain}
Inception	epochs = 1500 batch size = 64 $\eta = 0.001$ kernel sizes = {10, 20, 40}
DTW	full warping window

Table 4.4: Parameters and configuration for TSCAs

As mentioned earlier in section 3.3, there are two ways to work on multivariate time series data sets [Rui+20]. The first way is by using bespoke multivariate classifiers which can deal with the multiple dimensions available in the data. While the other is by using ensembles of univariate classifiers, each fit to a separate dimension under the assumption of independence between the dimensions. Our framework applies the second technique for handling multivariate data sets.

The model preparation agent checks the metadata of the data sets that the management console communicated, then it determines the number of dimensions each data set comprises. If the data set is univariate then the model preparation agent initializes the default classifier implementations provided by the libraries. In case the data set was found to be multivariate, it initializes a special classifier type which extends the usage of univariate classifiers for multivariate problems. This special classifier is referred to as *ColumnEnsembleClassifier* in the *sktime* library implementation and as the *MultivariateClassifier* in *pyts*.

There are two assumptions that we make at this step in our framework. First, we assume that to extend a specific univariate classifier to a problem, we should fit the same type of the classifier on all dimensions of the data. This goes back to our goal, that we want to compare classifiers from different groups and not to create an ensemble which makes use of different classifiers applying different techniques, like HIVE-COTE does. Second, we assume that we always have to fit one classifier per dimension, which have proved to be not a feasible solution for high dimensional data sets. We will discuss this in more details in our results.

After the model preparation agent finishes, it passes the initialized models to the experiment executor for the learning processes to actually start.

### The Chopper Implementation

The chopper is the data sets handling module of the framework. It is not only designed to apply chopping of the data based on the configuration, but also to load the data sets and preprocess it. The chopper collects information about the data sets to be used, the number of splits (chops) to be applied, the chopping direction and the chops filtration. These information are passed from the management console in the form of the parameters *Dataset*, *Splits*, *FromBeg* and *ChunksToKeep* respectively.

The first task of the chopper is to process input data. It is configured to process data sets in either of two ways; automatic download from the UCR/UEA data archives, or manually provided data sets that comply with some specifications. Since the data loading module from *pyts* [FJ20] offers integrated tools with the UCR/UEA data archives for fetching and downloading data, we have configured our chopper to extend *pyts* by transparently checking both data archives when provided with a data set name. The other option would be to provide a data set folder holding the name of the data set. Inside the folder there should exist one training data set file with the name “*FolderName\_TRAIN*” and one testing data set file with the name “*FolderName\_TEST*”. Which is the case for most of the data sets from both archives. For multivariate data sets, the chopper assumes that the training data set file, as well as the testing data set file, contains data of all dimensions. For maximum compatibility, the chopper extends the data processing utilities from *sktime* [Lön+19] which can read data sets that comply with the guidelines provided by the UCR and the UEA archives. We consider “arff” format as the primary expected data type.

Based on the assumption that all data instances are of the same length, we implement our early classification context using a time series chopping algorithm. The

algorithm is provided with some time series data set  $T$  with instances of length  $L$  and three user defined parameters. Parameter  $s$  decides the number of splits that the length  $L$  should be split to. *FromBeg* is a flag which indicates whether the data chopping should occur from the beginning or the end of the time series  $T$ . *ChunksToKeep* which provides the user with the option to use only specific chunks of interest and exclude the processing for all other chunks. Algorithm 1 demonstrates how the chopping algorithm works.

The result of the algorithm is then used to create new copies of  $T$ ;  $[T_1, T_2, \dots, T_s]$ , each copy  $T_i$  containing all instances from  $T$  but revealing a subsequence of  $L$  which we call the chunk. A chunk is a sequence of  $p$  data points starting from the beginning or the end of the instances in  $T$  based on the value of the parameter *FromBeg*. Each data set copy  $T_{i+1}$  reveals one more chunk than the previous data set  $T_i$  and the data set  $T_s$  represents the set with instances of the full length  $L$ . It should be noted that our algorithm tries to have an equal value for  $p$  used for each chunk, but this is limited by the length of the time series and the provided number of splits. The resulting data sets are then passed to the training process, where each algorithm applies its own technique to learn on the data set.

**Algorithm 1** The Chopping Algorithm

```

1: function GetSplitIndexes( $T, s, ChunksToKeep$ )
2:    $L \leftarrow ExtractLength(T)$ 
3:    $ChunksSizes \leftarrow GetChunkSizes(L, s)$ 
4:    $SplitIndexes \leftarrow CumSum(ChunksSizes)$ 
5:   if  $KeepChunks$  is not Null then
6:      $FilteredSplitIndexes \leftarrow KeepChunks(SplitIndexes, ChunksToKeep)$ 
7:   else
8:      $FilteredSplitIndexes \leftarrow SplitIndexes$ 
9:   end if
10:  return  $FilteredSplitIndexes$ 
11: end function

```

To better understand how the chopping algorithm works, consider the next example. Let's assume we have a data set  $T$  with instances of length  $L = 100$ , we set the value of  $s = 10$  and  $FromBeg = \text{True}$ . The algorithm starts by calculating the value of  $p$  for each of the required 10 chunks using the function 2.

The resulting list *ChunksSizes* from the algorithm will contain the value 10 for each of the chunks (*ChunksSizes* = [10, 10, 10, 10, 10, 10, 10, 10, 10, 10]); this is because the length of the time series is divisible by the number of splits provided, so it results in perfect splits of the data and thus chunks of equal sizes. If we assume the same scenario again but with a different *s* value like 8, the algorithm will try to provide as equal values of *p* as possible giving back a *ChunksSizes* of values [13, 13, 13, 13, 12, 12, 12, 12].

---

**Algorithm 2** Function to Get Chunks Sizes

---

```

1: function GetChunkSizes(L, s)
2:   ChunksSizes  $\leftarrow$  []
3:   for i  $\leftarrow$  0 to s – 1 do
4:     if i  $<$  L (mod s) then
5:       p = (L  $\div$  s) + 1
6:     else
7:       p = L  $\div$  s
8:     end if
9:     ChunksSizes.append(p)
10:   end for
11:   return ChunksSizes
12: end function

```

---

*SplitIndexes* = [10, 20, 30, 40, 50, 60, 70, 80, 90, 100]. This means that the corresponding data set  $T_1$  will contain all its data instances represented by the subsequences from the 1<sup>st</sup> time point till the 10<sup>th</sup> time points, and the data set  $T_3$  will contain subsequence from the 1<sup>st</sup> time point till the 30<sup>th</sup> time point. When we apply the same translation to the list *ChunksSizes* = [13, 13, 13, 13, 12, 12, 12, 12], it gets translated into *SplitIndexes* = [13, 26, 39, 52, 64, 76, 88, 100].

Finally, the algorithm filters out the list *SplitIndexes*, if the user provided a list of specific chunks to keep using the variable *ChunksToKeep*. The values passed through *ChunksToKeep* corresponds to the indexes of the chunks the user is interested in, all the other chunks are excluded. Note that this process is different from reducing the number of splits; as the number of splits affects the sizes of the chunks being created. While *ChunksToKeep* is used only to filter out the chunks after their sizes are already determined and translated into *SplitIndexes*. The space of possible values for *ChunksToKeep* is between [1, *s*]. The logic of the function is demonstrated in function 3

---

**Algorithm 3** Function to filter out Chunks

---

```

1: function KeepChunks(SplitIndexes, ChunksToKeep)
2:   FilteredSplitIndexes  $\leftarrow$  []
3:   for i  $\in$  ChunksToKeep do
4:     FilteredSplitIndexes.append(SplitIndexes[i])
5:   end for
6:   return FilteredSplitIndexes
7: end function

```

---

If we consider the list *ChunksToKeep* = [1, 3, 5] for our example *SplitIndexes* = [10, 20, 30, 40, 50, 60, 70, 80, 90, 100]. The resulting filtered list will be *FilteredSplitIndexes* = [10, 30, 50], which will result in only 3 copies of  $T$ ;  $T_1$ ,  $T_3$  and  $T_5$ , each represented by its respective subsequence length.

If we have a closer look at the parameter *s*, we can recognise that it contributes to the granularity of the earliness factor that we are calculating. The value of *s* decides the total number of chunks we will use, which subsequently decides also the ratio

represented by each chunk to the total length of the time series  $T$ . That means the greater the value of  $s$ , the smaller the chunk sizes and consequently the more granular results we can test for the algorithms. If we consider our previous example, when we set the value of  $s = 10$  then the length of each subsequence would be close to 10% of the total length. If we decided to change the value of  $s = 20$ , then each chunk will roughly represents 5% of the total length and so on. Increasing the number of splits to be used comes with the cost of time and resources, each extra split adds one extra run per algorithm per data set.

The final result of the chopper is a set of copies of the input data sets. Each copy contains the exact number of instances as the original data set, but the length of its time series instances is decided based on the number of splits that were applied. Some of these copies might be excluded from the results based on the chunks filtration process, while the other remaining ones are passed to the experiment executor to be combined with the prepared models for the actual execution of the experiments.

### **The Experiment Executor Implementation**

The experiment executor is where the actual learning and testing happens. It communicates with the management console for information regarding the learning and testing processes, then sets up the learning environment accordingly. This includes the sampling of hyperparameter space, the number of cross validation folds and the evaluation criteria on which the testing data sets will be evaluated. After the preparation of the chunks for the data sets by the chopper and the initialization of the classifiers is done by the model preparation agent, the experiment executor starts multiple learning processes.

A learning process involves running one classifier on one data set for a specific chunk. During training, the number of hyperparameter sampling is decided by the framework parameter *NumIterations*. Each of these training processes is then followed by a testing process, where performance scores are calculated using the function provided by the parameter *ScoringFunction* from the management console. After the testing processes are finished, it extracts patterns from the testing results and provides analysis on the performance across and within classifiers in the early classification context.

For the training processes, it is possible to choose whether to optimize hyperparameters for the classifiers through sampling and cross validation or to use the default parameters. Our experiments executor uses a randomized search over the hyperparameters space of classifiers. The framework parameter *NumIterations* represents the number of random samples to consider when optimizing the hyperparameters. If the space of the hyperparameters is less than *NumIterations*, then the space is exhausted by applying a grid search. Finally, if the experiment is configured to do hyperparameter optimization, then the classifier version that attained the highest validation score is used for calculating the performance score on the testing data set.

For the testing processes, the experiment executor allows the usage of any chosen performance metric which is supported by *sklearn* library [Ped+11]. The chosen performance metric is utilized for the evaluation during both; the cross validation processes and for calculating performance on the test data set. For our implementation, we considered “balanced accuracy” as our default performance metric; to account for problems where data sets have imbalanced classes distribution.

Our evaluation is not only dependent on the performance of classifiers on testing data sets, but also considers the amount of data that was used for learning. To answer

our research question *how to evaluate the classifiers in the early classification context*? We explain why we choose using the  $F_\beta$ -measure as an evaluation measure for the classifiers and how do we calculate its value.

We previously defined the  $F_\beta$ -measure in section 2.4, it has been used, in the form of  $F_1$ -Score, by [SL20] as the objective function for evaluating their algorithm *TEASER*. The  $F_\beta$ -measure is a popular choice for calculating a score that combines 2 objectives at the same time, which are earliness and performance for eTSC problems. Likewise, we use  $F_\beta$ -measure as the means of evaluating classifiers in our experiments; due to its capacity to incorporate both objectives in one score. Which allows us to compare the trained classifiers in terms of a one score value without worrying about fixing the value of either the earliness or the performance, as well as changing assigning different weights for the objectives if needed.

On the other hand, the  $F_\beta$ -measure equation that was discussed in section 2.4 doesn't fit directly into our evaluation framework. If we try to substitute the results for a classifier trained on the full length data into this equation, we would always get a value of 0. The reason is that the equation is designed to do reverse scoring to the value of earliness, by re-coding its value so that the less seen data points, the higher the contribution of earliness in the equation. This is simply achieved by using  $(1 - \text{earliness})$  instead of using the value of earliness directly. In case of calculating scores at full length, this would yield a value of 0 regardless of the performance of the classifier; since the earliness value at full length is 1, then the equation  $(1 - \text{earliness})$  would always give a 0 score causing the numerator to be 0.

Our modification to the  $F_\beta$ -measure is a two fold procedure. First, we reverse its scale, so that the lower the value the better. Then we reverse score the result again so that it returns to the original scale. We represent our 2 steps modification in equations 4.1 and 4.2.

$$F_\beta - \text{reversed} = \begin{cases} 1 & \text{if earliness} = 1 \wedge \text{performance} = 0 \\ 1 & \text{if earliness} = 0 \wedge \text{performance} = 0 \\ \frac{(1 - \text{performance}) \text{earliness}}{\beta^2(1 - \text{performance}) + \text{earliness}} & \text{otherwise} \end{cases} \quad (4.1)$$

$$F_\beta = 1 - F_\beta - \text{reversed} \quad (4.2)$$

In the first step, we reverse the scale of the original  $F_\beta$ -measure equation by changing its equation. Instead of reverse scoring earliness and using the actual performance score value, we use the actual earliness value and reverse the performance metric score. We name this measure the  $F_\beta$ -reversed, to clarify that it represents a reversed value of the  $F_\beta$  used in other papers. Using  $F_\beta$ -reversed, if we consider the case of classifiers learning at full length, it would solve the problem of attaining a score of 0 and replace it with a real value.

There are 2 special cases to be handled during this step. First, if the classifier doesn't predict anything correctly but at the earliest point of time. Second, if the classifier doesn't predict anything correctly while seeing the full length. These two cases are then assigned the lowest value, which is 1. In the second step, we reverse score the value from the first step again; to math the expected input for generating critical difference diagrams. We feed these equation with the values of the features *Revealed%* and *Testscore*, for the attributes earliness and performance respectively.

After the successful running of the experiments and calculating performance scores using the provided scoring function, we create a single report for each run experiment.

Item	Description
Classifier	The classifier name
Train time	The total training time (CPU Time)
Test time	The total testing time (CPU Time)
Test score	Performance score on testing data set
Params	List of hyperparameters used and their values
Revealed%	Percent of data points used for training from the total length
F-score	$F_\beta$ – measure combining test score and revealed%
Data set	The name of the data set

Table 4.6: Analysis collected for each experiment

We produce one report per classifier per data set per chunk. These reports collect the essential data and statistics that form the basis of our results, they also form the input for the recommender. Table 4.6 shows the structure of the output for the analysis report.

Most of these information are already available before the learning process starts; because they are metadata about the process which is already known like the *Classifier*, *Dataset* and *Revealed%*. While other information are available once the testing process finishes. For example *Traintime* and *Testtime* represent the total duration of the learning and testing processes. The *Testscore* represents performance of the classifier on testing data set. The last measure to calculate is the *F – score*; that's because it needs the testing process to be completely finished and the *Testscore* to be calculated. In the end, all of these information are used to compare the classifiers in the early classification context and used as an input for the recommender.

#### 4.3.4 The Recommender Implementation

One of the goals of our framework is to be able to recommend good performing algorithms to use on unseen data sets. The Recommender is the component which makes the final decision regarding which algorithms are best to use. It actually is responsible for answering our research question ***How to evaluate the proposed solution?***. A recommender which is able to learn how the different classifiers perform on data sets given chopped training data, would be able to suggest which classifier to use for unseen data sets. Our recommender is not a recommender system in the strict sense, but rather a group of learners that use the performance results of the classifiers from the testbed and predict their performance for new data sets. For the recommender learning, we use the results from all data sets, whether they were finished by all classifiers or not.

#### The Chunk Learners Implementation

The chunk learners are a group of regression algorithms that learn for each data set for each chunk for each classifier the *F – score* value that was attained. We learn for each chunk a separate regressor, which will be used to predict the scores of the

classifiers for that specific chunk. For example, the regressor that learns on the first chunk results, will be used for predicting  $F - scores$  on the new data sets assuming that the classifiers were trained using only the first chunk of the data. The regressor that learns on the second chunk results predicts for the results assuming the classifiers were trained using the 2<sup>nd</sup> chunk and so on. When predicting for a new data set, each of the chunk learners predicts the  $F - score$  for each of the 6 TSCA, including the dummy on this data set. Then these results are passed on to the classifier validator for the final prediction and for the insights extractor for getting insights.

For the implementation of the chunk learners, we have chosen random forests as our regression models. A Random Forest can be either a classifier or a regressor. It was first introduced in [Bre01]. A random forest is an ensemble of  $k$  tree predictors, in which every tree is grown using a randomly sampled input vector independent of the other trees. This technique is called Bagging, in which the training data is generated by randomly selecting  $N$  examples with replacement where  $N$  is the original training data set size. During the building of the tree, at each split, only a sample of the features are selected as candidates for the splitting [CPB18]. Since the output is numerical in case of regressor random forests, the mean squared generalization error for each tree is calculated and then the final prediction of the forest is the average error over all  $k$  trees [SSS17]. The first reason for choosing random forests, is using a classifier which can offer insights about the how the predictions were made. Although random forests are considered as black-box algorithms; due to the large number of trees which makes it hard to get insights about the prediction rule. However, they offer interpretable insights about the selection of features in the prediction and an importance measure for each feature [CPB18]. Also random forests are fast learning algorithms, due to the independence between the individual trees, which allows for parallel training. Finally, random forests have been applied in many scientific fields and have proved to be a competent and accurate algorithm [Jog+17; CPB18].

### The Classifier Validator Implementation

The classifier validator is the module that produces the final recommendations. It processes the results from the chunk learners to compare classifiers performances and give the final decision. It applies a logic for deciding which are the good performing classifiers on the data sets for each of the early classification chunks, as well as the full length chunk. We define a good classifier based on its performance in comparison to the dummy classifier on a specific data set. After the calculation of all the results for all data sets is finished, a ranking procedure is carried out for each data set separately. The classifier with the highest  $F - score$  is given the lowest (best) rank on the data set, while the classifier with the lowest  $F - score$  is given the highest (worst) rank. If two, or more, classifiers have equal scores, they are given the minimum rank. For example if the second and third classifiers are equal they are both given the rank 2. Then the ranking of each of the classifiers is compared to that of the dummy classifier. All the classifiers with ranks equal to or more (worse) than the dummy classifier are flagged as bad performing classifiers, while the others are considered good performing classifiers. In the end, the bad performing classifiers are filtered for each data set and the final recommendation is presented.

### The Insights Extractor Implementation

The insights extractor is responsible for providing analysis results and insights about the results of the chunk learners processes. If more than one learning process is started by the chunk learners, the insights extractor reports insights for each process in a separate folder. Then it aggregates the performance metrics scores across all the run experiments and provide a summary report on the overall performance of the recommender. Since we are using random forest regressors, it provides insights about the feature importance for each of the chunks. A specific feature which might be important for learning on the 1<sup>st</sup> chunk data sets, might be less (more) important for another chunk.

# Chapter 5

# Data Sets and Experimental Design

In order to address our research questions, this thesis considered a group of data sets from the UCR data archive [Dau+18] and the UEA data archive [Bag+18]. This Chapter will describe the data sets that were used to conduct our experiments and their general properties. Then we will describe our experimental setup.

## 5.1 UCR/UEA Archives

The University of California Riverside (UCR) and the University of East Anglia (UEA) data archives have been used as references in literature for many experiments on time series problems [AML19; Faw+20; Bag+17; YD19; Rui+20; Faw+19a] and specially for benchmarking the performance of algorithms.

### 5.1.1 UCR Archive

The UCR archive previously had 85 univariate data sets, but in 2018 it was expanded to reach 128 data sets. Despite sharing the property of being univariate, the data sets cover a variety of values for other aspects like; training data size [16, 8926]; testing data size [20, 16800]; number of classes [2, 60]; length of time series [15, 2844] and for some data sets length of instances vary within the data set; and the nature of the data being collected. The archive offers a single predefined train/test split for each of the data sets to facilitate reproducability of results. Finally, the data is z-normalized; to remove offset and scaling. Z-normalization is the transformation of the data to have a mean of zero and one unit of standard deviation.

### 5.1.2 UEA Archive

The other archive we used is the UEA archive for multivariate data sets. Like the UCR archive, it was developed and expanded over time. It started as a small archive of 12 data sets collected by Mustafa Baydogan, then later on, in 2018, was expanded to reach 30 data sets as a collaboration between researchers of UEA and UCR. The data

sets vary in their characteristics; training data size [12, 30000]; testing data size [15, 20000]; number of classes [2, 39]; length of time series [8, 17984]; number of dimensions [2, 1345] and six different groups based on their application area. The archive was processed so that the data instances have equal lengths, instances of missing data points were excluded and a single predefined train/test split is provided.

### 5.1.3 Used Data Sets

For the scope of our experiment, a total of 77 data sets were selected from both archives to be used within the framework. We have based our selection criteria for data sets on:

- The data belongs to a real time series structure behind it and not converted from another data format.
- All instances in the data sets has the same length of time series.
- The data sets should not contain any missing values.

Based on the previously mentioned criteria, we excluded data sets that belonged to either of the types; Simulated, Image or Motion type; to comply with our first selection criteria. For the other two criteria; we avoided any further preprocessing than what is offered by the data archives. As mentioned before, the univariate data sets from the UCR are already Z-normalized. Z-normalization is a common and recommended practice for time series data as it has shown to have a positive impact on the performance of classifiers [Bag+17; Faw+19a]. This form of preprocessing is not T for the multivariate data sets from the UEA archive.

We excluded data sets with missing data and avoided filling these missing values; as we believe that data imputation is a whole another topic. We also avoided any interpolation of data. A technique which has been previously used in [RK05; Faw+19a] to make all instances of the data set of the same length as the longest instance. Our goal in this experiment is to focus on comparing the performance of the classifiers in our proposed context based while ensuring, as much as possible, that the main contributing factor is the technique behind the classifiers. This doesn't mean that preprocessing of time series data is of a less importance, but we simply excluded it because we believe that it doesn't fit into our current scope.

The remaining collection of data sets still represented a wide variety in terms of data criteria. We present the statistics about the remaining data sets all together in figure 5.1. Since the smallest data set (StandWalkJump) and the largest data set (InsectWingbeat) are the same as the UEA archive; training and testing data sizes are the same as previously mentioned ranges in the original archive, [12, 30000] and [15, 20000] respectively. The data sets covers from binary classification problems up to problems with 52 classes, the range [2, 52]. The range of the time series length is [24, 3000], with Chinatown being the shortest and MotorImagery the longest data set. For the 22 multivariate data sets, the number of dimensions still covered the whole range [2, 1345] like the UEA archive. Finally after our exclusion criteria, the remaining 14 problem types are AUDIO, DEVICE, ECG, EEG, EOG, EPG, HAR, HEMO, MISC, OTHER, SENSOR, SOUND, SPECTRO and Traffic.

Although the data sets that we used covered a variety of properties, but within some properties there were huge imbalances in the distribution of values. Within grouping by training data set size, the group of data sets with sizes [501, 1000] was the smallest, with only 4 data sets which is less than half of the next smallest group. Grouping by series length was the best, among the different groupings, in terms of

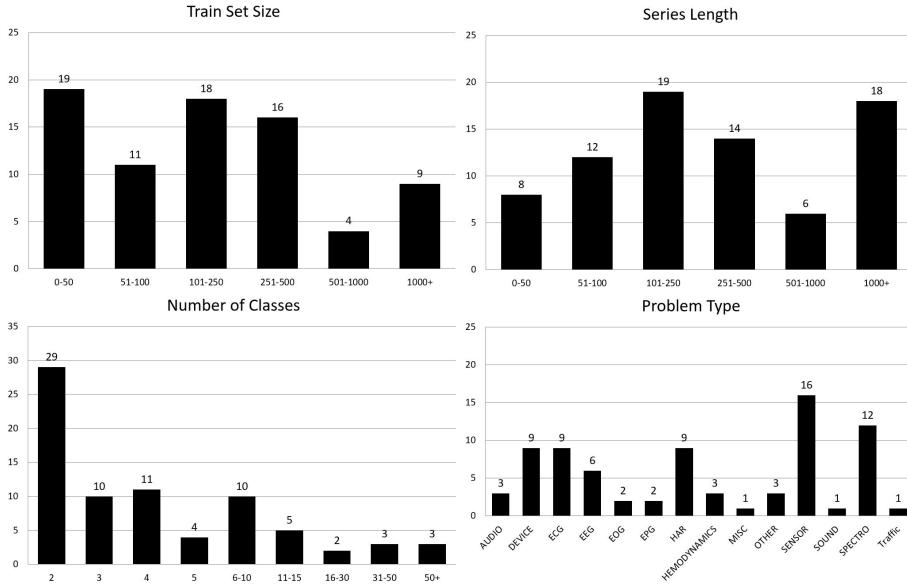


Figure 5.1: Metadata of the used data sets

imbalance. Data sets of lengths [501, 1000] were the smallest group with a count of 6 compared to the largest group, length [101, 250], of 19 data sets. Number of classes was highly imbalanced; the binary classification group contained 29 data sets, which was more than double the second largest group [6, 10] of 10 data sets. This is expected because all the data sets from the UCR archive fall into the same bin, while the multivariate data sets, which are already less in number, are distributed among multiple bins. Finally, there was an over representation of type SENSOR than any other type with a total of 16 data sets. While multiple other groups were represented by very few data sets like MISC, SOUND and TRAFFIC which only had 1 data set; or like EOG and EPG which had 2 data sets.

Since the final goal of our work is to learn a recommender which will be able to suggest the best performing classifier to unseen data sets based on their properties. It is important to include all necessary and useful metadata about the data sets, which will help the recommender correctly give its recommendations. Imbalanced data sets can easily produce misleading results in terms of performance, specially if using a simple performance metrics like accuracy. In cases of extreme imbalance, one can attain very high accuracy scores, simply by always predicting the majority class. Handling the data imbalance is a two fold solution. The first step to help learn about the skewness is to have a definition for imbalance. Then based on this definition flag the data sets as either having balanced classes or not. The second step is to use a suitable performance metric which can give actual performance results and not be mislead by the imbalance.

We discuss the first step here as it is a feature of the data set, the other step we consider in the design of our experiments in section 5.2. To be able to inform our recommender about the skewness property of a data set, we create a new feature, in addition to the previously discussed ones ; to indicate whether a data set is balanced in terms of class distribution or skewed. We define a data set to be skewed if there

exists a class in this data set which is, at least, as twice as likely as any other class.

To better explain how this property works, consider the two data sets;  $DS_1$  and  $DS_2$  each of which contains 5 classes.  $DS_1$  has a class distribution:  $P(C_1) = 20\%$ ,  $P(C_2) = 15\%$ ,  $P(C_3) = 15\%$ ,  $P(C_4) = 25\%$  and  $P(C_5) = 25\%$ . While  $DS_2$  has a class distribution:  $P(C_1) = 20\%$ ,  $P(C_2) = 15\%$ ,  $P(C_3) = 15\%$ ,  $P(C_4) = 30\%$  and  $P(C_5) = 20\%$ . In this case we would flag  $DS_1$  as a balanced data set; because there is no prior class probability which is double any other. On the other hand,  $DS_2$  would be flagged as skewed/imbalanced because class  $C_4$  has a probability of 30% which is double the value of  $C_2$  and  $C_3$ . We denote this feature with the parameter *Balanced*, which is a boolean flag.

Dataset	Train Size	Test Size	Length	#Classes	Type	#Dim	Balanced
ACSF1	100	100	1460	10	DEVICE	1	T
AtrialFibrillation	15	15	640	3	ECG	2	T
BasicMotions	40	40	100	4	HAR	6	T
Beef	30	30	470	5	SPECTRO	1	T
Car	60	60	577	4	SENSOR	1	T
Chinatown	20	345	24	2	Traffic	1	T
CinCECGTorso	40	1380	1639	4	ECG	1	F
Coffee	28	28	286	2	SPECTRO	1	T
Computers	250	250	720	2	DEVICE	1	T
Cricket	108	72	1197	12	HAR	6	T
DuckDuckGeese	60	40	270	5	AUDIO	1345	T
Earthquakes	322	139	512	2	SENSOR	1	F
ECG200	100	100	96	2	ECG	1	F
ECG5000	500	4500	140	5	ECG	1	F
ECGFiveDays	23	861	136	2	ECG	1	T
ElectricDevices	8926	7711	96	7	DEVICE	1	F
EOGHorizontalSignal	362	362	1250	12	EOG	1	T
EOGVerticalSignal	362	362	1250	12	EOG	1	T
Epilepsy	137	138	207	4	HAR	3	T
ERing	30	270	65	6	HAR	4	T
EthanolConcentration	261	263	1751	4	OTHER	3	T
EthanolLevel	504	500	1751	4	SPECTRO	1	T
FaceDetection	5890	3524	62	2	EEG	144	T
FingerMovements	316	100	50	2	EEG	28	T
FordA	3601	1320	500	2	SENSOR	1	T
FordB	3636	810	500	2	SENSOR	1	T
FreezerRegularTrain	150	2850	301	2	SENSOR	1	T
FreezerSmallTrain	28	2850	301	2	SENSOR	1	T
Fungi	18	186	201	18	OTHER	1	T

Dataset	Train Size	Test Size	Length	#Classes	Type	#Dim	Balanced
Ham	109	105	431	2	SPECTRO	1	T
HandMovementDirection	160	74	400	4	EEG	10	T
Handwriting	150	850	152	26	HAR	3	F
Heartbeat	204	205	405	2	AUDIO	61	F
HouseTwenty	34	101	3000	2	DEVICE	1	T
InsectEPGRegularTrain	62	249	601	3	EPG	1	F
InsectEPGSmallTrain	17	249	601	3	EPG	1	F
InsectWingbeatSound	25000	25000	600	10	AUDIO	1	T
ItalyPowerDemand	67	1029	24	2	SENSOR	1	T
LargeKitchenAppliances	375	375	720	3	DEVICE	1	T
Libras	180	180	45	15	HAR	2	T
Lightning2	60	61	637	2	SENSOR	1	F
Lightning7	70	73	319	7	SENSOR	1	F
LSST	2459	2466	36	14	OTHER	6	F
Meat	60	60	448	3	SPECTRO	1	T
MoteStrain	20	1252	84	2	SENSOR	1	T
MotorImagery	278	100	3000	2	EEG	64	T
NATOPS	180	180	51	6	HAR	24	T
NonInvasiveFetalECGThorax1	1800	1965	750	42	ECG	1	T
NonInvasiveFetalECGThorax2	1800	1965	750	42	ECG	1	T
OliveOil	30	30	570	4	SPECTRO	1	F
PEMS-SF	267	173	144	7	MISC	963	T
Phoneme	214	1896	1024	39	SOUND	1	F
PigAirwayPressure	104	208	2000	52	HEMO	1	T
PigArtPressure	104	208	2000	52	HEMO	1	T
PigCVP	104	208	2000	52	HEMO	1	T
Plane	105	105	144	7	SENSOR	1	F
PowerCons	180	180	144	2	DEVICE	1	T
RacketSports	151	152	30	4	HAR	6	T

Dataset	Train Size	Test Size	Length	#Classes	Type	#Dim	Balanced
RefrigerationDevices	375	375	720	3	DEVICE	1	T
Rock	20	50	2844	4	SPECTRO	1	T
ScreenType	375	375	720	3	DEVICE	1	T
SelfRegulationSCP1	268	293	896	2	EEG	6	T
SelfRegulationSCP2	200	180	1152	2	EEG	7	T
SemgHandGenderCh2	300	600	1500	2	SPECTRO	1	T
SemgHandMovementCh2	450	450	1500	6	SPECTRO	1	T
SemgHandSubjectCh2	450	450	1500	5	SPECTRO	1	T
SmallKitchenAppliances	375	375	720	3	DEVICE	1	T
SonyAIBORobotSurface1	20	601	70	2	SENSOR	1	F
SonyAIBORobotSurface2	27	953	65	2	SENSOR	1	T
StandWalkJump	12	15	2500	3	ECG	4	T
StarlightCurves	1000	8236	1024	3	SENSOR	1	F
Strawberry	613	370	235	2	SPECTRO	1	T
Trace	100	100	275	4	SENSOR	1	T
TwoLeadECG	23	1139	82	2	ECG	1	T
UWaveGestureLibrary	2238	2241	315	8	HAR	3	T
Wafer	1000	6164	152	2	SENSOR	1	F
Wine	57	54	234	2	SPECTRO	1	T

Table 5.1: Summary of the 77 data sets used in the experiments

## 5.2 Experimental Design

In this section we discuss the design of the experiments for the 2 components of our framework; the testbed and the recommender. The experiments for the testbed include running the implemented TSCAs on the data archives using different data chunks and calculating their performance using our implementation for the  $F_\beta$  – measure previously discussed in subsection 4.3.3. While the experiments for the recommender include learning classifiers on the results of the testbed to predict the best performing classifiers for unseen data sets.

### 5.2.1 The Testbed Experimental Design

Our experiments were conducted on univariate data sets from the UCR archive as well as multivariate data sets from the UEA archive. The data sets were chosen based on the criteria discussed in the subsection 5.1.3. Each of the data archives offers a train and test split of the data which we have used unchanged for all the classifiers. All our experiments were run on a LINUX server with an AMD Ryzen 7 1700 Processor and 32GB RAM using Python version 3.7.10.

We used the same testbed configuration through out all our experiments, which we will discuss as follows. We fixed the value of the parameter  $Splits = 10$ , that is for every data set the data would be revealed for the classifiers incrementally in batches of 10% from the total length. The chopping algorithm was set to always reveal data from the beginning of the time series.

Due to time limitations, we imposed a constraint on the number of chunks to be used and a time constraint for each experiment. We restricted our experiments to run only on the 1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup> and 10<sup>th</sup> chunks of the data sets by setting  $ChunksToKeep = \{1,2,3,10\}$ . The 10<sup>th</sup> chunk was included to represent the baseline performance of each classifier if shown the full length data. While the first 3 chunks were used to represent the classifiers in the early context scenarios. We defined a time constraint of 2 days for each experiment. Experiments that didn't finish within the time constraint were stopped to make space for other experiments to run.

Hyperparameters optimization was always carried out on all data sets with a  $NumIterations$  maximum sampling value of 50 iterations, unless proven unfeasible for any of the classifiers due to memory shortage or time constraint. In this case the experiment is repeated for all classifiers without hyperparameters optimization. We fixed the number of cross validation folds to 5. The selected performance metric for our experiments was *BalancedAccuracy*. Balanced accuracy is a performance metric which can handle data sets with skewed class distributions by avoiding inflated performance estimates. To calculate balanced accuracy, the *Recall* value is computed for each class then averaged over the total number of classes.

### Classifiers that failed the experiments

Some of the included classifiers in the framework were excluded from our experiments; either because they couldn't operate in the early classification context or because they couldn't attain comparable results to the published performances by previously published frameworks.

For example, KNNED and InceptionTime were excluded due to the nature of their techniques which couldn't handle our created context. Both algorithms are clearly able to learn on chopped training data sets. However once the testing phase is reached,

they would fail to classify instances of the testing data, showing errors related to mismatches between the expected length of instances and the provided length. Yet they would finish the last chunk where the data is provided in its original full length. The reason why KNNED fails such scenario, is that it uses ED which is a point-wise comparison distance measure that cannot compare time series of unequal lengths [Tan+19]. On the other hand, InceptionTime is a deep learning model whose input layer architecture, represented by number of nodes, depends on the length of the input time series [Faw+19a]. Since we use full length instances for testing, this caused an overflow of input data than what the model structure was expecting. There have been literature discussing adapting TSCAs to unequal time series, but this is out of the scope of our experiments. For more details refer to [CCP09; Tan+19; Faw+19a]

Although KNNDTW has been a competent time series classifier for decades; thanks to the exploitation of the elastic distance measure DTW. We couldn't get either implementation of KNNDTW from *sktime* and *pyts* to work on the chopped data; due to errors in the data representation needed by lower level internal libraries. KNNDTW would have been able to operate in the early classification context if used with a full warping window; which would have been successful in handling extreme classification cases like classifying full length data even when learning on the 10% chunk data set.

Two classifiers were excluded because they attained significantly inferior results compared to the published scores; these are KNNMSM and LS. Specially on the InsectWingbeatSound data set, for which we attained a difference in performance of -45.71% for KNNMSM and -20.71% for LS than the results published by [Bag+17]

### Adjustments applied to remaining classifiers

Our experiments proceeded with the remaining 5 classifiers; PForest, TSF, ST, CBOSS and WEASEL. To focus only on the comparison of performance between classifiers and not preprocessing, we consider only data sets where all instances have the same length and no attributes have missing data. We also focus our interest on data which is originally a timeseries data. This means that it is a collection of specific measurements across a span of time. Which lead us to exclude some data set types from taking part in the experiment, more details are presented in subsection 5.1.3. We cover a total of 77 data sets from both the UR and the UEA archives, out of which 55 data sets are univariate and 22 data sets are multivariate.

Time and resources limitations played a big role shaping our experiments. Running classifiers is computationally expensive [SL20] and takes up to hundreds of processing days. This lead us to modify some algorithms so that they can finish promptly; like in the case of PForest, or to use an enhanced version of the original classifier with an option to control the learning phase time; like BOSS and ST.

Due to time limitations; We used the contractable implementations of BOSS and ST offered by *sktime*. These implementations allows passing a time parameter which controls the sampling space for BOSS and the shapelets searching time for ST, beside other performance enhancements. We set the value of the time contract for both to 60 minutes. We have discussed the literature details of these enhancements for ST in subsection 2.3.3 and for BOSS in subsection 2.3.4.

We applied 2 adjustments to the PForest classifier. First, we excluded ED from the pool of distance measures that PForest selects from; since the ED has previously proven not to be compatible with our early classification context when used with the KNN classifier. Second, we excluded TWE from the pool of distance measures when the size of the training data set exceeded 150 instances, which is the median of all

data sets training sizes. TWE is a very slow distance measure. It has been reported by [Bag+17] to be the slowest among all elastic distance measures. When measured on classifying 10 test instances from the StarlightCurves data set, it performed 10,000 times slower than ED.

Also we relaxed the time constraint for PForest; due to its quadratic complexity with length [Luc+19] which made it very slow as we used longer chunks. Since the other algorithms already had a bit of an advantage either for using simple tranformations or for using time contracts, we felt that it would be unfair for PForest to abide by the time constraint for the 10<sup>th</sup> chunk. Our exception for PForest was if the first 3 chunks were able to finish within the time constraint, we would leave the 10<sup>th</sup> chunk to run even if it exceeded it. More details about runtime durations is provided in our results.

### Learning weights for the $F_\beta$ – measure

We use the  $F_\beta$  – measure as an objective score for the performance of the classifiers. The  $F_\beta$  – measure has been a popular choice for evlauating eTSC algorithms [SL20] because of its capacity to combine both earliness and accuracy. Previous literature has always considered using a special case of the scoring function in which both earliness and accuracy contribute with the same weights. This case is referred to as the Harmonic mean or the F-1 score and can be simply achieved by assigning the value of  $\beta$  to 1.

For our experiments we conducted some experiments to decide if the use of equal weights for accuracy and earliness best suits our newly created context. We experimented with the harmonic mean (i.e.  $\beta = 1$ ) like the previous papers and plotted the distribution of values for  $F_\beta$  – measure. Figures 5.2 show a comparison between the distribution of values of  $F_1$  and the distribution values of *BalancedAccuracy* for CBoss. The boxplots from the accuracy image shows that CBoss on the full length of the data has a median value of **0.829**, while on the 10% data it shows a median value of **0.44**. The two boxes do not overlap which means that there is a significant difference in the scores of the two versions. What we would like for the  $F_\beta$  – measure in this scenario, is to give a privilege for the 10% version of the classifier as it learns on less training data, such that it would have a high score if it can score a balanced accuracy close to that of the 100% version.

By checking the boxplots for the  $F_1$ , we find that the median value for the 10% version is **0.83**, while the median value of the 100% version is **0.708**. Also the distribution of the values for the 10% version is very close to the median and the minimum value is higher than the median of the 100% version. Which means that even for the data sets on which the 10% version of CBoss scores very low accuracy scores, it will still get a high  $F_1$  relative to the 100% version. The harmonic mean over-compensates for the lower accuracy scores with very high values for earliness causing misleading values.

To decide which  $\beta$  value is best to use, we carried out a sequence of experiments using  $\beta = [0.1, 0.9]$ . The value that gave the most reasonable results was 0.5. Figure 5.3 shows the same comparison between the boxplots of *BalancedAccuracy* and  $F_\beta$  – measure for CBoss but with  $\beta = 0.5$ . Still  $F_1$  allows for compensation between accuracy and earliness for the 10% version, but it refines the values. The median value for the 10% version became **0.705** while that of the 100% is **0.795**. The distribution of values for the 10% is spread more and the maximum value is falling near the median value of the 100%.

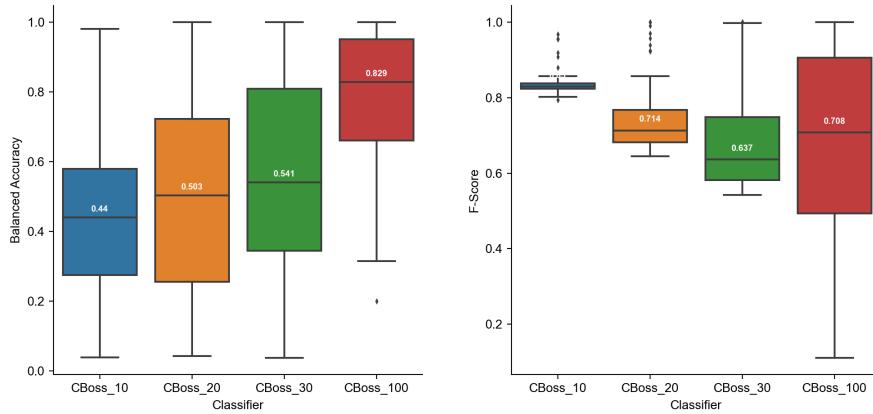


Figure 5.2: Comparison between Balanced Accuracy (left) and  $F_1$  (right) for CBoss

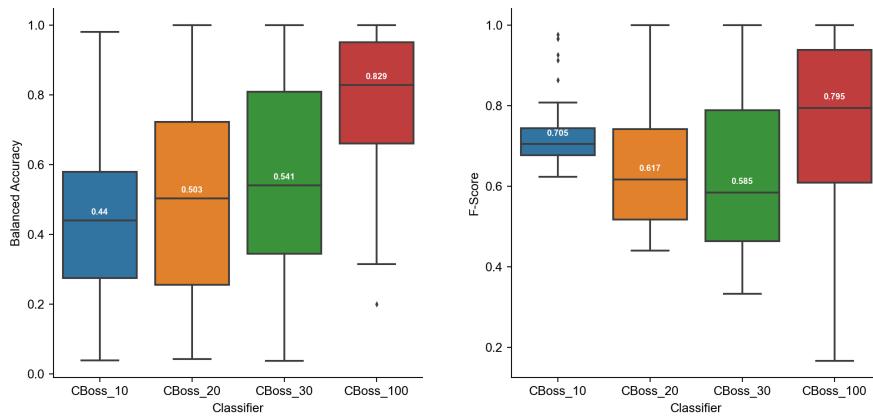


Figure 5.3: Comparison between Balanced Accuracy (left) and  $F_{0.5}$  (right) for CBoss

We believe that selecting the best value of  $\beta$  is not trivial. It is problem dependant and not one specific value can fit for all problems. For our experiments we decided to give earliness less priority than performance, but in other situations it might be the case that both of them contribute with the same importance or even earliness is more critical.

### 5.2.2 The Recommender Experimental Design

As previously discussed in the recommender design, we train one random forest regressor per chunk, which we call the chunk learner. Each chunk learner is trained on the results for a specific chunk and is responsible for predicting the performance score,  $F_\beta - measure$ , for all classifiers on unseen data sets for that specific chunk.

To train the chunk learners, we use the results of experiments from the testbed combined with metadata about the data set. Each training instance consists of metadata features from the table 5.1 in addition to 3 features from the testbed results in table 4.6; the classifier name, the revealed% and the  $F_\beta - measure$  value. The results of all data sets are included in the learning process, regardless of whether they were finished by all the classifiers and for all the chunks or not. We only exclude the three data sets PEMS-SF, DuckDuckGeese and FaceDetection; because they were only finished by one classifier. Which leaves us with 74 data sets.

In the beginning, we randomly split the data into training and testing data sets. Then for each chunk learner, we start a grid search process, on the training data, for learning hyperparameters and model selection. In the end, the chunk learners are evaluated using the best selected model on the unseen testing data sets for the final evaluation.

For the creation of training and testing data sets, we use an 80/20 split with random selection. A total of 59 data sets are used for the learning of chunk learners, while 15 data sets are left out for the final evaluation.

For the training of the chunk learners we use a grid search; to select the best combination of hyperparameters for each chunk learner model for its designated chunk. During the grid search operation, a 10-fold cross validation process is carried out in which each training model is trained and evaluated 10 times. The final score for each model is represented by its average across all 10-folds. We define the hyperparameter space for 2 parameters of the random forests; the maximum depth of trees =  $\{3, 5, 10, \infty\}$  and the number of trees per forest =  $\{2, 3, 4, 5, 10, 50, 100, 250, 500\}$ . After the grid search is finished, the best performing model is selected to be used as the final chunk learner on the unseen testing data sets. To avoid biases in the results of the chunk learners from the selection of training and testing data sets, this whole process is repeated 50 times using random selection of the data sets for training and testing. We use the seeds = [0, 49] for the random selection of data sets and for the bootstrapping and sampling of features in the building for trees; to allow for reproducability of the results.

Since we are using regression random forests, one mandatory preprocessing step has to be carried out; to transform the categorical features into numerical features. We do *one-hot-encoding* on the two features *classifier* and *type*, which represent the name of the classifier that was used and the type of the data set respectively.

# Chapter 6

# Results

This chapter analyzes the results of our experiments.

## 6.1 Testbed Results

We will divide our analysis as follows; analysis for the results of the same classifier across different revealed percentages of data (chunks), analysis for the results of different classifiers across the same revealed percentages of data (chunks) and analysis for different classifiers and different percentages all together using the  $F_\beta$  – measure.

There were two classifiers that we excluded because they attained significantly inferior results compared to their published scores; these are KNNMSM and LS. Specially on the InsectWingbeatSound data set, for which we attained a difference in performance of -45.71% for KNNMSM and -20.71% for LS than the results published by [Bag+17].

For the remaining classifiers, not all of them were able to finish on all data sets within the time constraint. Specially the multivariate archive, which was challenging for most of the algorithms in terms of memory and time complexity. Our choice of adapting univariate classifiers for multivariate problems using column ensembling might have made things even more challenging; because this technique fits for each dimension in the data sets a classifier, which is a memory and space extensive solution. This was clearly evident for the runs as the number of dimensions of the data sets increased.

TSF was the most successful algorithm in finishing the biggest number of datasets. It was able to finish all the data sets from both archives, with the exception of FaceDetection only. WEASEL was the second best, finishing all data sets from both archives, except for the largest three multivariate data sets; DuckDuckGeese, FaceDetection and PEM-SF. WEASEL is known for its quality and speed but not for scalability [Luc+19], it failed these data sets due to out of memory errors. CBoss failed to run on DuckDuckGeese, FaceDetection, PEM-SF and MotorImagery within our time constraint. ST is known to be one of the slowest algorithms, but we use an enhanced version to accelerate its performance. However, ST was not able to finish within our time limit on 15 multivariate datasets; due to its linear time complexity with the number of dimensions. It hit the limit of the time constraint once the number dimensions reached 6. It has also failed on the Fungi dataset, because it couldn't find high quality

shapelets during the contract time that we configured it to use. The slowest of all and which achieved the least progress on the data sets was PForest. It finished a total of 7 multivariate data Sets and 35 univariate data sets. Although PForest can attain quasi linear training times with the size of training data, but it has a training complexity, which can be quadratic with length based on its randomly selected distance measure. This is why PForest hit the time constraint easily as the length of the time series increased. All 5 classifiers completed 36 common data sets. TSF, WEASEL and CBoss finished more than 72 data sets, ST finished 61 data sets and PForest finished 42 data sets.

Tables 6.1 and 6.2 present the results of our experiments for the used 5 classifiers on all 77 data sets. We display values for balanced accuracy,  $F_\beta - measure$  value and revealed% for each classifier. The default train/test split of the data sets was used for all data runs.

Table 6.1:  $F_\beta$  scores for CBoss, PForest and ST on 77 data sets for all chunks

Data set	CBoss				PForest				ST			
	10%	20%	30%	100%	10%	20%	30%	100%	10%	20%	30%	100%
ACSF1	77.78	74.77	70.81	87.80	65.27	48.19	37.80		80.30	81.63	84.33	78.47
AtrialFibrillation	66.67	48.00	43.10	16.67	70.00	60.00	40.00	28.57	71.43	48.00	53.85	34.78
BasicMotions	71.05	84.21	100.00	100.00	67.39	61.54	55.88	70.59				
Beef	67.65	50.00	58.00	76.19	67.14	51.06	37.10	25.53	66.67	55.81	53.85	76.19
Car	68.46	53.93	51.89	80.00	71.43	53.33	44.74	42.80	68.75	59.26	61.34	78.09
Chinatown	63.88	65.26	71.10	93.88	65.48	54.84	97.48	96.74	66.30	73.11	72.93	96.38
CinCECGTorso	68.54	54.33	45.28	75.78					67.52	55.62	56.06	89.88
Coffee	73.14	95.73	95.66	100.00	76.42	65.12	71.15	100.00	79.17	84.85	95.66	91.23
Computers	77.17	71.68	68.11	64.25	74.27	62.89	57.14	55.37	75.25	67.80	66.28	63.82
Cricket	65.69	74.23	81.71	100.00								
Earthquakes	80.68	76.06	73.98	70.39	80.68	76.06	73.98	71.19	75.36	76.06	73.98	59.50
ECG200	69.42	56.74	47.78	79.62	76.67	77.67	54.83	87.80	76.32	80.81	75.75	83.09
ECG5000	69.12	67.39	81.34	92.83	74.69	65.85	61.37	55.33	81.47	87.15	89.73	92.48
ECGFiveDays	74.41	61.68	56.06	100.00	72.86	65.43	60.06	82.08	72.29	67.34	65.21	99.85
ElectricDevices	70.24	57.31	55.78	65.58					69.61	55.75	43.04	32.92
EOGHorizontalSignal	65.06	46.82	40.58	42.00	65.09	46.14	33.93		65.22	46.59	40.58	41.73
EOGVerticalSignal	65.38	47.04	37.77	36.93					65.06	46.74	40.95	33.53
Epilepsy	96.62	100.00	100.00	100.00					79.93	79.20	86.78	98.19
ERing	68.00	53.40	51.35	87.36	66.17	50.15	40.61	45.18				
EthanolConcentration	67.41	51.27	42.35	35.79					69.54	55.90	47.29	55.41
EthanolLevel	67.79	51.81	43.46	39.37					68.66	56.18	47.19	30.85
FingerMovements	71.05	62.02	55.26	51.47								
FordA	80.03	83.41	82.37	68.33					70.87	61.50	59.44	49.67
FordB	77.65	74.65	72.32	56.72					72.85	60.17	60.37	48.13
FreezerRegularTrain	72.22	97.15	99.87	99.56	78.90	62.50	89.63	90.49	72.37	94.33	99.87	99.87
FreezerSmallTrain	72.26	94.57	99.74	96.56	75.18	63.60	72.00	66.11	71.91	91.27	99.82	99.96

Table 6.1:  $F_\beta$  scores for CBoss, PForest and ST on 77 data sets for all chunks (continued)

Data set	CBoss				PForest				ST			
	10%	20%	30%	100%	10%	20%	30%	100%	10%	20%	30%	100%
Fungi	65.16	45.53	33.24		64.92	47.57	43.73	89.47				
Ham	71.88	61.31	55.00	61.54	72.11	62.69	55.59	65.63	71.88	60.87	58.00	71.91
HandMovementDirection	68.14	50.08	43.28	31.49	66.70	51.84	47.37	32.75				
Handwriting	64.69	44.49	37.53	41.10					64.69	44.65	35.96	31.14
HouseTwenty	71.72	78.55	84.89	96.87	81.90	75.44	71.84	87.70	76.27	77.91	82.39	80.81
InsectEPGRegularTrain	96.72	94.32	100.00	100.00	100.00	99.50	100.00	100.00	69.18	55.92	48.55	98.50
InsectEPGSmallTrain	91.29	100.00	94.67	100.00	99.50	90.05	98.51	100.00	69.18	55.77	48.33	96.51
InsectWingbeatSound	65.22	48.98	45.17	49.64					65.21	48.51	39.14	56.54
ItalyPowerDemand	69.43	57.66	78.91	86.41	70.71	71.06	58.02	94.95	68.66	58.52	69.13	94.47
LargeKitchenAppliances	72.73	61.86	57.14	72.89	69.79	56.07	50.45	37.83	69.04	54.95	47.26	61.54
Libras	62.41	44.04	41.36	79.36	62.69	42.74	32.38	5.41	62.99	45.24	39.66	83.87
Lightning2	73.28	64.38	58.50	69.23	70.89	66.12	58.50	80.31	73.28	59.66	55.37	65.65
Lightning7	69.41	55.41	52.18	69.45	66.19	49.74	39.63	25.66	65.84	47.71	42.80	64.97
LSST	62.41	51.74	45.33	45.29								
Meat	86.36	70.59	85.00	85.83	73.53	57.14	46.43	59.77	80.77	92.31	80.12	95.87
MoteStrain	71.39	74.58	63.24	93.20	70.60	78.17	76.93	86.88	73.41	59.66	71.72	89.92
NATOPS	64.69	51.86	46.40	80.00	63.76	47.53	37.95	18.60				
NonInvasiveFetalECGThorax1	71.16	56.47	63.73	78.51					72.98	70.56	62.04	28.48
NonInvasiveFetalECGThorax2	68.38	70.68	57.54	82.27					67.95	74.08	75.43	51.36
OliveOil	69.35	85.71	92.11	83.87	66.67	50.00	50.00	51.13	77.27	92.31	88.46	87.80
Phoneme	66.27	50.61	42.38	26.83	64.43	44.79	32.33	1.35	66.47	50.18	41.30	23.45
PigAirwayPressure	66.05	50.82	41.85	91.15					65.81	50.06	51.30	75.43
PigArtPressure	97.71	99.40	98.22	97.61					88.81	90.73	95.34	95.24
PigCVP	80.77	87.86	92.57	94.06	64.43	45.29			68.81	58.92	72.51	72.19
Plane	67.21	48.84	52.23	100.00	69.09	53.16	56.37	100.00	67.36	63.16	63.57	100.00
PowerCons	74.32	71.64	75.95	90.46	76.47	67.92	68.47	99.31	73.68	61.80	66.73	92.48

Table 6.1:  $F_\beta$  scores for CBoss, PForest and ST on 77 data sets for all chunks (continued)

Data set	CBoss				PForest				ST			
	10%	20%	30%	100%	10%	20%	30%	100%	10%	20%	30%	100%
RacketSports	67.49	50.96	43.89	77.24	69.14	53.19	46.54	38.04				
RefrigerationDevices	72.80	60.00	56.47	48.83	70.00	58.48	55.30	40.66	72.46	60.36	57.14	50.45
Rock	70.83	58.82	61.11	87.80	68.18	53.33	42.78	40.53	67.54	51.28	44.74	85.44
ScreenType	70.16	58.59	52.12	43.79	69.33	54.05			69.08	56.82	50.45	43.00
SelfRegulationSCP1	76.14	70.52	68.62	74.49								
SelfRegulationSCP2	72.09	64.00	62.11	47.20	70.81							
SemgHandGenderCh2	78.10	69.57	76.02	85.63					69.11	70.90	82.35	77.52
SemgHandMovementCh2	66.33	48.98	39.21	59.18					66.89	55.73	47.47	41.83
SemgHandSubjectCh2	67.89	51.65	49.14	70.96					68.37	62.18	61.72	52.72
SmallKitchenAppliances	79.64	73.53	68.42	73.78	69.38	58.37	48.54		71.13	55.35	53.21	74.08
SonyAIBORobotSurface1	72.28	80.16	52.72	41.26	81.68	83.07	84.85	79.08	70.60	71.04	70.17	92.02
SonyAIBORobotSurface2	70.00	66.69	85.02	89.47	75.47	66.44	73.24	85.48	75.54	67.12	64.01	89.47
StandWalkJump	70.00	50.00	46.43	28.57	68.75	52.17	46.43		68.75	57.14	46.43	28.57
StarlightCurves	68.05	52.71	45.23	96.88					74.31	69.86	79.00	91.51
Strawberry	76.43	69.16	95.73	96.98	76.43	69.16	65.63	90.39	80.33	77.69	76.86	92.68
Trace	68.35	57.55	70.00	100.00	72.73	61.54	54.65	47.43	68.18	50.63	39.55	100.00
TwoLeadECG	72.36	60.84	65.52	99.56	75.81	70.10	67.14	93.07	72.23	61.76	63.57	99.78
UWaveGestureLibrary	65.46	48.39	39.30	81.57					65.69	47.76	36.88	85.95
Wafer	92.71	98.07	98.47	99.61	93.19	95.00	89.55	99.35	94.50	100.00	99.49	100.00
Wine	75.96	63.34	71.80	69.57	71.77	64.29	70.30	53.78	74.77	65.26	65.99	71.62

Table 6.2:  $F_\beta$  scores for TSF and WEASEL on 77 data sets for all chunks

Data set	TSF				WEASEL			
	10%	20%	30%	100%	10%	20%	30%	100%
ACSF1	77.40	68.38	66.13	59.77	78.99	76.19	68.42	89.00
AtrialFibrillation	66.67	54.55	34.38	22.54	71.43	54.55	43.10	22.54
BasicMotions	82.00	96.97	100.00	100.00	75.00	100.00	100.00	100.00
Beef	78.57	72.73	72.73	65.12	65.79	51.06	50.00	76.19
Car	75.00	71.64	66.13	66.93	68.46	53.93	45.58	83.87
Chinatown	95.61	98.20	97.84	95.66	63.88	47.42	39.12	93.53
CinCECGTorso	72.35	67.48	73.58	88.76	76.16	81.60	68.91	95.16
Coffee	89.44	84.85	80.57	100.00	73.14	91.80	65.60	100.00
Computers	78.41	72.99	70.97	66.42	71.93	70.42	70.32	67.29
Cricket	84.45	85.21	91.79	89.80	71.56	65.75	77.81	98.27
Earthquakes	80.68	76.06	73.98	70.39	80.68	75.54	73.37	70.39
ECG200	84.48	85.11	86.32	81.93	74.39	59.26	51.19	81.93
ECG5000	92.19	91.23	90.26	92.18	74.10	80.70	74.52	92.64
ECGFiveDays	75.16	67.28	62.13	90.32	73.20	63.15	58.82	99.85
ElectricDevices	73.15	63.54	61.22	63.55	75.13	66.37	60.86	72.62
EOGHorizontalSignal	65.44	46.59	39.58	45.81	64.96	47.43	36.72	42.27
EOGVerticalSignal	65.25	46.82	40.45	39.85	64.90	46.74	35.13	27.14
Epilepsy	91.69	97.35	97.33	96.40	98.25	100.00	100.00	100.00
ERing	75.47	76.89	82.66	94.05	68.56	58.20	58.76	94.05
EthanolConcentration	70.45	58.38	49.96	36.88	67.70	51.77	46.49	50.73
EthanolLevel	69.42	54.64	45.54	39.37	67.42	51.68	42.40	43.85
FingerMovements	70.62	62.50	55.88	44.44	72.47	62.02	55.88	53.51
FordA	79.40	75.70	75.85	77.74	83.29	79.70	91.60	95.77
FordB	74.88	67.22	64.00	63.52	76.30	76.96	64.17	78.30
FreezerRegularTrain	96.45	99.82	99.82	98.91	72.22	92.87	73.29	97.51
FreezerSmallTrain	88.93	97.98	99.78	80.85	72.22	89.55	64.56	95.05

Table 6.2:  $F_\beta$  scores for TSF and WEASEL on 77 data sets for all chunks (continued)

Data set	TSF				WEASEL			
	10%	20%	30%	100%	10%	20%	30%	100%
Fungi	65.54	52.80	54.89	93.37	64.86	45.95	34.29	94.68
Ham	74.71	76.36	72.73	74.04	72.58	62.22	55.00	59.52
HandMovementDirection	67.70	54.21	43.94	49.85	67.29	50.08	38.92	31.49
Handwriting	65.70	49.52	43.41	26.49	64.71	44.16	33.38	26.39
HouseTwenty	86.28	92.25	89.27	77.91	88.64	94.07	92.04	95.83
InsectEPGRegularTrain	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
InsectEPGSmallTrain	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
InsectWingbeatSound	71.15	64.89	63.83	56.33	65.66	53.33	48.13	58.39
ItalyPowerDemand	75.92	71.42	76.09	95.31	69.43	58.41	55.66	91.17
LargeKitchenAppliances	72.12	60.61	54.17	51.54	68.16	54.05	46.29	63.82
Libras	67.90	64.65	68.53	76.82	63.52	48.98	50.78	83.87
Lightning2	73.28	68.93	70.39	65.65	73.28	66.12	64.14	62.12
Lightning7	67.55	61.86	58.96	69.45	68.67	52.42	41.50	67.95
LSST	68.85	57.77	52.68	44.89	68.23	58.44	51.08	53.18
Meat	83.33	82.76	83.33	91.80	76.09	61.54	46.43	81.93
MoteStrain	87.89	86.66	86.98	84.52	77.16	78.17	79.46	94.18
NATOPS	64.43	52.07	51.71	76.82	64.26	52.69	54.07	80.00
NonInvasiveFetalECGThorax1	83.62	79.59	78.29	84.53	64.54	49.42	40.60	
NonInvasiveFetalECGThorax2	86.97	83.53	80.52	87.84	64.53	46.37	36.14	90.45
OliveOil	80.00	80.00	88.46	83.87	65.79	48.98	44.74	87.80
Phoneme	65.76	49.36	39.78	16.52	66.84	51.35	43.60	25.78
PigAirwayPressure	65.34	47.22	35.17	18.51	68.07	57.30	49.68	58.65
PigArtPressure	66.24	50.67	43.69	35.15	93.70	96.52	95.91	100.00
PigCVP	65.99	49.03	38.80	23.63	71.42	68.09	81.55	91.73
Plane	91.18	97.67	95.38	98.81	66.15	48.84	86.92	100.00
PowerCons	81.03	77.42	77.96	99.31	73.23	64.57	62.59	95.19

Table 6.2:  $F_\beta$  scores for TSF and WEASEL on 77 data sets for all chunks (continued)

Data set	TSF				WEASEL			
	10%	20%	30%	100%	10%	20%	30%	100%
RacketSports	75.17	72.55	78.21	87.96	69.14	55.88	52.32	84.85
RefrigerationDevices	72.94	63.29	58.00	53.17	69.69	57.92	52.90	44.84
Rock	71.74	64.52	63.92	69.48	67.24	57.97	39.11	76.19
ScreenType	69.89	57.80	52.58	38.08	71.01	57.58	52.27	47.76
SelfRegulationSCP1	80.97	83.89	82.65	76.04	76.73	86.37	78.45	76.04
SelfRegulationSCP2	72.93	63.16	58.36	42.26	72.09	63.16	61.72	55.12
SemgHandGenderCh2	81.50	89.05	87.58	91.20	76.67	70.07	66.25	74.87
SemgHandMovementCh2	73.56	67.92	71.26	81.28	66.33	48.98	39.60	25.13
SemgHandSubjectCh2	75.86	78.26	76.96	89.40	66.67	50.00	45.29	58.72
SmallKitchenAppliances	77.27	70.59	66.38	79.54	74.35	70.26	58.70	67.29
SonyAIBORobotSurface1	78.73	77.20	74.60	64.45	72.12	88.41	86.21	80.03
SonyAIBORobotSurface2	80.88	79.26	77.85	74.13	79.73	77.04	89.81	93.40
StandWalkJump	71.43	54.55	53.85	34.78	71.43	57.14	58.00	28.57
StarlightCurves	92.50	94.13	95.21	95.28	69.11	52.94	43.76	93.75
Strawberry	89.57	90.24	94.78	93.66	76.43	69.16	65.63	95.65
Trace	77.40	73.39	68.42	98.75	67.24	54.42	50.00	100.00
TwoLeadECG	80.11	82.76	93.48	70.18	72.23	62.30	64.93	99.78
UWaveGestureLibrary	70.06	67.72	71.01	81.93	65.24	49.71	42.90	77.61
Wafer	99.48	99.06	99.31	99.31	99.32	88.10	87.63	99.96
Wine	77.27	66.26	61.98	67.53	72.22	61.54	60.71	75.77

We divide our analysis into three parts: a comparison of the effect of data chopping within the same classifier, comparison on the effect of data chopping across different classifiers and unified comparison between all classifiers across all data chops. Since some of the analysis results we present include comparisons between the same classifier trained on different percentages of data, we use a special naming convention. We name the classifiers using the format *NameOfClassifier\_Revealed%*, where the first part of the name refer to the name of the classifier and the second part refers to the percentage of data used to train the classifier. For example, *TSF\_10* refers to classifier TSF which was trained on 10% of the data length. This naming convention helps clarify the type of the classifier and the context of comparison regardless of the type of analysis being discussed.

### 6.1.1 Results Across Classifiers

An aspect to investigate within the early classification context is the relative performance of classifiers to each other while more data becomes available to train on. We fix the percentage of data revealed, *revealed%*, and compare the performance of classifiers, which is represented by their ranking based on the value of  $F_\beta$  – *measure*. We will dig deeper into the properties of the data sets, to explore the relation between the characteristics of the data and the best performing algorithm.

To better understand the actual performance of the classifiers on the data sets and assess the quality of their results, we introduce a new baseline classifier to the comparison. The baseline classifier always predicts the majority class of a given data set. We name the classifier *dummy*; as an indication of its role in the comparisons. To ensure the integrity of our comparisons, we calculate balanced accuracy and  $F_\beta$  – *measure* for the dummy classifier on all data sets and for the same 4 chunks. Although the dummy classifier will always predict the majority class and thus score the same balanced accuracy, but its  $F_\beta$  – *measure* will differ based on the *revealed%*.

Figure 6.1 shows the critical difference diagram across all the classifiers and for a fixed revealed percentage of data. In the beginning when only 10% of the data is revealed, TSF is significantly better than all the other classifiers. The remaining four classifiers form a clique together and are not significantly different from the dummy classifier. Inside that clique, CBoss is the best classifier, WEASEL is the second best and both PForest and ST are almost identical in their ranking.

When the revealed percentage is increased to 20%, TSF still prevails as the first classifier but with no significant different to CBoss and WEASEL. There is still no significant difference between the PForest, ST and WEASEL and between the dummy classifier. ST is able to learn better shapelets on the data and beats PForest. On the third chunk, CBoss and ST are not significantly different from TSF and join it to form the first clique. While WEASEL and PForest are not significantly different from the dummy classifier. The middle clique shows no significant difference between CBoss, ST, WEASEL and PForest. However WEASEL and PForest are both significantly worse than TSF. On the full length of data, the 5 classifiers are all not significantly different from each other. WEASEL attains the best ranking followed by ST and CBoss which are identical, while TSF loses its edge on full length data. PForest ranks the last among the 5 classifiers. All the classifiers form one clique which is significantly better than the dummy classifier.

The most apparent conclusion from the graphs is that TSF is able to do better on data sets during the earlier stages of the early classification context. We believe the reason why TSF is always better on the shorter chunks of data; is because it uses

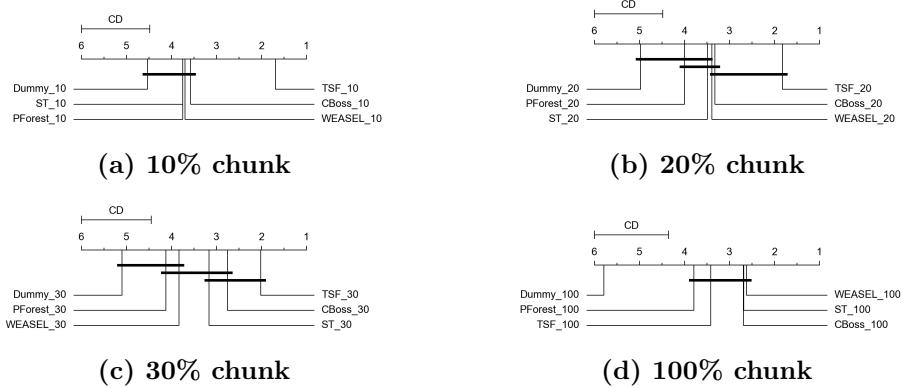


Figure 6.1: Critical Difference diagrams per chunk across all classifiers

three simple aggregations as features and doesn't depend on advanced features like the other classifiers do. On the other side, the BOP technique applied by WEASEL and CBoss is challenged by the chopping, which notably decreases the chances of repetition of distinctive patterns in the data within such a small length. PForest is also expected to perform badly as it needs to stretch the extent of its elastic distance measures; by measuring distances between testing instances which are 10 times the length of the training instances. ST gets the lowest scores because in cases like our experiments where majority of the data is chopped, it is very hard to find shapelets that can uniquely identify the different classes by learning on only the first 10% of the data. An obvious change is visible when more data is made available for training, the techniques of; WEASEL; CBoss and ST are able to make use of the more available information to find better and more distinctive features in the data; which helps them be more competent. Also PForest is able to score get better performances; because its elastic distance measure do not have to compensate for the big difference in length.

### 6.1.2 Data Set characteristics and performance in early classification context

The goal of our research is to learn a recommender which can suggest the best classifiers to unseen data sets. The recommender analyzes the characteristics of the data set and then based on the performance from previous experiments, it suggests the most suitable classifiers. Thus it could be very convenient if the characteristics of the data sets we used; like length of series, number of classes and train size are able to provide insights about the best algorithm for a data set. These comparison aspects were first defined in [Bag+17] and have also been adopted in [Faw+19a]. The group of data sets we use for this analysis is very small and for some analysis we break them down by certain criteria. We believe that our results should be interpreted carefully and not generalized for all data sets.

### Length of Series

The first feature of the data sets we investigate is the length of the time series. Like the findings of [Bag+17] on univariate data sets and the results of [Faw+19a] on multivariate data sets, our results in the early classification context, didn't show information about the performance of the classifiers in accordance with the length of time series. We show the results of the 5 classifiers on data sets grouped by their time series length in table 6.3 for the 10% length. The table display for each classifier in percentage how many times it was able to score higher  $F_\beta - measure$  than the baseline classifier. In the case where multiple classifiers are able to score higher than the baseline on one data set, then the score of this data set is counted for each one of them.

Train Length	#Data sets	CBoss	PForest	ST	TSF	WEASEL
1-50	3	0.00%	0.00%	0.00%	100.00%	0.00%
51-100	5	40.00%	60.00%	40.00%	100.00%	80.00%
101-250	7	71.43%	71.43%	100.00%	100.00%	57.14%
251-500	8	62.50%	75.00%	50.00%	87.50%	37.50%
501-1000	10	60.00%	70.00%	60.00%	70.00%	50.00%
1001+	3	66.67%	33.33%	66.67%	100.00%	66.67%

Table 6.3: Performance of classifiers with respect to length in comparison to baseline on 10% chunk data sets

We couldn't infer a strong relation between the length of time series and the performance of classifiers except for TSF. Our findings for TSF shows that it is consistently better than the baseline on all length groups for all data chunks. It was able to score good results on all the data sets in the groups (51-100), (101-250) and (1001+) even by training on only the 10% chunk. Unlike TSF, none of the other classifiers was able to beat the baseline for the data set group with the least length (1-50). There is a clear enhancement in the performance of all the classifiers as they learn on more data chunks as we have mentioned before, but no clear relation between the length of the time series and the performance across the different chunks. The results for the other chunks are in appendix A.2.

### Training Data Set Size

The second feature of the data sets we investigate is the size of the training data set, and how the number of instances affects the performance of classifiers in our introduced context. We show the results of the 5 classifiers on data sets grouped by their train size in table 6.4 for the 10% chunk, table 6.5 for the 20% chunk, table 6.6 for the 30% chunk and table 6.7 for the 100% chunk. The table display for each classifier in percentage how many times it was able to score higher  $F_\beta - measure$  than the baseline classifier.

TSF consistently scores better than, or at least the same as, the other classifiers on all train size groups for the first 3 chunks. The scores for TSF on all the train size

Train Size	#Data sets	CBoss	PForest	ST	TSF	WEASEL
1-50	14	50.00%	64.29%	50.00%	92.86%	57.14%
51-100	9	66.67%	55.56%	44.44%	77.78%	55.56%
101-250	7	57.14%	57.14%	71.43%	100.00%	42.86%
251-500	4	50.00%	75.00%	75.00%	75.00%	25.00%
501-1000	2	50.00%	50.00%	100.00%	100.00%	50.00%

Table 6.4: Performance of classifiers with respect to train size in comparison to baseline on 10% chunk data sets

Train Size	#Data sets	CBoss	PForest	ST	TSF	WEASEL
1-50	14	71.43%	78.57%	64.29%	92.86%	71.43%
51-100	9	77.78%	88.89%	55.56%	100.00%	66.67%
101-250	7	71.43%	71.43%	71.43%	100.00%	85.71%
251-500	4	75.00%	75.00%	75.00%	75.00%	50.00%
501-1000	2	50.00%	50.00%	100.00%	100.00%	0.00%

Table 6.5: Performance of classifiers with respect to train size in comparison to baseline on 20% chunk data sets

Train Size	#Data sets	CBoss	PForest	ST	TSF	WEASEL
1-50	14	92.86%	71.43%	85.71%	92.86%	71.43%
51-100	9	77.78%	55.56%	55.56%	100.00%	55.56%
101-250	7	85.71%	57.14%	100.00%	100.00%	85.71%
251-500	4	75.00%	50.00%	75.00%	75.00%	50.00%
501-1000	2	100.00%	50.00%	100.00%	100.00%	0.00%

Table 6.6: Performance of classifiers with respect to train size in comparison to baseline on 30% chunk data sets

groups are better than the baseline across the different chunks except for the second group (51-100) on the 10% chunk. TSF levels with the baseline classifier on the two data sets Lightning2 and Lightning7 on the 10% chunk, but beats it once more data is available. During the 20% and the 30% chunks, WEASEL cannot beat the baseline classifier on the data sets in the largest train size group. However on the 100% chunk it catches up and beats it on both data sets. Our results don't show any significance from comparing algorithms based on the training size across the different early classification runs. There is also no significant pattern for the performance of the same classifiers across the runs.

### Type of Problem

The third feature of the data sets we investigate is the type of problem. The intention is to try to find an evidence if there is a classifier which is suitable for certain types of problems. We display the results of the 5 classifiers on data sets grouped by problem type in table 6.8 for the 10% chunk, table 6.9 for the 20% chunk, table 6.10 for the 30% chunk and table 6.11 for the 100% chunk. The table display for each classifier in percentage how many times it was able to score higher  $F_\beta$ -measure than the baseline classifier.

Our results show that TSF is able to achieve better results than the baseline classifier on all data sets of the type SPECTRO. This performance is consistently kept across all the different chops, till the other classifiers catch up on the full length data chop. PForest is not able to beat the baseline on types HAR and SOUND during any of the experiments.

However, since both types consist only of one data set, we cannot conclude that PForest will not be able to learn on such types. During the full length experiment, all the classifiers are able to beat the baseline classifier on all the data sets for the groups; SPECTRO, Traffic, DEVICE and EPG.

### Number of classes

The fourth feature of the data sets we investigate is the number of classes. The goal is to try to find a connection between the number of classes of data sets and if it has an impact on the performance of certain classifiers.

We display our results for the number of classes feature in table 6.12 for the 10% chunk, table 6.12 for the 20% chunk, table 6.12 for the 30% chunk and table 6.12 for the 100% chunk. The table display for each classifier in percentage how many times it was able to score higher  $F_\beta$ -measure than the baseline classifier.

Like our previous observations for the other features, TSF is constantly better than the baseline classifier on all groups across all the chunks. It starts with a good performance on all groups in the 10% experiment and maintains a slight improvement till it reaches the 100%. PForest and ST generally perform better on data sets with small number of classes. As they move towards the 100% chunk they tend to enhance their performance gradually. WEASEL and CBoss on the contrary learn better on the data sets with larger number of classes at the earlier chunks. PForest cannot beat the baseline classifier on the group of (11+) across all chunks. Since the group consists of only two data sets, we cannot generalize this conclusion.

Train Size	#Data sets	CBoss	PForest	ST	TSF	WEASEL
1-50	14	92.86%	92.86%	100.00%	92.86%	92.86%
51-100	9	100.00%	100.00%	100.00%	100.00%	100.00%
101-250	7	100.00%	71.43%	100.00%	100.00%	100.00%
251-500	4	75.00%	100.00%	75.00%	75.00%	75.00%
501-1000	2	100.00%	100.00%	100.00%	100.00%	100.00%

Table 6.7: Performance of classifiers with respect to train size in comparison to baseline on 100% chunk data sets

Type	#Data sets	CBoss	PForest	ST	TSF	WEASEL
ECG	5	40.00%	100.00%	80.00%	80.00%	60.00%
SPECTRO	8	50.00%	37.50%	62.50%	100.00%	12.50%
SENSOR	13	53.85%	53.85%	46.15%	76.92%	61.54%
Traffic	1	0.00%	0.00%	0.00%	100.00%	0.00%
DEVICE	5	80.00%	100.00%	100.00%	100.00%	60.00%
EPG	2	100.00%	100.00%	0.00%	100.00%	100.00%
HAR	1	0.00%	0.00%	0.00%	100.00%	0.00%
SOUND	1	100.00%	0.00%	100.00%	100.00%	100.00%

Table 6.8: Performance of classifiers with respect to problem type in comparison to baseline on 10% chunk data sets

Type	#Data sets	CBoss	PForest	ST	TSF	WEASEL
ECG	5	40.00%	100.00%	80.00%	80.00%	60.00%
SPECTRO8		62.50%	62.50%	75.00%	100.00%	37.50%
SENSOR	13	76.92%	76.92%	53.85%	92.31%	76.92%
Traffic	1	100.00%	100.00%	100.00%	100.00%	0.00%
DEVICE	5	100.00%	100.00%	100.00%	100.00%	80.00%
EPG	2	100.00%	100.00%	0.00%	100.00%	100.00%
HAR	1	0.00%	0.00%	0.00%	100.00%	100.00%
SOUND	1	100.00%	0.00%	100.00%	100.00%	100.00%

Table 6.9: Performance of classifiers with respect to problem type in comparison to baseline on 20% chunk data sets

Type	#Data sets	CBoss	PForest	ST	TSF	WEASEL
ECG	5	60.00%	40.00%	100.00%	80.00%	60.00%
SPECTRO8		87.50%	25.00%	87.50%	100.00%	37.50%
SENSOR	13	84.62%	76.92%	69.23%	92.31%	69.23%
Traffic	1	100.00%	100.00%	100.00%	100.00%	0.00%
DEVICE	5	100.00%	100.00%	100.00%	100.00%	80.00%
EPG	2	100.00%	100.00%	0.00%	100.00%	100.00%
HAR	1	100.00%	0.00%	100.00%	100.00%	100.00%
SOUND	1	100.00%	0.00%	100.00%	100.00%	100.00%

Table 6.10: Performance of classifiers with respect to problem type in comparison to baseline on 30% chunk data sets

Type	#Data sets	CBoss	PForest	ST	TSF	WEASEL
ECG	5	80.00%	80.00%	100.00%	80.00%	80.00%
SPECTRO	8	100.00%	100.00%	100.00%	100.00%	100.00%
SENSOR	13	92.31%	100.00%	92.31%	92.31%	92.31%
Traffic	1	100.00%	100.00%	100.00%	100.00%	100.00%
DEVICE	5	100.00%	100.00%	100.00%	100.00%	100.00%
EPG	2	100.00%	100.00%	100.00%	100.00%	100.00%
HAR	1	100.00%	0.00%	100.00%	100.00%	100.00%
SOUND	1	100.00%	0.00%	100.00%	100.00%	100.00%

Table 6.11: Performance of classifiers with respect to problem type in comparison to baseline on 100% chunk data sets

#Classes	#Data sets	CBoss	PForest	ST	TSF	WEASEL
2	20	40.00%	55.00%	55.00%	90.00%	40.00%
3	6	83.33%	100.00%	66.67%	83.33%	83.33%
4-5	6	66.67%	66.67%	66.67%	100.00%	33.33%
6-10	2	100.00%	50.00%	50.00%	50.00%	100.00%
11+	2	50.00%	0.00%	50.00%	100.00%	50.00%

Table 6.12: Performance of classifiers with respect to number of classes in comparison to baseline on 10% chunk data sets

#Classes	#Data sets	CBoss	PForest	ST	TSF	WEASEL
2	20	65.00%	85.00%	70.00%	95.00%	60.00%
3	6	83.33%	100.00%	50.00%	83.33%	66.67%
4-5	6	83.33%	66.67%	83.33%	100.00%	66.67%
6-10	2	100.00%	50.00%	50.00%	100.00%	100.00%
11+	2	50.00%	0.00%	50.00%	100.00%	100.00%

Table 6.13: Performance of classifiers with respect to number of classes in comparison to baseline on 20% chunk data sets

#Classes	#Data sets	CBoss	PForest	ST	TSF	WEASEL
2	20	80.00%	75.00%	90.00%	95.00%	65.00%
3	6	83.33%	66.67%	66.67%	83.33%	50.00%
4-5	6	100.00%	33.33%	66.67%	100.00%	66.67%
6-10	2	100.00%	50.00%	50.00%	100.00%	50.00%
11+	2	100.00%	0.00%	100.00%	100.00%	100.00%

Table 6.14: Performance of classifiers with respect to number of classes in comparison to baseline on 30% chunk data sets

#Classes	#Data sets	CBoss	PForest	ST	TSF	WEASEL
2	20	95.00%	100.00%	95.00%	95.00%	95.00%
3	6	83.33%	83.33%	100.00%	83.33%	83.33%
4-5	6	100.00%	100.00%	100.00%	100.00%	100.00%
6-10	2	100.00%	100.00%	100.00%	100.00%	100.00%
11+	2	100.00%	0.00%	100.00%	100.00%	100.00%

Table 6.15: Performance of classifiers with respect to number of classes in comparison to baseline on 100% chunk data sets

### 6.1.3 Results Within Classifier

In our first experiment, we carry out a comparison between copies of the same classifier but trained using different chunks of data; to explore how their performances evolve in the early classification context. We combine results that are based on both  $F_\beta - \text{measure}$  and balanced accuracy. We exclude the baseline classifier from this comparison.

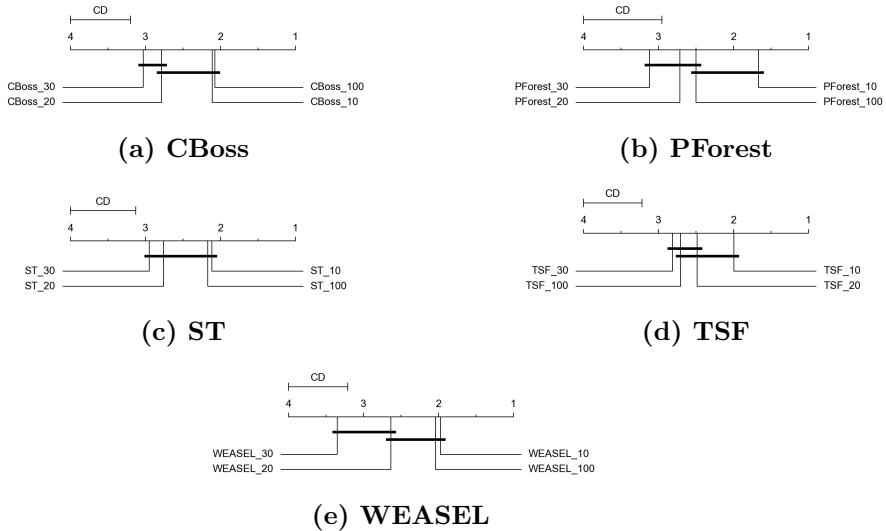


Figure 6.2: Critical difference diagrams between versions of the same classifier trained on different chunks ( $F_\beta$ )

Figure 6.2 shows the critical difference diagram based on the  $F_\beta$  for each classifier separately. We compare each classifier on all the data sets that it finished across all the different chunks, regardless of the other classifiers. The best performing algorithms lie on the right side of the diagram and the worse performing algorithms on the left side. A bar is used to connect a clique of classifiers which are not significantly different based on their ranking.

All classifiers show no significant difference between their 100% classifier version and the 10% version based on their  $F_\beta - \text{measure}$  scores. For CBoss, CBoss\_100 comes in the first place and CBoss\_10 comes in the second place. While for PForest, ST and WEASEL; their respective 10% versions; ST\_10, PForest\_10 and WEASEL\_10 rank first. Although TSF shows no significant difference as well, but its 100% version TSF\_100 comes in the third place. The 20% and 30% versions for all the classifiers always show no difference in ranking.

If we switch to the critical difference diagrams of accuracies, we find rather different results. Figure 6.3 shows the critical difference diagrams based on the balanced accuracy for each classifier separately.

For all the classifiers, the 100% and the 10% version are always significantly different in their ranking. WEASEL\_100, CBoss\_100 and PForest\_100 are all significantly better than any of their other versions. TSF\_100 and ST\_100 are ranked better than

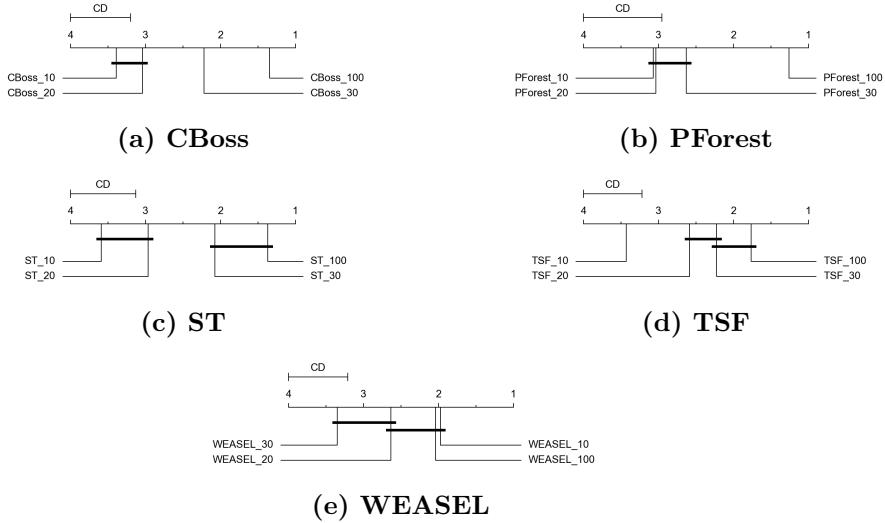


Figure 6.3: Critical difference diagrams between versions of the same classifier trained on different chunks (balanced accuracy)

their 30% versions TSF\\_30 and ST\\_30, but they are not significantly different. This conforms with the concept of eTSC that as more data becomes available the better a classifier performs and also demonstrates the role that earliness plays in the scores of  $F_\beta - measure$ .

#### 6.1.4 Results on Runtime Duration

Another aspect we investigate for our experiments is the runtime of the classifiers. It is not possible to carry out a comparison between the different algorithms we have used; because of the different adaptations we have applied to them. For some classifiers, we have used an enhanced version which allows setting a contract time for feature extraction phases; like in the case of CBoss and ST. This parameter allows us to restrict the amount of time given for the classifier. For PForest, we have excluded two distance measures; euclidean distance because it cannot handle the early classification context and TWE distance because it is the slowest elastic distance measure. These enhancements make our results not representative for the original classifiers runtime performance. Instead we provide information about the CPU time that was needed by each classifier to complete training on the data sets. Figures 6.4, 6.5, 6.6, 6.7 and 6.8 show the CPU time, in seconds, for the classifiers on all the used data sets. The values represent the training time of the best model selected from the 5 cross validation folds. If no cross validation was carried out, then the value represents the CPU time for the training process directly.

TSF and WEASEL have short training durations. TSF doesn't use any complex features and instead depends on simple summary statistics on intervals. The stochastic selection of intervals speeds its learning process even faster. WEASEL on the other hand is designed for speed and quality, but this comes on the expense of memory.

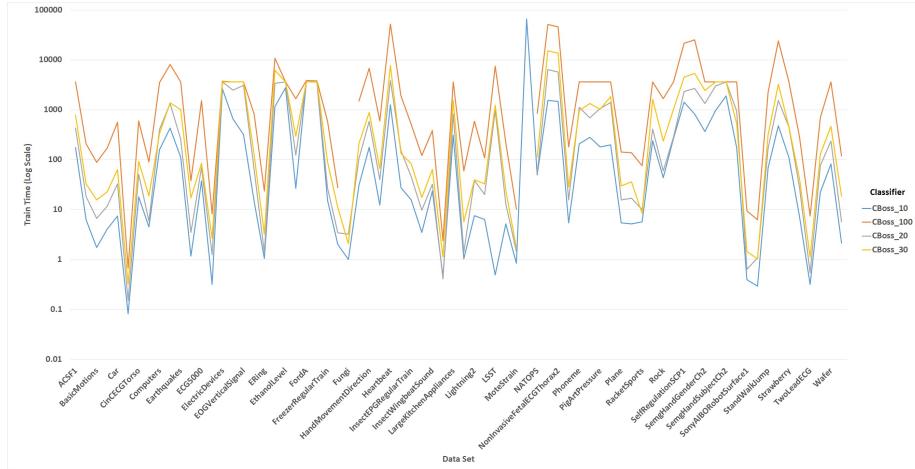


Figure 6.4: Train Time (CPU Time in Log Scale) for CBoss across all chunks

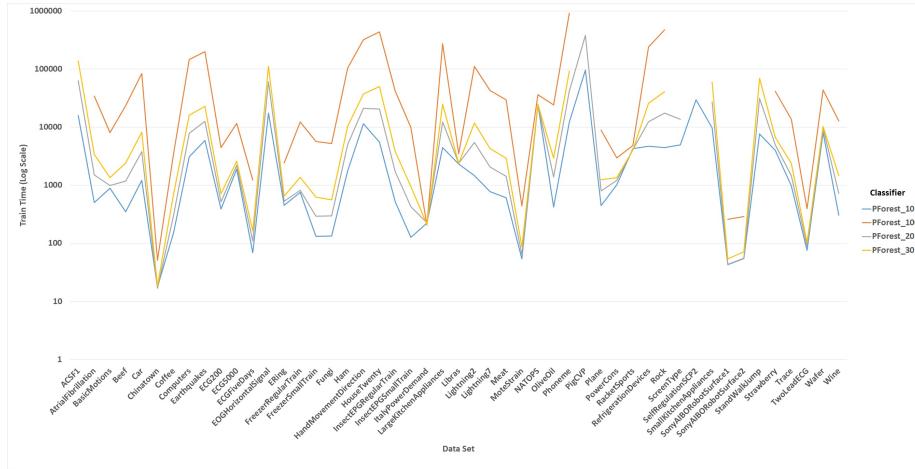
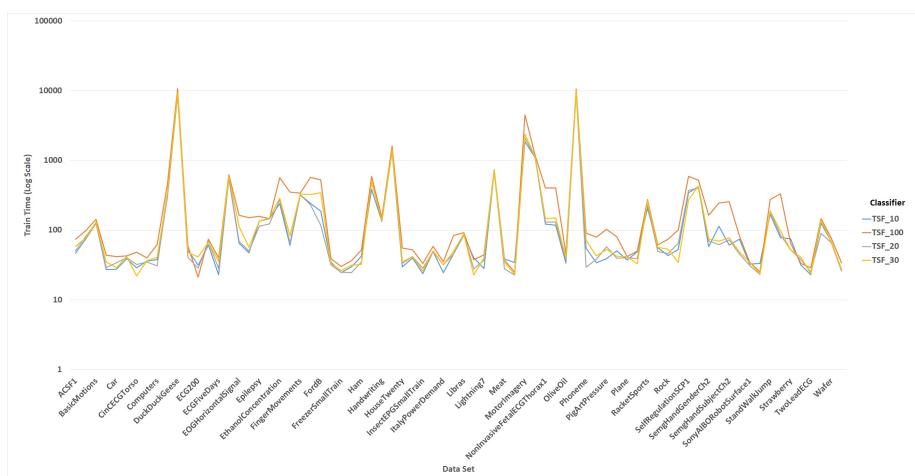
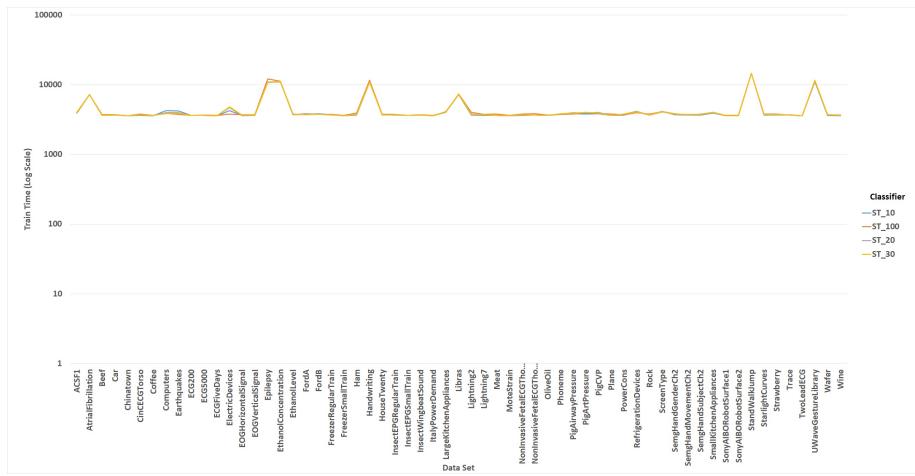


Figure 6.5: Train Time (CPU Time in Log Scale) for PForest across all chunks

Although WEASEL takes short training times, it has failed to complete on some data sets due to its large memory print. ST is known to be a strong but rather a significantly slow classifier. The use of the contractable version of ST clearly pays off, it stabilizes the training time across the different chunks by limiting the time during which ST looks for shapelets. The slowest classifier among all experiments was PForest. We provide train time results and breakdown by data characteristics in appendix A.4



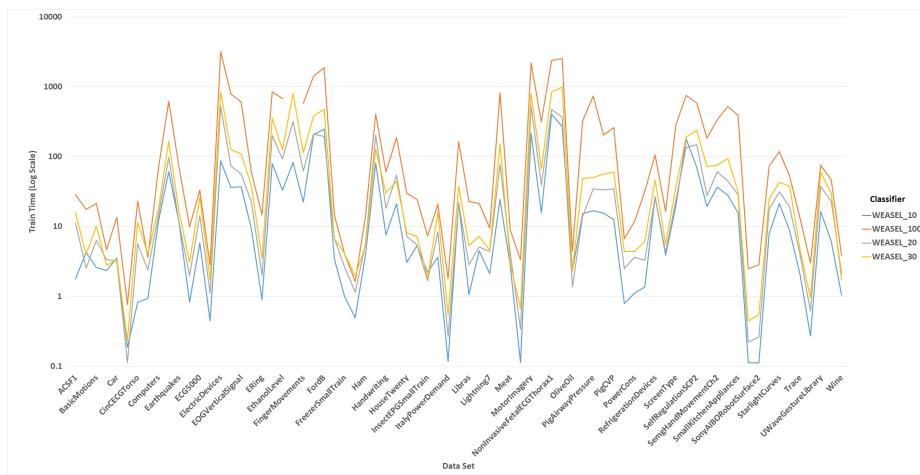


Figure 6.8: Train Time (CPU Time in Log Scale) for WEASEL across all chunks

## 6.2 Recommender Results

We divide our results for the recommender into two main parts. The first part presents the results for the chunk learners, while the second part presents the results for the final recommendation of the whole recommender. In the chunks learners results part, we discuss the performance of each chunk learner and the important features for building the random forests. In the final recommendation part, we discuss the overall performance of the recommender based on its ability to correctly suggest good performing classifiers.

### 6.2.1 Are the chunk learners able to learn F-scores for the classifiers ?

Table 6.16 shows the overall performance of the chunk learners across all 50 runs. In general, the chunk learners are able to learn about the performance of classifiers on the different chunks, but their mean errors are relatively high. As evident from the results, there is an inverse relation between the performance of the chunk learners and the revealed%. The chunk learner that trains on the 10% data is able to score the best scores across all the chunk learners. The more revealed data, the harder it becomes for the chunk learners to accurately predict performance of classifiers.

Revealed%	Avg MAE	Avg RMSE
10	4.82	6.71
20	9.57	12.42
30	12.03	15.46
100	15.58	19.52

Table 6.16: Overall performance of chunk learners in predicting actual  $F_\beta$  scores over 50 runs

Figure 6.9 gives a more detailed view of the distribution of RMSE for the chunk learners. It is clear that not only the curves of the higher chunk learner shift towards higher error values, but also the spread of the values is wider. This indicates that few of the predictions made by the lower chunk learner are significantly worse than their actual values. While the higher chunk learner is generally worse and makes larger mistakes by predicting values that are clearly worse than the actual performance of the classifiers.

Figures 6.10 and 6.11 visualize the MEA and RMSE of each chunk learner per classifier. It is evident from the results that for CBoss, ST and WEASEL; the chunk learners are able to score better results as more data is revealed. It is the opposite scenario for PForest. The first 3 chunk learners score the least errors for PForest, but on the 100% data the errors significantly increase to be the worst among all. This explains the the long tail that appears in the 100% chunk learner RMSE histogram due to the large differences between the actual and predicted F-scores for PForest on the 100% data.

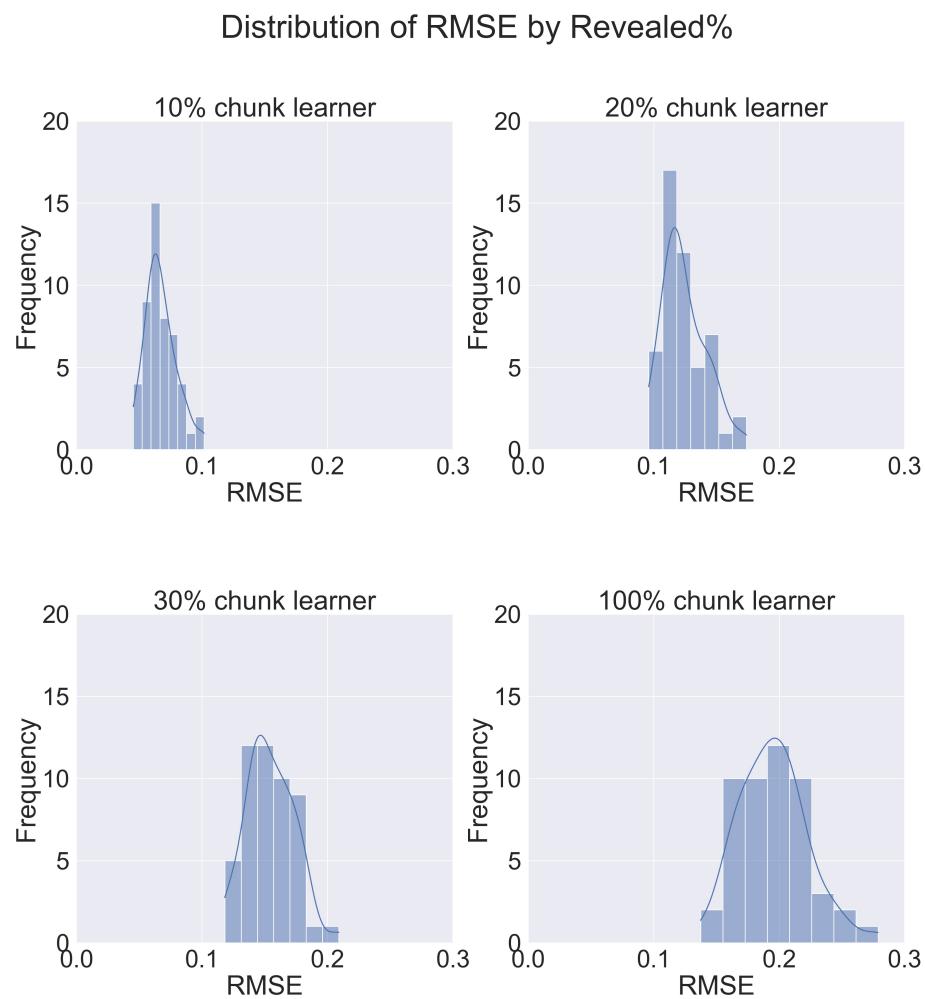


Figure 6.9: Histograms for distribution of RMSE per chunk over 50 runs

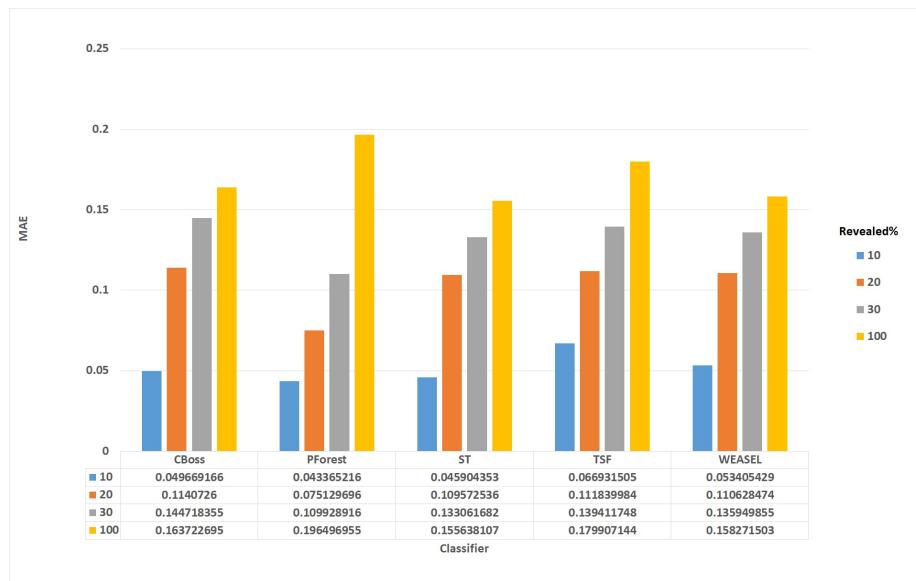


Figure 6.10: Avg MAE for classifiers per chunk over 50 runs

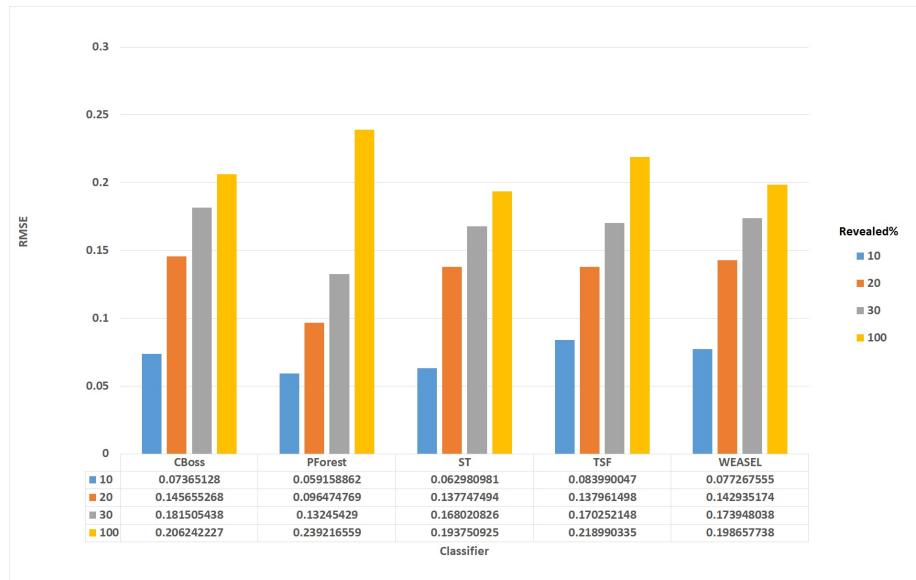


Figure 6.11: Avg RMSE for classifiers per chunk over 50 runs

### 6.2.2 Which features contribute the most to the F-Score Prediction ?

One of the main reasons why we chose to use random forests for the chunk learners, is that they give a notion of importance for features that were used to build the trees. This helps understand what features play role in the prediction of scores for the classifiers for each chunk.

Figure 6.12 shows the average importance of the features across all the 50 runs. The feature importance is calculated using the default method from the *sklearn* library; the mean decrease of impurity. For each selected feature in the internal nodes of the regression trees, the impurity of the split is calculated using variance reduction. After the whole forest is built, the mean variance reduction for each feature is calculated over all the trees.

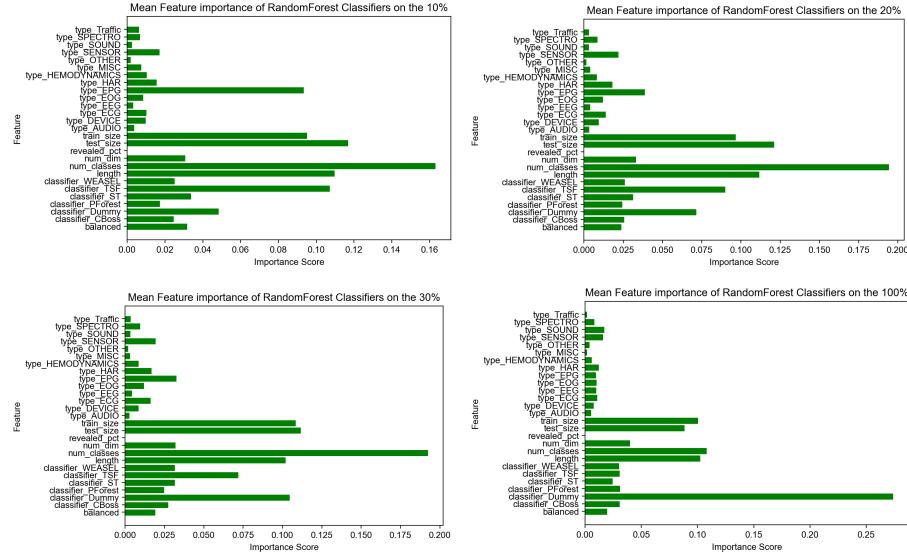


Figure 6.12: Mean Feature Importance per chunk across 50 runs

For the lower chunk learners, the main metadata of the data sets; number of classes, length, train size and test size have the greatest importances among all the features. The number of classes persistently has the highest feature importance over all other features. It is also visible that TSF has high importance on the 10% chunk results, which conforms with our CD diagram on the 10% data; due to its significantly better performance. This importance of TSF fades away as more data is revealed for the classifiers, while on the other hand the importance of the dummy classifier increases. For the 100% chunk learner, the Dummy classifier has the highest importance. This conforms with our results from the CD diagram from Figure 6.1, where all the classifiers are much better than the Dummy classifier; which makes it easy to purely separate from all other results. This then leads the importance of all the other classifiers to become closer to each other.

### 6.2.3 Are the final recommendations of the framework consistent with actual results ?

At the end, after all the experiments are finished, the final result of our framework is the suggestions made by the recommender for the user. For this analysis part, we focus on whether the chunk learners were able to correctly recommend classifiers based on their performance in comparison to the dummy baseline classifier. A good classifier would be one that can beat the baseline, while a bad classifier is one that doesn't.

By aggregating the performance metrics for each chunk learner across the 50 experiments, we can get an overview of their performance. Table 6.17 shows the average performance of the 4 chunk learners over all experiments.

Revealed%	Avg Accuracy	Avg Recall	Avg Precision
10	73.62	0.9520	0.6920
20	80.19	0.9597	0.7769
30	86.15	0.9925	0.8346
100	95.80	0.1000	0.9481

Table 6.17: Overall Performance of chunk learners in predicting good and bad performers over 50 runs

In general, our chunk learners are able to achieve acceptable results in terms of accuracy specially on the full length data. But accuracy alone is not the best measure to depend on for assessing the results, thus we have a look into the other evaluation metrics.

By looking into recall, we score very high results. This means that our chunk learners are able to detect a high majority of the existing good classifiers. For example, It achieves a recall of 1 on the 100% data. This indicates that on the full length data, the recommender is able to suggest all the good classifiers. Although this is a quite impressive result, but we cannot contribute the score to the ability of the classifier to learn on the data. From our CD diagrams results, we know that the dummy classifier is significantly worse than all the classifiers on the 100% chunk. We also know from the F-Score results that the chunk learner on the 100% data has the highest average RMSE values across all the runs. Thus, we believe that the reason such a high score is achieved; is because the dummy classifier is the baseline that decides whether a certain classifier is good or bad. The dummy classifier is a simple classifier that predicts the majority class and it allows for a wide span of error without affecting the results of the recommendation. So even if we underestimate the F-scores of the classifiers compared to their actuals, we are still able to get the final suggestions correctly; because their F-scores are way better than that of the dummy classifier.

On the other hand, we can see that our results are not as good when considered from a precision point of view. This goes back to the nature of the precision measure which focuses on the False Positives in the data. The same reason that makes the dummy classifier weaker on the full length data, is the one that makes it hard to beat on the shorter length data. The dummy classifier always predicts the majority class regardless of the length of data. This means that its balanced accuracy scores are always the same and their total F-scores keep getting better on the lower chunks; because their earliness scores improve as well. From our CD diagrams results, we know that the dummy classifier is not significantly different than most of the classifiers on

the lower chunks. Thus, the residuals between the actual and predicted F-scores have larger effect and the overestimation of the performance of classifiers is more prominent.

Table 6.18 gives a breakdown of the misestimations that the chunk learners did through out the experiments on the classifier level. From the previous analysis of precision and recall, it is quite obvious that our chunk learners tend to overestimate the performance of the classifiers than to underestimate them. The misestimation is clearer on the lower chunk predictions, but less obvious on the higher chunks. This can obviously be recognized in Figure 6.13 by comparing the vertical distances between the actual and predicted F-scores for each classifier. TSF is the least misestimated classifier; it is the least overestimated classifier and almost never underestimated. ST is the second least overestimated classifier, while PForest is the most overestimated.

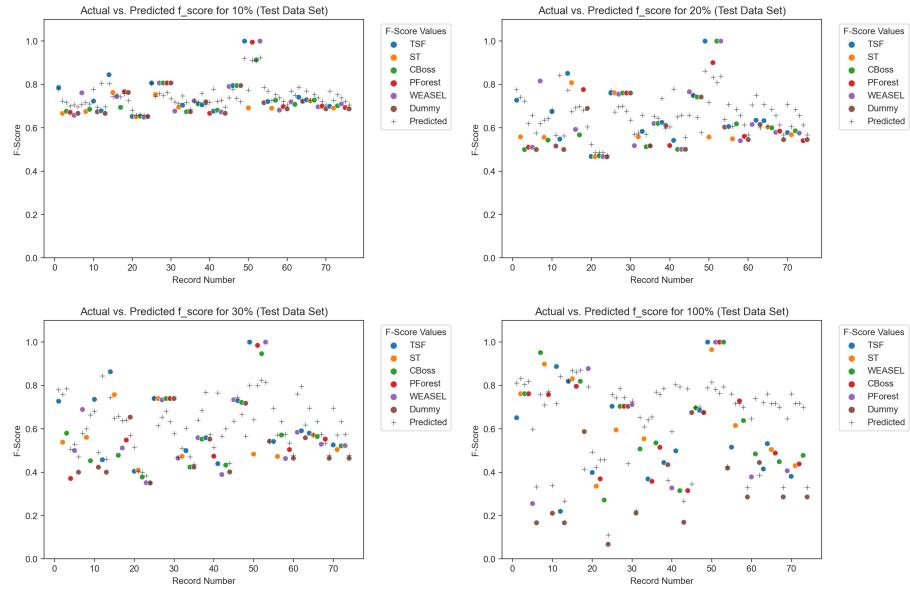


Figure 6.13: Actual vs. Predicted F-score for chunk learners (experiment seed = 22)

Classifier	Revealed %	Total Tests	#Overest- imations	#Underest- imations	Overest- imations%	Underest- imations%	#Misest- imations	Misest- imations%
CBoss	10	726	226	36	31.13%	4.96%	262	36.09%
CBoss	20	717	204	20	28.45%	2.79%	224	31.24%
CBoss	30	717	108	2	15.06%	0.28%	110	15.34%
CBoss	100	717	46	0	6.42%	0.00%	46	6.42%
PForest	10	477	204	16	42.77%	3.35%	220	46.12%
PForest	20	477	104	24	21.80%	5.03%	128	26.83%
PForest	30	454	154	5	33.92%	1.10%	159	35.02%
PForest	100	417	32	0	7.67%	0.00%	32	7.67%
ST	10	618	165	31	26.70%	5.02%	196	31.72%
ST	20	618	155	26	25.08%	4.21%	181	29.29%
ST	30	618	78	3	12.62%	0.49%	81	13.11%
ST	100	618	19	0	3.07%	0.00%	19	3.07%
TSF	10	745	74	4	9.93%	0.54%	78	10.47%
TSF	20	745	46	2	6.17%	0.27%	48	6.44%
TSF	30	745	26	1	3.49%	0.13%	27	3.62%
TSF	100	745	35	0	4.70%	0.00%	35	4.70%
WEASEL	10	731	277	31	37.89%	4.24%	308	42.13%
WEASEL	20	731	193	25	26.40%	3.42%	218	29.82%
WEASEL	30	731	169	10	23.12%	1.37%	179	24.49%
WEASEL	100	726	34	0	4.68%	0.00%	34	4.68%

Table 6.18: Misestimation results per classifier per chunk learner over 50 runs



## Chapter 7

# Conclusion and Future Work

The primary goal of this thesis is to compare the performance of TSCAs in an early classification context. A context which is inspired by the problem of eTSC. Our research questions investigated adapting workflows for the early classification context, how to evaluate the classifiers in the context while including earliness and performance and how to evaluate the proposed solution.

We introduced a new framework that simulates the early classification context and compares TSCAs on the UCR/UEA data archive using a testbed. The framework is then evaluated by building a recommender which learns on the results and predicts for unseen data sets which are the best classifiers to use.

Our proposed implementation of the context cannot be handled by all classifiers. They either have to be able to compare instances of different lengths or involve a feature extraction step that creates a fixed feature vector size. Classifiers that are bound by the shape of data fail to operate in the context. Our final experiment results included the classifiers CBoss, PForest, ST, TSF and WEASEL compared on chosen data sets from the UCR/UEA time series data archives. Along with a baseline classifier that predict majority class.

We use the  $F_\beta$  – measure as an overall performance metric; as it combines both performance of classifier and earliness. The value of  $\beta$  plays a fundamental role in the results and should be adjusted according to the needs of the experiments. Using equal weights for performance and earliness, overemphasizes the importance of earliness, which consequently overcompensates for the lower performace scores.

Overall, our testbed results for the across classifiers comparison show that TSF is, on average, better than the other 4 classifiers on early chunks of data. It is clearly superior on the 10% chunk data, while the others are no better than the majority class baseline. As more data is revealed to learn on, the other classifiers; CBoss, WEASEL, ST and PForest, start to catch up with TSF. On the full length data, the 5 classifiers are not significantly different from each other and are all better than the baseline, yet WEASEL has the highest rank.

For the within classifier comparison based on  $F_\beta$  – measure, our testbed results show that there is no statistical difference between the 10% version of any of the classifiers and their respective 100% versions. This contrasts with the within comparison

based on balanced accuracy, where the 100% versions are significantly better than the 10% versions for the classifiers. This means that the earliness factor plays a big role in shaping the results; it gives an advantage to the lower chunks that compensates for the lower performance. A research interest that we couldn't infer useful insights for is the effect of the different data characteristics on the performance of classifiers across the different chunks.

To evaluate the created framework, we built a recommender that predicts the performance of TSCAs for unseen data sets in the early classification context. The recommender is better in predicting the  $F_\beta$ -measure of classifiers for data sets on the lower chunks scoring an average 6.71% RMSE error, but gets worse when predicting the scores on higher chunks till the RMSE reaches 19.52%. The train size, test size, length and number of classes play an important role in predicting the chunk learners results. Other data set characteristics like; problem type and number of dimensions don't seem to have a precipitable relation with the results.

For every chunk, the recommender is good at judging whether a classifier is a good performer or not, despite the residuals between the actual and predicted  $F_\beta$ -measure values. This goes back to the central role that the baseline classifier plays in the recommendation part. The baseline classifier is generally worse than all TSCAs on all chunks, so it is easy to beat on many data sets. This causes our recall values to be very high, specially on the 100% chunk data when the baseline classifier is significantly worse than all the others.

There are multiple weaknesses in our study that we freely acknowledge. The selection criteria of the data set problem types is cherry picked to suit our motivating goal, the UCR/UEA data archives contains other problem types that we have not included which cover a variety of other time series problems. There are multiple components and parameters of the framework that have a great effect on the final results. For example, the evaluation metric to use for comparing the classifiers. We emphasize that  $F_\beta$ -measure is not the only suitable measure for assessing classifiers performance. But if used, the choice of weights has a great effect that should be considered in interpreting the results. In handling multivariate data sets, our choice for the column ensembling technique by repeating the same classifier for each dimension might not be the best way for evaluating the real performance of the classifiers on these data sets. Our used configurations for the classifiers and the adaptations applied might not be fair for all of them. For example, we use the contractable versions of CBoss and ST and used a 1 hour time contract. This might have hindered their actual performance on the data sets and might have scored better if given more time. Also the exclusion of the TWE elastic distance metric from the PForest classifier to meet time constraints, might have introduced a bias in the results for PForest.

We believe there are a lot of areas that are open for more research and improvement in the early classification context. Handling multivariate data sets is one of the areas. For instance, the inclusion of bespoke multivariate classifiers to learn on multivariate data sets while capturing relations between dimensions. Running experiments using multiple resampling of data instead of the default train/test split; to avoid uncertainty about the performance of classifiers on data sets. Deeper and more granular analysis of the features used by the classifiers; either as shapelets for ST or temporal importance curves for TSF. Using a different baseline in the recommender for example KNNDTW for deciding good and bad performing classifiers.

## Appendix A

# Appendix - Testbed Results

### A.1 Results Across Classifiers

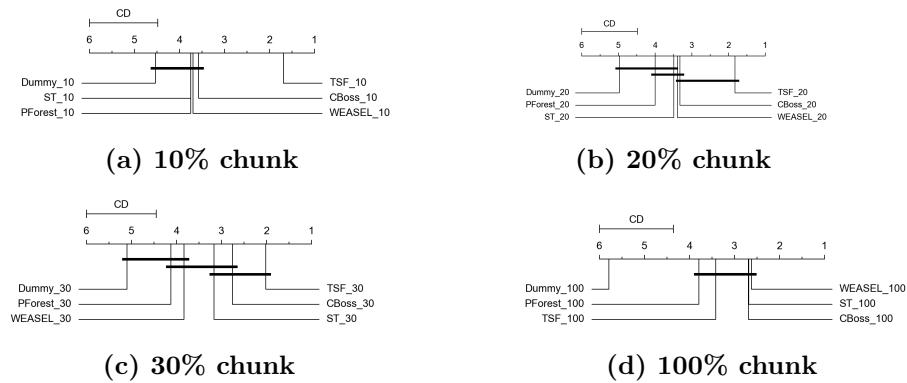


Figure A.1: Critical Difference diagrams per chunk across all classifiers

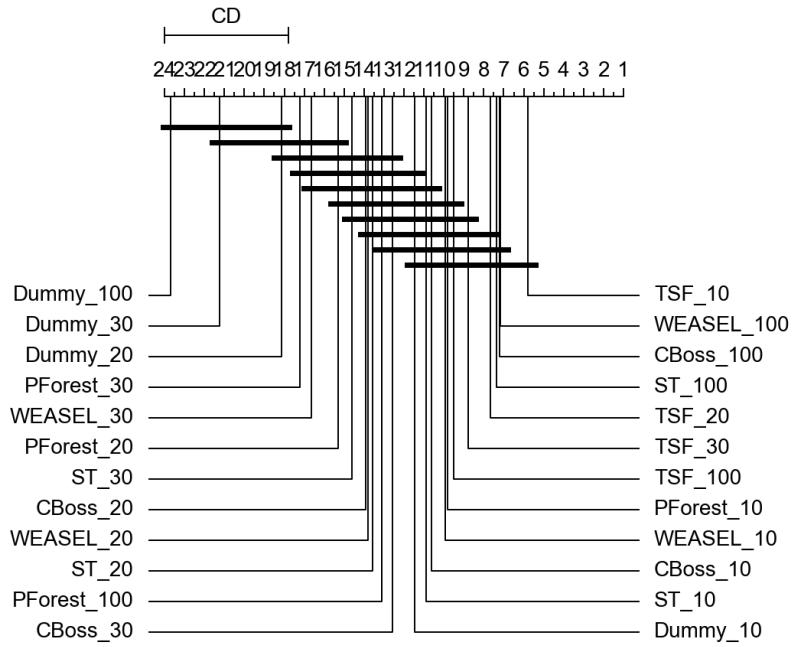
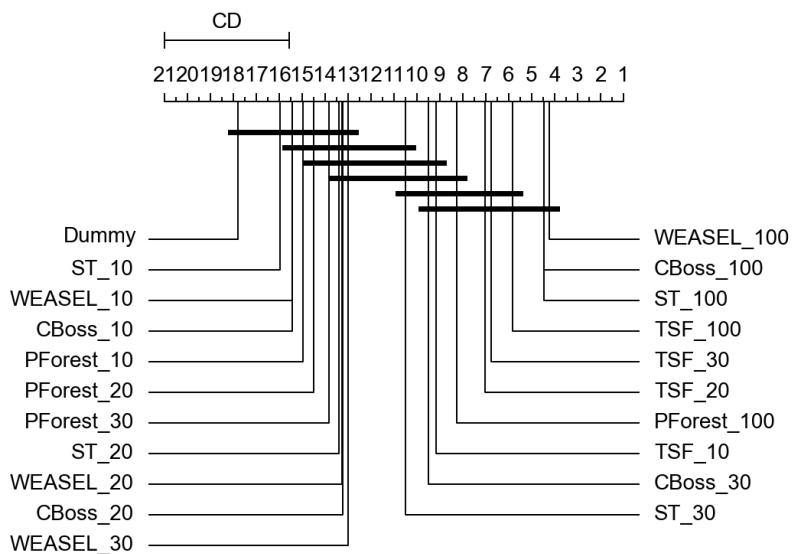
Figure A.2: Critical Difference diagram for all classifiers and all chunks ( $F_\beta$ )

Figure A.3: Critical Difference diagram for all classifiers and all chunks (balanced accuracy)

## A.2 Data Set characteristics and performance in early classification context

Train Length	#Data sets	CBoss	PForest	ST	TSF	WEASEL
1-50	3	33.33%	66.67%	33.33%	100.00%	33.33%
51-100	5	40.00%	80.00%	60.00%	100.00%	80.00%
101-250	7	85.71%	85.71%	100.00%	100.00%	57.14%
251-500	8	75.00%	87.50%	75.00%	100.00%	87.50%
501-1000	10	80.00%	80.00%	50.00%	80.00%	60.00%
1001+	3	100.00%	33.33%	66.67%	100.00%	66.67%

Table A.2: Performance of classifiers with respect to length in comparison to baseline on 20% chunk data sets

Train Length	#Data sets	CBoss	PForest	ST	TSF	WEASEL
1-50	3	100.00%	66.67%	100.00%	100.00%	33.33%
51-100	5	80.00%	80.00%	100.00%	100.00%	80.00%
101-250	7	100.00%	71.43%	100.00%	100.00%	71.43%
251-500	8	87.50%	50.00%	75.00%	100.00%	62.50%
501-1000	10	70.00%	60.00%	60.00%	80.00%	60.00%
1001+	3	100.00%	33.33%	66.67%	100.00%	66.67%

Table A.4: Performance of classifiers with respect to length in comparison to baseline on 30% chunk data sets

<b>Train Length</b>	<b>#Data sets</b>	<b>CBoss</b>	<b>PForest</b>	<b>ST</b>	<b>TSF</b>	<b>WEASEL</b>
1-50	3	100.00%	66.67%	100.00%	100.00%	100.00%
51-100	5	100.00%	100.00%	100.00%	100.00%	100.00%
101-250	7	100.00%	100.00%	100.00%	100.00%	100.00%
251-500	8	100.00%	100.00%	100.00%	100.00%	100.00%
501-1000	10	80.00%	90.00%	90.00%	80.00%	80.00%
1001+	3	100.00%	66.67%	100.00%	100.00%	100.00%

Table A.6: Performance of classifiers with respect to length in comparison to baseline on 100% chunk data sets

<b>Train Size</b>	<b>#Data sets</b>	<b>CBoss</b>	<b>PForest</b>	<b>ST</b>	<b>TSF</b>	<b>WEASEL</b>
1-50	14	50.00%	64.29%	50.00%	92.86%	57.14%
51-100	9	66.67%	55.56%	44.44%	77.78%	55.56%
101-250	7	57.14%	57.14%	71.43%	100.00%	42.86%
251-500	4	50.00%	75.00%	75.00%	75.00%	25.00%
501-1000	2	50.00%	50.00%	100.00%	100.00%	50.00%

Table A.8: Performance of classifiers with respect to train size in comparison to baseline on 10% chunk data sets

<b>Train Size</b>	<b>#Data sets</b>	<b>CBoss</b>	<b>PForest</b>	<b>ST</b>	<b>TSF</b>	<b>WEASEL</b>
1-50	14	71.43%	78.57%	64.29%	92.86%	71.43%
51-100	9	77.78%	88.89%	55.56%	100.00%	66.67%
101-250	7	71.43%	71.43%	71.43%	100.00%	85.71%
251-500	4	75.00%	75.00%	75.00%	75.00%	50.00%
501-1000	2	50.00%	50.00%	100.00%	100.00%	0.00%

Table A.10: Performance of classifiers with respect to train size in comparison to baseline on 20% chunk data sets

Train Size	#Data sets	CBoss	PForest	ST	TSF	WEASEL
1-50	14	92.86%	71.43%	85.71%	92.86%	71.43%
51-100	9	77.78%	55.56%	55.56%	100.00%	55.56%
101-250	7	85.71%	57.14%	100.00%	100.00%	85.71%
251-500	4	75.00%	50.00%	75.00%	75.00%	50.00%
501-1000	2	100.00%	50.00%	100.00%	100.00%	0.00%

Table A.12: Performance of classifiers with respect to train size in comparison to baseline on 30% chunk data sets

Train Size	#Data sets	CBoss	PForest	ST	TSF	WEASEL
1-50	14	92.86%	92.86%	100.00%	92.86%	92.86%
51-100	9	100.00%	100.00%	100.00%	100.00%	100.00%
101-250	7	100.00%	71.43%	100.00%	100.00%	100.00%
251-500	4	75.00%	100.00%	75.00%	75.00%	75.00%
501-1000	2	100.00%	100.00%	100.00%	100.00%	100.00%

Table A.14: Performance of classifiers with respect to train size in comparison to baseline on 100% chunk data sets

Type	#Data sets	CBoss	PForest	ST	TSF	WEASEL
ECG	5	40.00%	100.00%	80.00%	80.00%	60.00%
SPECTRO	8	50.00%	37.50%	62.50%	100.00%	12.50%
SENSOR	13	53.85%	53.85%	46.15%	76.92%	61.54%
Traffic	1	0.00%	0.00%	0.00%	100.00%	0.00%
DEVICE	5	80.00%	100.00%	100.00%	100.00%	60.00%
EPG	2	100.00%	100.00%	0.00%	100.00%	100.00%
HAR	1	0.00%	0.00%	0.00%	100.00%	0.00%
SOUND	1	100.00%	0.00%	100.00%	100.00%	100.00%

Table A.15: Performance of classifiers with respect to problem type in comparison to baseline on 10% chunk data sets

Type	#Data sets	CBoss	PForest	ST	TSF	WEASEL
ECG	5	40.00%	100.00%	80.00%	80.00%	60.00%
SPECTRO8		62.50%	62.50%	75.00%	100.00%	37.50%
SENSOR	13	76.92%	76.92%	53.85%	92.31%	76.92%
Traffic	1	100.00%	100.00%	100.00%	100.00%	0.00%
DEVICE	5	100.00%	100.00%	100.00%	100.00%	80.00%
EPG	2	100.00%	100.00%	0.00%	100.00%	100.00%
HAR	1	0.00%	0.00%	0.00%	100.00%	100.00%
SOUND	1	100.00%	0.00%	100.00%	100.00%	100.00%

Table A.16: Performance of classifiers with respect to problem type in comparison to baseline on 20% chunk data sets

Type	#Data sets	CBoss	PForest	ST	TSF	WEASEL
ECG	5	60.00%	40.00%	100.00%	80.00%	60.00%
SPECTRO8		87.50%	25.00%	87.50%	100.00%	37.50%
SENSOR	13	84.62%	76.92%	69.23%	92.31%	69.23%
Traffic	1	100.00%	100.00%	100.00%	100.00%	0.00%
DEVICE	5	100.00%	100.00%	100.00%	100.00%	80.00%
EPG	2	100.00%	100.00%	0.00%	100.00%	100.00%
HAR	1	100.00%	0.00%	100.00%	100.00%	100.00%
SOUND	1	100.00%	0.00%	100.00%	100.00%	100.00%

Table A.17: Performance of classifiers with respect to problem type in comparison to baseline on 30% chunk data sets

Type	#Data sets	CBoss	PForest	ST	TSF	WEASEL
ECG	5	80.00%	80.00%	100.00%	80.00%	80.00%
SPECTRO	8	100.00%	100.00%	100.00%	100.00%	100.00%
SENSOR	13	92.31%	100.00%	92.31%	92.31%	92.31%
Traffic	1	100.00%	100.00%	100.00%	100.00%	100.00%
DEVICE	5	100.00%	100.00%	100.00%	100.00%	100.00%
EPG	2	100.00%	100.00%	100.00%	100.00%	100.00%
HAR	1	100.00%	0.00%	100.00%	100.00%	100.00%
SOUND	1	100.00%	0.00%	100.00%	100.00%	100.00%

Table A.18: Performance of classifiers with respect to problem type in comparison to baseline on 100% chunk data sets

#Classes	#Data sets	CBoss	PForest	ST	TSF	WEASEL
2	20	40.00%	55.00%	55.00%	90.00%	40.00%
3	6	83.33%	100.00%	66.67%	83.33%	83.33%
4-5	6	66.67%	66.67%	66.67%	100.00%	33.33%
6-10	2	100.00%	50.00%	50.00%	50.00%	100.00%
11+	2	50.00%	0.00%	50.00%	100.00%	50.00%

Table A.19: Performance of classifiers with respect to number of classes in comparison to baseline on 10% chunk data sets

#Classes	#Data sets	CBoss	PForest	ST	TSF	WEASEL
2	20	65.00%	85.00%	70.00%	95.00%	60.00%
3	6	83.33%	100.00%	50.00%	83.33%	66.67%
4-5	6	83.33%	66.67%	83.33%	100.00%	66.67%
6-10	2	100.00%	50.00%	50.00%	100.00%	100.00%
11+	2	50.00%	0.00%	50.00%	100.00%	100.00%

Table A.20: Performance of classifiers with respect to number of classes in comparison to baseline on 20% chunk data sets

#Classes	#Data sets	CBoss	PForest	ST	TSF	WEASEL
2	20	80.00%	75.00%	90.00%	95.00%	65.00%
3	6	83.33%	66.67%	66.67%	83.33%	50.00%
4-5	6	100.00%	33.33%	66.67%	100.00%	66.67%
6-10	2	100.00%	50.00%	50.00%	100.00%	50.00%
11+	2	100.00%	0.00%	100.00%	100.00%	100.00%

Table A.21: Performance of classifiers with respect to number of classes in comparison to baseline on 30% chunk data sets

#Classes	#Data sets	CBoss	PForest	ST	TSF	WEASEL
2	20	95.00%	100.00%	95.00%	95.00%	95.00%
3	6	83.33%	83.33%	100.00%	83.33%	83.33%
4-5	6	100.00%	100.00%	100.00%	100.00%	100.00%
6-10	2	100.00%	100.00%	100.00%	100.00%	100.00%
11+	2	100.00%	0.00%	100.00%	100.00%	100.00%

Table A.22: Performance of classifiers with respect to number of classes in comparison to baseline on 100% chunk data sets

### A.3 Results Within Classifier

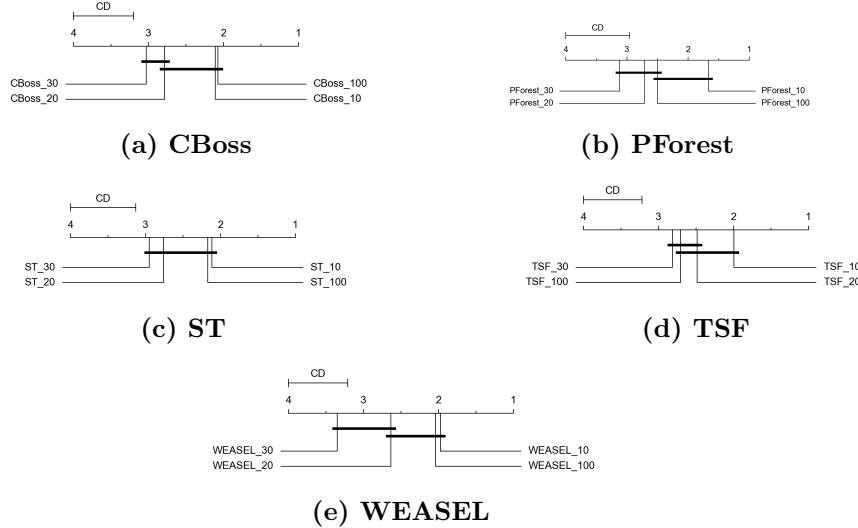


Figure A.4: Crtitical difference diagrams between versions of the same classifier trained on different chunks ( $F_\beta$ )

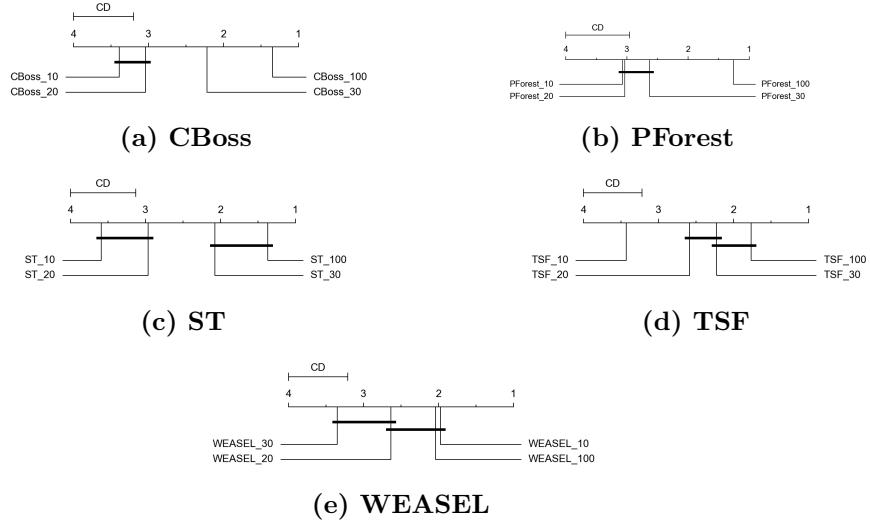


Figure A.5: Crtitical difference diagrams between versions of the same classifier trained on different chunks (balanced accuracy)

#### **A.4 Results on Runtime Duration**

Table A.23: Train Time (CPU Time) for CBoss, PForest and ST on 77 data sets for all chunks

Data set	CBoss				PForest				ST			
	10%	20%	30%	100%	10%	20%	30%	100%	10%	20%	30%	100%
ACSF1	175	428	779	3601	15972	63002	137622		3969	3991	4044	3919
AtrialFibrillation	6	18	33	207	506	1515	3381	34266	7205	7214	7213	7248
BasicMotions	2	7	16	88	894	993	1346	8063				
Beef	4	12	23	173	350	1178	2413	23595	3606	3630	3641	3715
Car	7	33	63	559	1209	3762	8142	83360	3647	3662	3676	3723
Chinatown	0	0	0	1	18	17	19	51	3603	3603	3603	3603
CinCECGTorso	18	51	93	596					3629	3654	3712	3804
Coffee	4	6	19	91	152	303	693	3456	3604	3606	3607	3628
Computers	159	407	334	3532	3071	7870	16136	146384	4242	3948	3852	3883
Cricket	433	1350	1387	8110								
DuckDuckGeese												
Earthquakes	112	261	1006	3603	5893	12552	22818	198829	4182	3982	3862	3750
ECG200	1	3	17	38	388	517	713	4431	3608	3611	3613	3626
ECG5000	38	70	85	1519	1911	2200	2594	11566	3638	3642	3637	3687
ECGFiveDays	0	1	3	8	69	109	165	1240	3602	3603	3605	3608
ElectricDevices	2622	3642	3634	3749					4257	4697	4830	3789
EOGHorizontalSignal	661	2487	3604	3618	17477	60769	112039		3608	3604	3673	3707
EOGVerticalSignal	317	3114	3604	3606					3609	3609	3670	3704
Epilepsy	15	63	126	826					10869	10916	10988	12083
ERing	1	1	3	24	456	534	648	2436				
EthanolConcentration	1154	3376	6188	10786					10985	11051	11105	11134
EthanolLevel	2774	3604	3606	3631					3691	3701	3775	3702
FaceDetection												
FingerMovements	27	124	293	1662								
FordA	3611	3607	3643	3832					3760	3816	3707	3841
FordB	3636	3711	3675	3810					3829	3800	3771	3733

Table A.23: Train Time (CPU Time) for CBoss, PForest and ST on 77 data sets for all chunks (continued)

Data set	CBoss				PForest				ST			
	10%	20%	30%	100%	10%	20%	30%	100%	10%	20%	30%	100%
FreezerRegularTrain	15	27	90	592	755	816	1379	12258	3715	3698	3661	3739
FreezerSmallTrain	2	3	11	28	131	296	627	5683	3605	3605	3604	3610
Fungi	1	3	2		133	297	564	5261				
Ham	31	105	211	1505	1772	5044	10404	104283	3665	3740	3765	3909
HandMovementDirection	177	593	889	6823	11414	21094	37411	321914				
Handwriting	12	40	67	595					10858	10932	10987	11578
Heartbeat	1257	3866	7667	52348								
HouseTwenty	28	152	134	1897	5564	20665	50006	435853	3697	3706	3680	3761
InsectEPGRegularTrain	16	46	83	514	524	1736	3818	41925	3674	3697	3696	3767
InsectEPGSmallTrain	3	10	17	122	126	426	956	9921	3607	3608	3614	3629
InsectWingbeatSound	24	33	64	382					3650	3678	3678	3714
ItalyPowerDemand	0	0	1	2	223	227	211	206	3604	3603	3602	3608
LargeKitchenAppliances	310	824	1544	3604	4433	12251	24973	270857	4106	4014	4038	4027
Libras	1	1	6	60	2340	2382	2390	3495	7211	7235	7259	7363
Lightning2	8	38	40	593	1478	5459	11789	111426	3665	3721	3792	3997
Lightning7	6	20	33	110	780	2058	4271	42482	3636	3672	3723	3757
LSST	0	956	1247	7538								
Meat	5	13	24	219	608	1428	2924	29454	3616	3626	3640	3781
MoteStrain	1	1	2	10	54	65	88	441	3602	3602	3604	3616
MotorImagery	65994											
NATOPS	49	52	114	851	24792	24483	24704	36247				
NonInvasiveFetalECGThorax1	1538	6368	15200	51086					3645	3678	3707	3812
NonInvasiveFetalECGThorax2	1479	5761	13591	45878					3649	3679	3680	3846
OliveOil	5	16	29	180	425	1389	2940	24060	3613	3615	3630	3674
PEMS-SF												
Phoneme	207	1134	969	3602	12344	42219	91768	915733	3760	3811	3815	3741

Table A.23: Train Time (CPU Time) for CBoss, PForest and ST on 77 data sets for all chunks (continued)

Data set	CBoss				PForest				ST						
	10%	20%	30%	100%	10%	20%	30%	100%	10%	20%	30%	100%			
PigAirwayPressure	283	688	1338	3601					3830	3967	3869	3777			
PigArtPressure	181	1049	1033	3601					3798	3936	4020	3964			
PigCVP	197	1405	1831	3601	96143	380569			3843	4007	3940	3886			
Plane	5	16	30	143	447	800	1248	8822	3606	3617	3632	3781			
PowerCons	5	17	36	138	1014	1218	1361	2969	3657	3792	3774	3712			
RacketSports	6	10	8	76	4296	4130	4072	4806							
RefrigerationDevices	239	410	1608	3605	4721	12419	25695	239726	4054	4148	4046	3972			
Rock	44	59	237	1678	4424	17372	41005	475073	3648	3649	3689	3789			
ScreenType	269	283	1025	3606	4964	13757			4149	4072	4118	4110			
SelfRegulationSCP1	1424	2340	4519	21574											
SelfRegulationSCP2	825	2693	5369	25174	29494										
SemgHandGenderCh2	367	1331	2429	3605					3695	3789	3774	3833			
SemgHandMovementCh2	963	2976	3605	3598					3647	3687	3731	3681			
SemgHandSubjectCh2	1900	3603	3611	3610					3648	3802	3742	3753			
SmallKitchenAppliances	175	941	525	3602	9800	27216	59998		3934	4007	3996	4022			
SonyAIBORobotSurface1	0	1	1	9	43	43	54	260	3603	3603	3604	3607			
SonyAIBORobotSurface2	0	1	1	6	55	56	72	292	3603	3605	3603	3607			
StandWalkJump	69	169	318	2246	7699	31252	69596		14443	14471	14491	14539			
StarlightCurves	478	1542	3258	24151					3676	3750	3752	3792			
Strawberry	112	453	448	3601	4020	4948	6687	41484	3649	3683	3741	3781			
Trace	7	27	42	284	977	1472	2405	13619	3654	3690	3657	3679			
TwoLeadECG	0	1	1	7	76	90	104	394	3602	3602	3602	3604			
UWaveGestureLibrary	22	77	129	698					11015	11399	11599	11477			
Wafer	82	235	466	3605	8088	9642	10263	43869	3643	3696	3711	3707			
Wine	2	6	19	120	305	725	1470	12827	3604	3620	3630	3676			

Table A.24: Train Time (CPU Time) for TSF and WEASEL on 77 data sets for all chunks

Data set	TSF				WEASEL			
	10%	20%	30%	100%	10%	20%	30%	100%
ACSF1	47	50	58	74	2	11	16	29
AtrialFibrillation	75	81	77	98	4	3	4	18
BasicMotions	126	128	131	143	3	6	10	21
Beef	27	29	35	44	2	3	3	5
Car	28	34	29	42	4	3	3	14
Chinatown	39	40	42	43	0	0	0	1
CinCECGTorso	29	32	22	49	1	6	12	23
Coffee	36	35	36	40	1	2	4	4
Computers	38	31	41	62	12	15	18	69
Cricket	293	308	337	450	61	98	167	623
DuckDuckGeese	8780	8996	9323	10706				
Earthquakes	52	41	48	59	12	12	17	71
ECG200	32	28	41	21	1	2	3	10
ECG5000	61	69	67	74	6	14	26	33
ECGFiveDays	23	28	36	40	0	1	2	3
ElectricDevices	557	567	590	622	87	524	836	3201
EOGHorizontalSignal	65	69	112	164	36	74	127	786
EOGVerticalSignal	47	50	57	151	37	57	109	611
Epilepsy	135	114	135	157	9	23	39	61
ERing	145	123	149	146	1	2	4	15
EthanolConcentration	241	271	287	566	80	202	355	844
EthanolLevel	61	65	83	350	33	93	126	679
FaceDetection					83	321	806	
FingerMovements	323	321	330	341	23	63	114	583
FordA	242	232	326	571	206	210	379	1425
FordB	188	118	342	527	250	195	480	1884

Table A.24: Train Time (CPU Time) for TSF and WEASEL on 77 data sets for all chunks (continued)

Data set	TSF				WEASEL			
	10%	20%	30%	100%	10%	20%	30%	100%
FreezerRegularTrain	35	32	33	39	3	7	7	15
FreezerSmallTrain	25	25	26	30	1	3	4	4
Fungi	30	24	32	37	0	1	2	2
Ham	42	35	32	52	4	6	5	14
HandMovementDirection	384	486	502	593	81	203	129	412
Handwriting	134	140	138	153	8	18	30	60
Heartbeat	1343	1384	1409	1606	21	55	45	189
HouseTwenty	30	33	35	55	3	7	8	30
InsectEPGRegularTrain	40	42	42	52	5	5	7	25
InsectEPGSmallTrain	24	28	26	33	2	2	2	7
InsectWingbeatSound	50	49	51	59	4	8	16	21
ItalyPowerDemand	25	32	32	35	0	0	1	2
LargeKitchenAppliances	46	45	50	84	22	22	38	164
Libras	85	88	88	92	1	3	5	23
Lightning2	40	28	23	38	5	5	7	21
Lightning7	28	37	40	44	2	4	5	10
LSST	698	698	709	733	25	78	152	826
Meat	39	28	35	37	3	4	3	9
MoteStrain	34	23	23	25	0	0	1	3
MotorImagery	1874	2115	2391	4483	217	533	806	2208
NATOPS	1103	1108	1125	1178	16	39	67	316
NonInvasiveFetalECGThorax1	121	131	146	404	408	477	842	2395
NonInvasiveFetalECGThorax2	118	130	150	402	272	372	991	2539
OliveOil	34	38	41	46	2	1	2	4
PEMS-SF	9541	9759	9849	10541				
Phoneme	56	29	72	91	15	14	48	326

Table A.24: Train Time (CPU Time) for TSF and WEASEL on 77 data sets for all chunks (continued)

Data set	TSF				WEASEL			
	10%	20%	30%	100%	10%	20%	30%	100%
PigAirwayPressure	34	38	43	80	17	35	51	741
PigArtPressure	39	58	53	102	16	34	57	203
PigCVP	50	39	42	81	13	35	60	261
Plane	37	40	42	42	1	3	4	7
PowerCons	49	39	33	50	1	4	4	12
RacketSports	208	230	266	276	1	3	6	32
RefrigerationDevices	56	49	56	61	27	27	47	106
Rock	43	46	54	74	4	4	6	17
ScreenType	53	65	34	101	20	24	37	288
SelfRegulationSCP1	342	373	272	589	181	136	192	755
SelfRegulationSCP2	419	405	433	524	71	149	238	591
SemgHandGenderCh2	58	70	75	163	19	28	72	184
SemgHandMovementCh2	114	62	69	246	37	61	76	338
SemgHandSubjectCh2	61	73	78	256	28	45	94	529
SmallKitchenAppliances	74	46	49	84	16	29	32	393
SonyAIBORobotSurface1	32	31	35	35	0	0	0	2
SonyAIBORobotSurface2	33	23	24	25	0	0	1	3
StandWalkJump	170	187	190	274	8	18	25	73
StarlightCurves	78	86	99	332	21	31	43	119
Strawberry	75	52	53	65	9	19	38	53
Trace	31	38	41	33	2	4	5	13
TwoLeadECG	23	25	24	29	0	1	1	3
UWaveGestureLibrary	126	90	138	145	17	38	61	76
Wafer	68	67	69	77	6	23	30	47
Wine	26	26	27	34	1	2	2	4

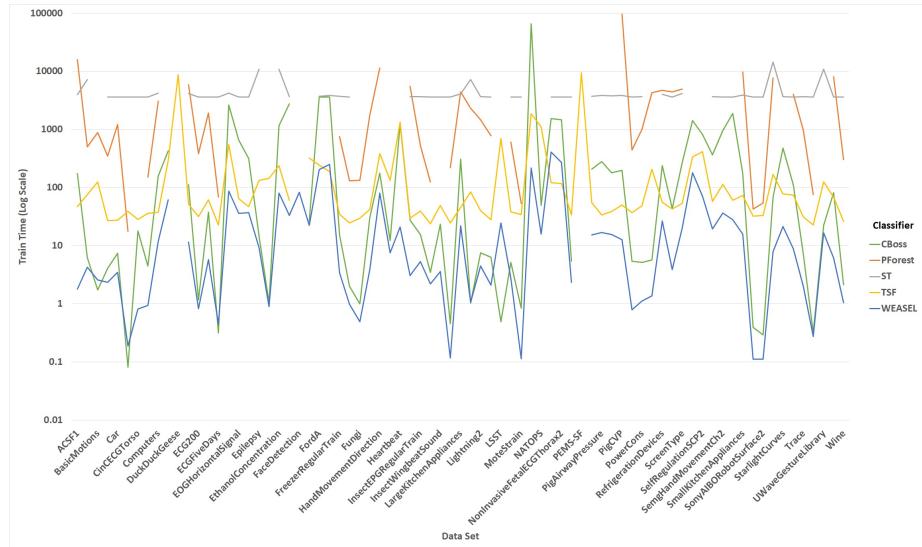


Figure A.6: Train Time (CPU Time in Log Scale) for all 5 classifiers for 10% chunk

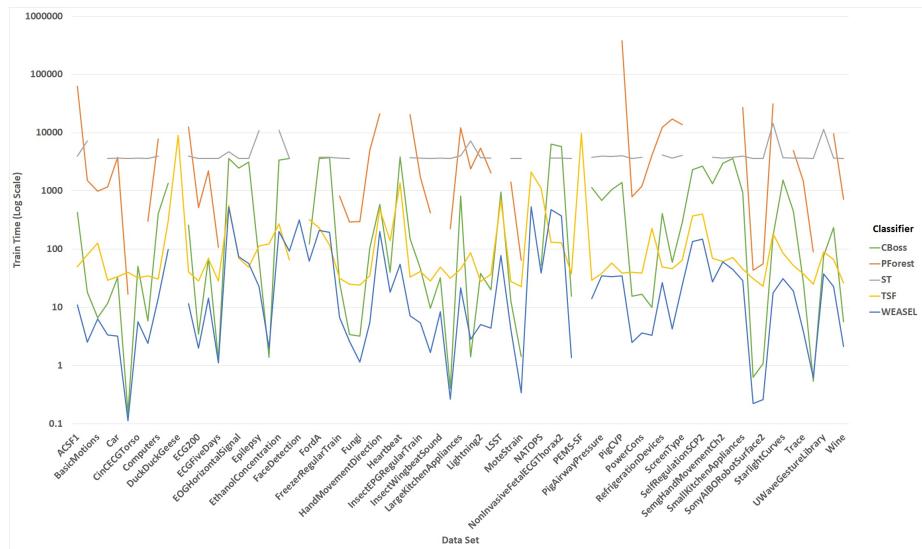


Figure A.7: Train Time (CPU Time in Log Scale) for all 5 classifiers for 20% chunk

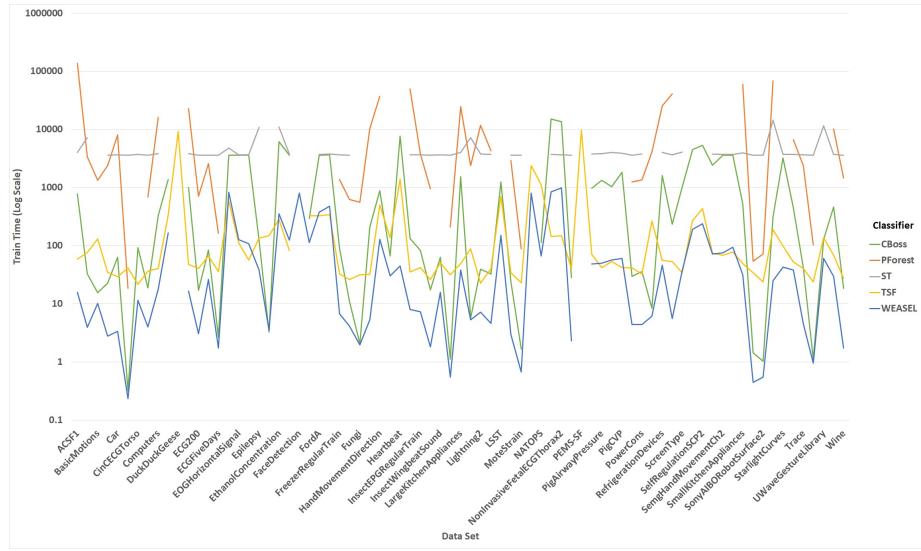


Figure A.8: Train Time (CPU Time in Log Scale) for all 5 classifiers for 30% chunk

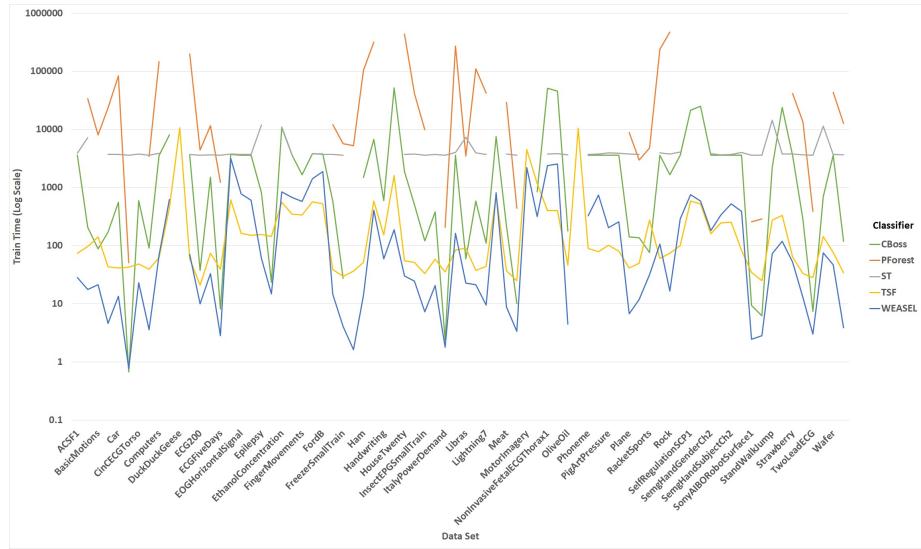


Figure A.9: Train Time (CPU Time in Log Scale) for all 5 classifiers for 100% chunk

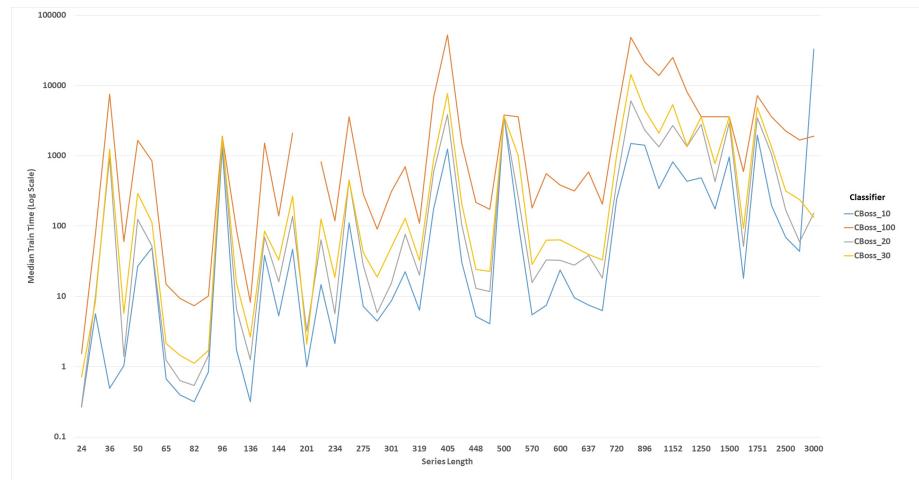


Figure A.10: Train Time (CPU Time in Log Scale) for CBOss per series length

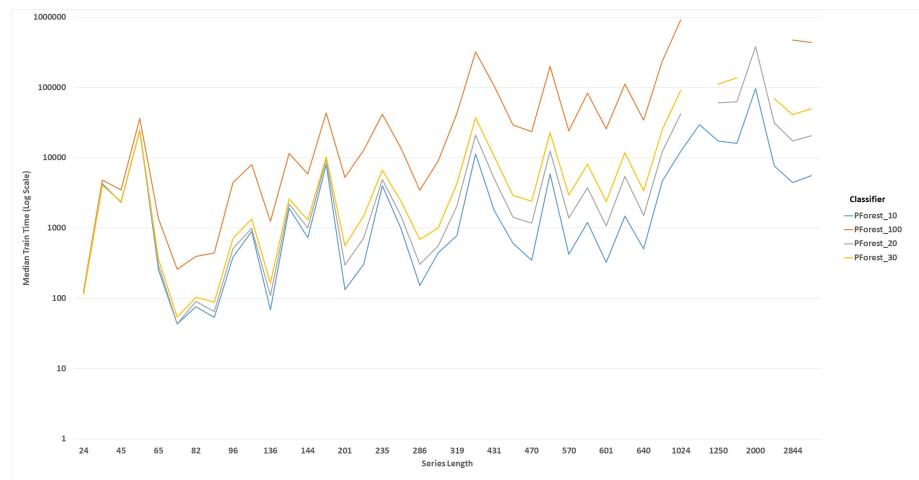


Figure A.11: Train Time (CPU Time in Log Scale) for PForest per series length

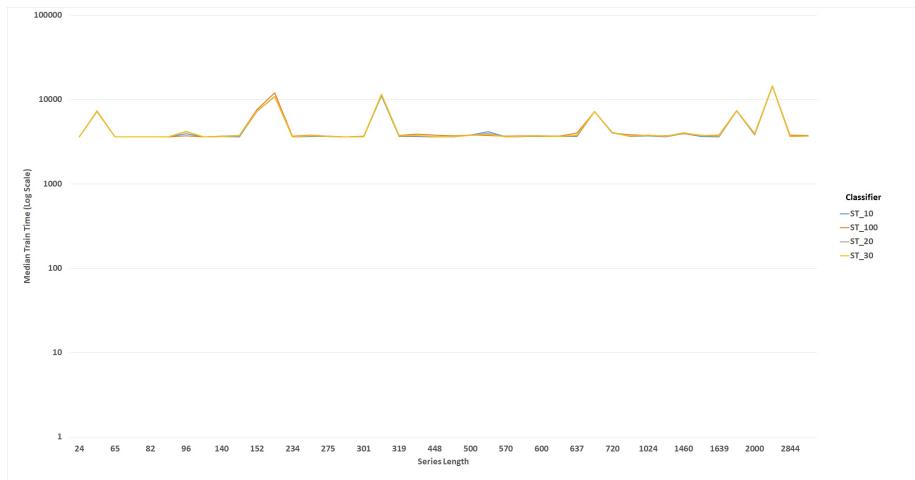


Figure A.12: Train Time (CPU Time in Log Scale) for ST per series length

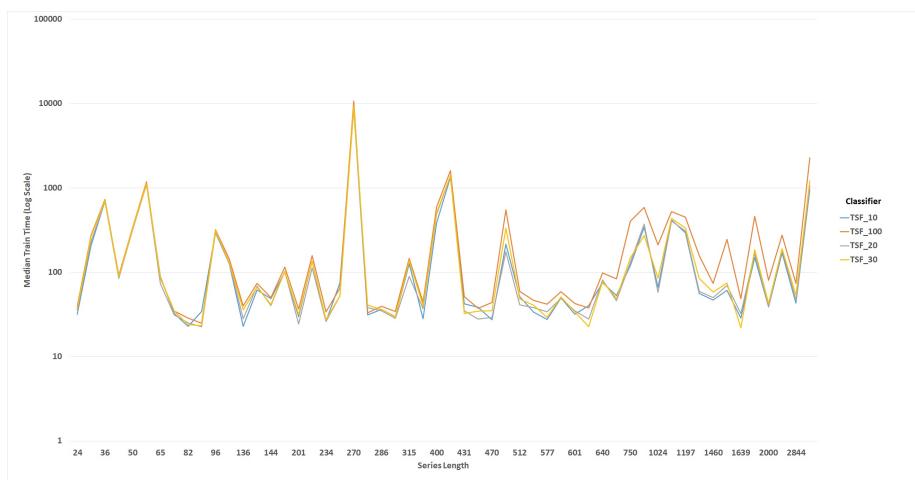


Figure A.13: Train Time (CPU Time in Log Scale) for TSF per series length

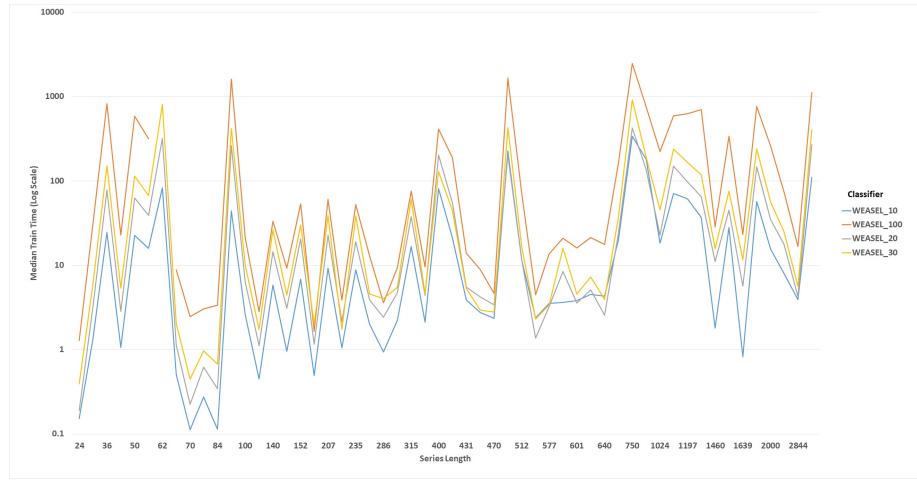


Figure A.14: Train Time (CPU Time in Log Scale) for WEASEL per series length

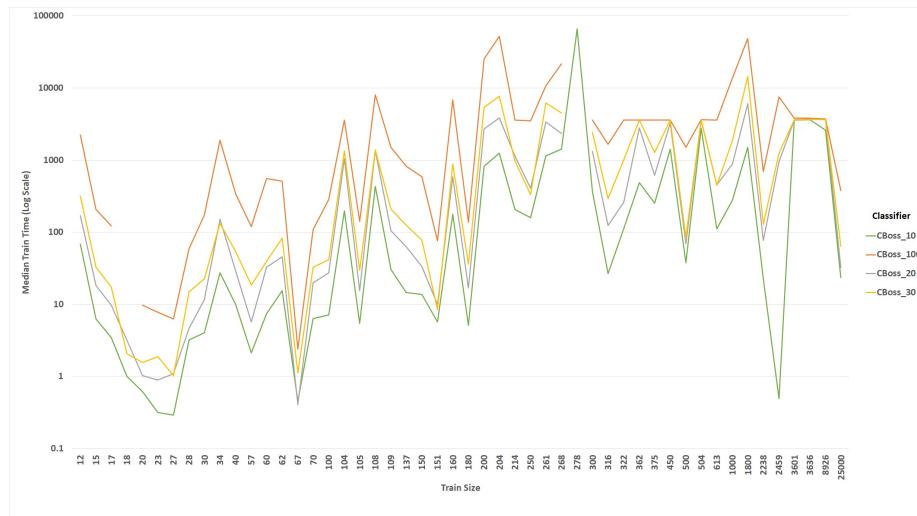


Figure A.15: Train Time (CPU Time in Log Scale) for CBOss per train size

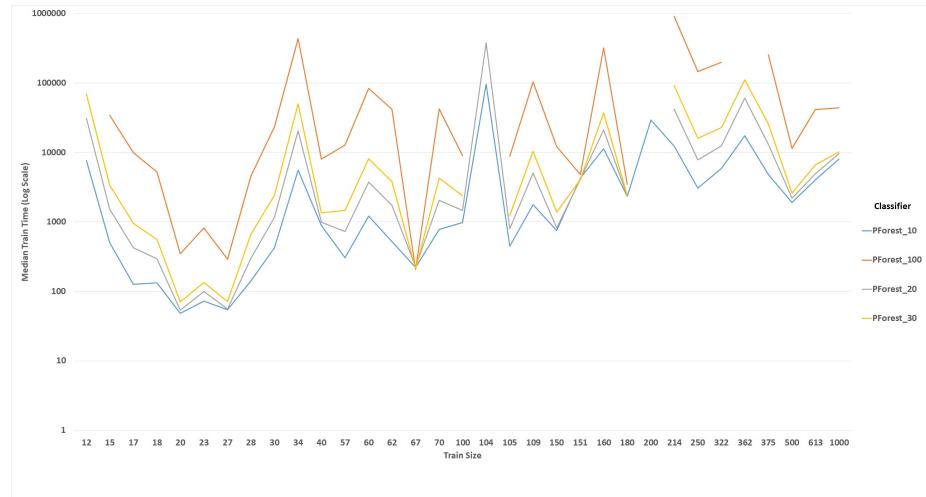


Figure A.16: Train Time (CPU Time in Log Scale) for PForest per train size

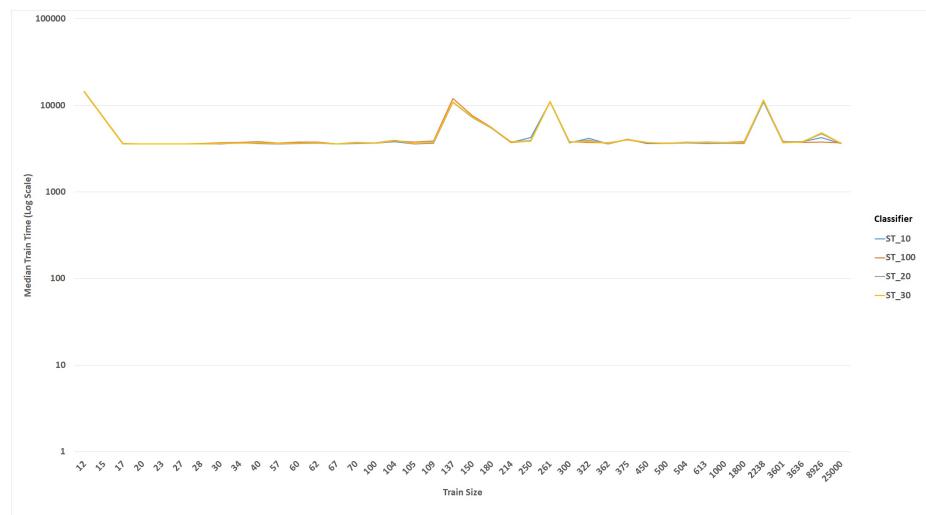


Figure A.17: Train Time (CPU Time in Log Scale) for ST per train size

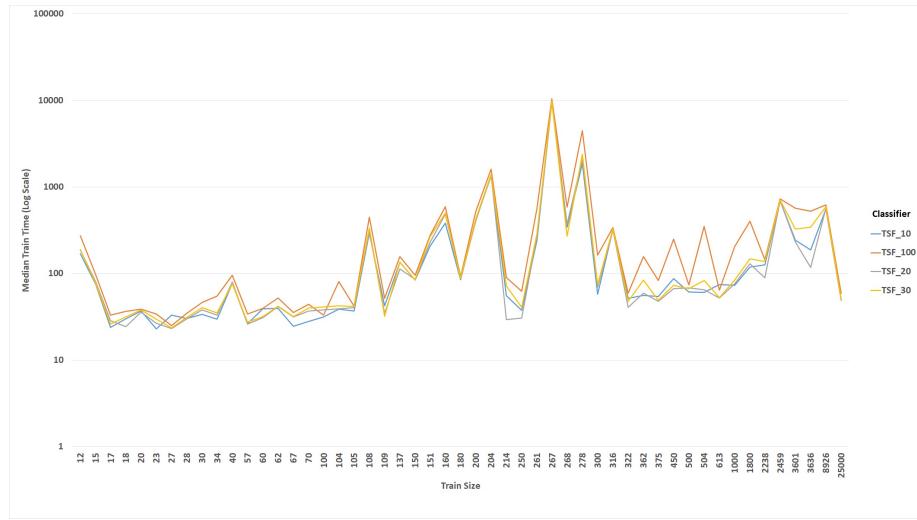


Figure A.18: Train Time (CPU Time in Log Scale) for TSF per train size

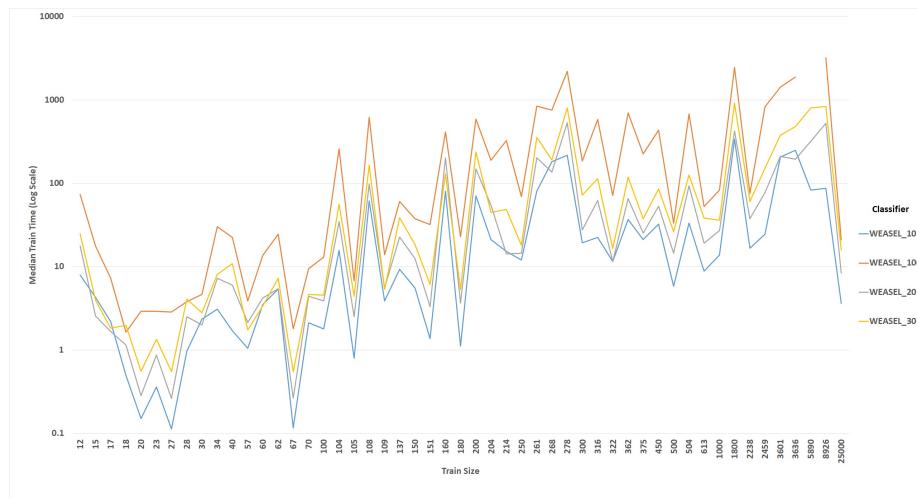


Figure A.19: Train Time (CPU Time in Log Scale) for WEASEL per train size

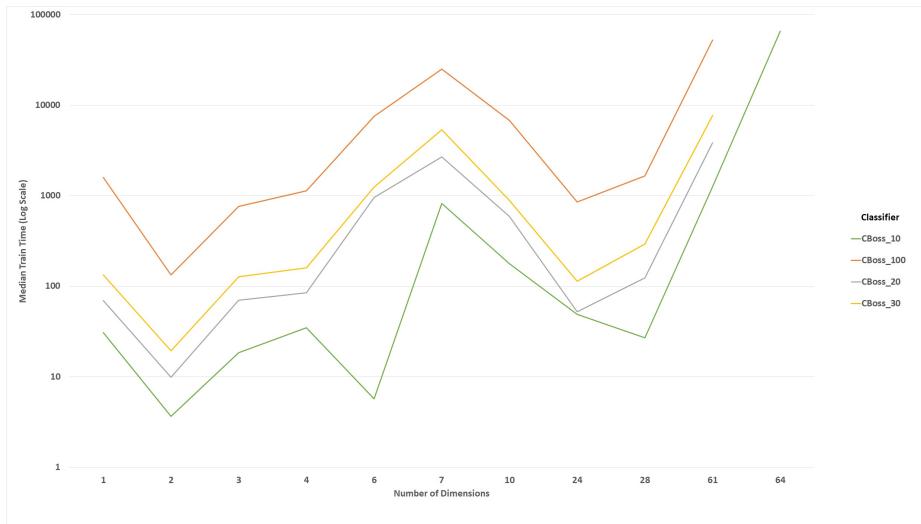


Figure A.20: Train Time (CPU Time in Log Scale) for CBOss per number of dimensions

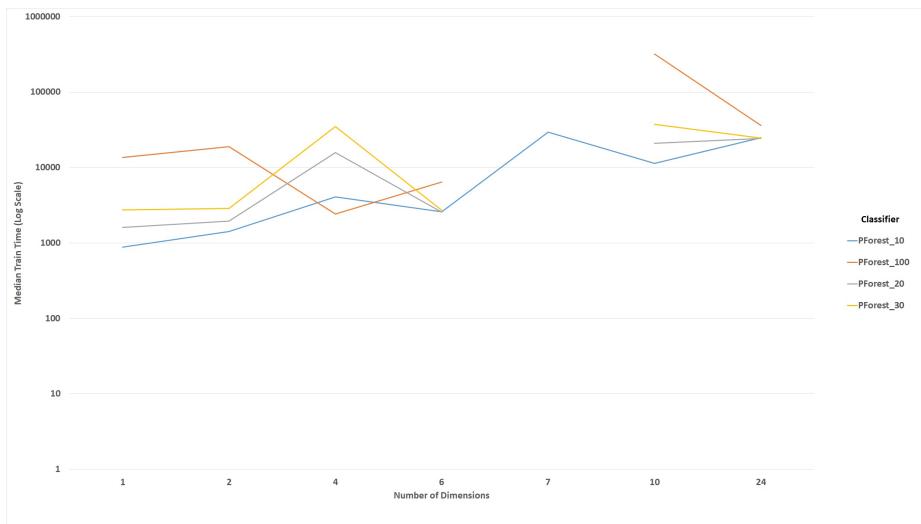


Figure A.21: Train Time (CPU Time in Log Scale) for PForest per number of dimensions

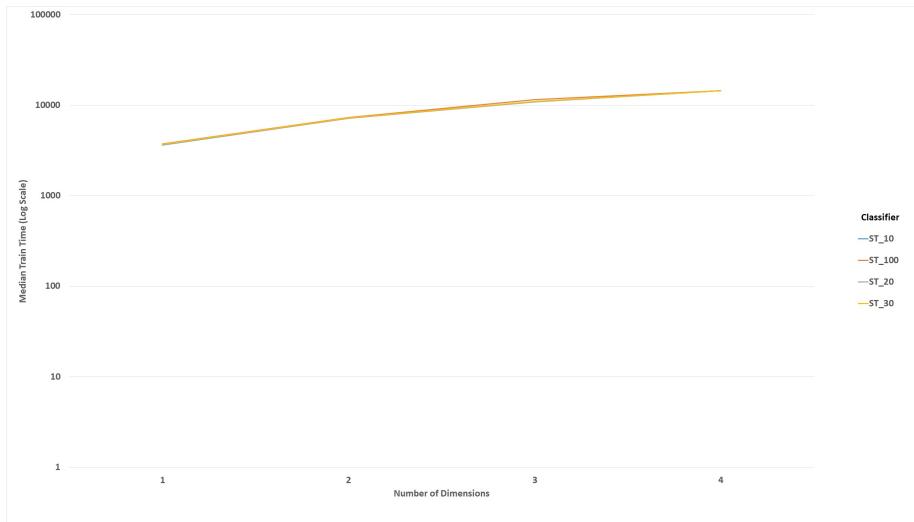


Figure A.22: Train Time (CPU Time in Log Scale) for ST per number of dimensions

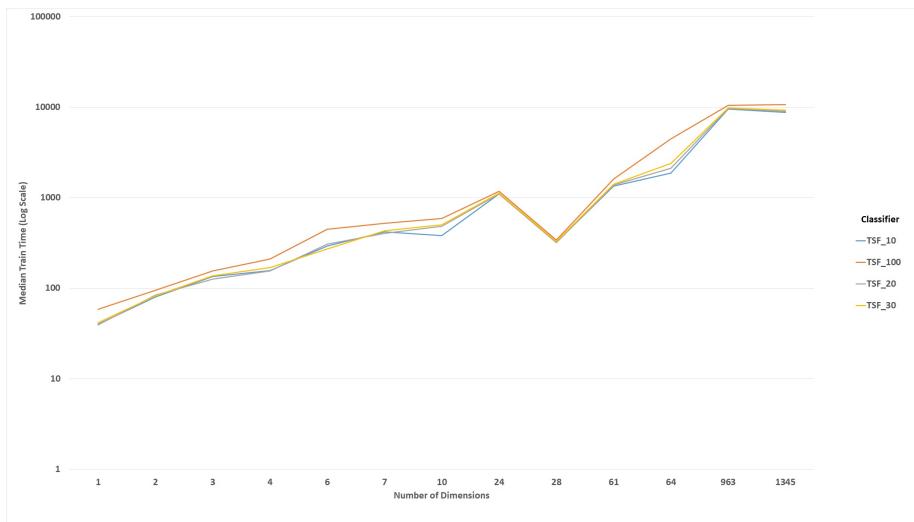


Figure A.23: Train Time (CPU Time in Log Scale) for TSF per number of dimensions

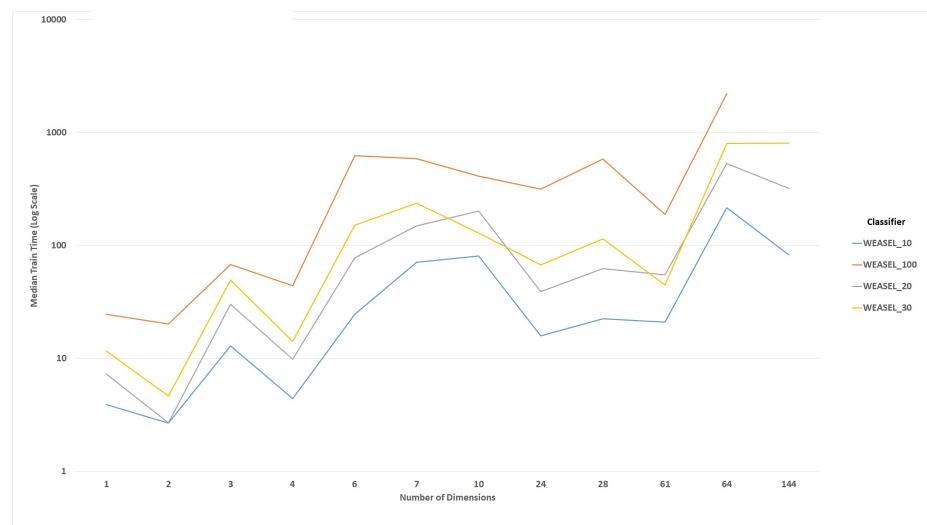


Figure A.24: Train Time (CPU Time in Log Scale) for WEASEL per number of dimensions

# Appendix B

# Recommender Results

## B.1 Are the chunk learners able to learn F-scores for the classifiers ?

#Seed	Revealed%	MAE	RMSE	$R^2$
0	10	0.059	0.087	0.046
0	20	0.098	0.129	0.340
0	30	0.109	0.149	0.462
0	100	0.125	0.169	0.664
1	10	0.057	0.079	0.108
1	20	0.122	0.150	0.079
1	30	0.134	0.165	0.254
1	100	0.168	0.207	0.371
2	10	0.035	0.049	0.314
2	20	0.086	0.111	0.385
2	30	0.110	0.143	0.321
2	100	0.150	0.189	0.544
3	10	0.070	0.102	0.155
3	20	0.110	0.149	0.268
3	30	0.123	0.177	0.321
3	100	0.176	0.217	0.511
4	10	0.049	0.068	0.010
4	20	0.104	0.145	0.110
4	30	0.137	0.186	0.065
4	100	0.164	0.230	0.366
5	10	0.048	0.059	0.529
5	20	0.093	0.113	0.443
5	30	0.118	0.139	0.462
5	100	0.167	0.189	0.546
6	10	0.049	0.063	0.396
6	20	0.098	0.119	0.309
6	30	0.131	0.173	0.085

#Seed	Revealed%	MAE	RMSE	$R^2$
6	100	0.136	0.169	0.556
7	10	0.058	0.083	0.011
7	20	0.135	0.174	-0.050
7	30	0.165	0.209	0.044
7	100	0.197	0.238	0.382
8	10	0.043	0.054	0.330
8	20	0.099	0.116	0.339
8	30	0.122	0.146	0.366
8	100	0.157	0.202	0.401
9	10	0.045	0.056	0.516
9	20	0.082	0.101	0.464
9	30	0.097	0.118	0.529
9	100	0.138	0.180	0.515
10	10	0.050	0.072	0.449
10	20	0.097	0.120	0.487
10	30	0.112	0.139	0.515
10	100	0.160	0.198	0.506
11	10	0.049	0.075	0.006
11	20	0.089	0.125	0.059
11	30	0.123	0.162	0.138
11	100	0.168	0.193	0.425
12	10	0.048	0.068	0.221
12	20	0.088	0.112	0.349
12	30	0.109	0.140	0.409
12	100	0.137	0.168	0.641
13	10	0.051	0.071	0.465
13	20	0.098	0.125	0.428
13	30	0.142	0.174	0.283
13	100	0.157	0.194	0.540
14	10	0.048	0.064	0.406
14	20	0.093	0.120	0.314
14	30	0.115	0.139	0.407
14	100	0.190	0.225	0.395
15	10	0.045	0.055	0.190
15	20	0.111	0.130	0.283
15	30	0.121	0.153	0.408
15	100	0.225	0.279	0.180
16	10	0.045	0.071	0.338
16	20	0.091	0.129	0.255
16	30	0.114	0.162	0.289
16	100	0.148	0.196	0.561
17	10	0.050	0.075	0.127
17	20	0.096	0.136	0.263
17	30	0.135	0.183	0.149
17	100	0.135	0.168	0.686
18	10	0.042	0.059	-0.215
18	20	0.092	0.117	-0.122

#Seed	Revealed%	MAE	RMSE	$R^2$
18	30	0.122	0.167	-0.345
18	100	0.166	0.214	0.230
19	10	0.040	0.054	0.588
19	20	0.076	0.097	0.577
19	30	0.096	0.122	0.594
19	100	0.136	0.181	0.624
20	10	0.058	0.088	0.265
20	20	0.111	0.148	0.246
20	30	0.139	0.180	0.291
20	100	0.175	0.213	0.524
21	10	0.045	0.066	0.302
21	20	0.088	0.110	0.509
21	30	0.121	0.146	0.460
21	100	0.172	0.249	0.360
22	10	0.038	0.045	0.638
22	20	0.084	0.103	0.352
22	30	0.098	0.124	0.350
22	100	0.165	0.208	0.221
23	10	0.053	0.062	0.480
23	20	0.093	0.117	0.488
23	30	0.119	0.151	0.465
23	100	0.138	0.172	0.678
24	10	0.037	0.067	0.040
24	20	0.081	0.112	0.349
24	30	0.102	0.143	0.366
24	100	0.148	0.192	0.549
25	10	0.046	0.060	0.275
25	20	0.100	0.129	0.314
25	30	0.141	0.171	0.175
25	100	0.159	0.195	0.504
26	10	0.049	0.072	0.395
26	20	0.092	0.118	0.527
26	30	0.119	0.146	0.502
26	100	0.180	0.211	0.535
27	10	0.046	0.064	0.026
27	20	0.089	0.117	0.165
27	30	0.113	0.147	0.243
27	100	0.127	0.155	0.660
28	10	0.040	0.064	0.363
28	20	0.086	0.120	0.436
28	30	0.105	0.140	0.483
28	100	0.134	0.169	0.647
29	10	0.040	0.047	0.187
29	20	0.091	0.110	0.278
29	30	0.108	0.142	0.377
29	100	0.144	0.188	0.605
30	10	0.047	0.062	0.489

#Seed	Revealed%	MAE	RMSE	$R^2$
30	20	0.095	0.116	0.488
30	30	0.128	0.153	0.408
30	100	0.147	0.190	0.596
31	10	0.050	0.062	0.176
31	20	0.095	0.118	0.327
31	30	0.127	0.154	0.369
31	100	0.189	0.227	0.462
32	10	0.049	0.063	0.434
32	20	0.097	0.120	0.410
32	30	0.130	0.158	0.360
32	100	0.159	0.204	0.480
33	10	0.049	0.063	-0.369
33	20	0.094	0.114	-0.014
33	30	0.129	0.155	-0.004
33	100	0.177	0.209	0.312
34	10	0.066	0.084	0.012
34	20	0.104	0.143	0.146
34	30	0.125	0.157	0.315
34	100	0.172	0.210	0.540
35	10	0.047	0.074	0.073
35	20	0.098	0.137	0.111
35	30	0.116	0.168	0.139
35	100	0.135	0.164	0.667
36	10	0.047	0.059	0.439
36	20	0.110	0.146	-0.019
36	30	0.115	0.142	0.366
36	100	0.151	0.190	0.461
37	10	0.039	0.055	0.028
37	20	0.088	0.115	0.054
37	30	0.132	0.175	-0.173
37	100	0.121	0.155	0.663
38	10	0.048	0.061	0.134
38	20	0.077	0.096	0.392
38	30	0.100	0.125	0.449
38	100	0.145	0.179	0.625
39	10	0.056	0.080	0.336
39	20	0.106	0.136	0.380
39	30	0.133	0.161	0.436
39	100	0.165	0.214	0.532
40	10	0.045	0.060	0.464
40	20	0.081	0.099	0.523
40	30	0.103	0.128	0.495
40	100	0.111	0.137	0.721
41	10	0.041	0.056	0.084
41	20	0.092	0.126	0.164
41	30	0.132	0.162	0.250
41	100	0.165	0.198	0.535

#Seed	Revealed%	MAE	RMSE	$R^2$
42	10	0.054	0.075	0.283
42	20	0.102	0.122	0.343
42	30	0.124	0.148	0.329
42	100	0.139	0.175	0.500
43	10	0.039	0.048	0.295
43	20	0.081	0.111	0.252
43	30	0.095	0.135	0.357
43	100	0.151	0.177	0.614
44	10	0.044	0.067	-0.021
44	20	0.103	0.144	-0.186
44	30	0.126	0.163	0.086
44	100	0.125	0.161	0.705
45	10	0.049	0.076	0.150
45	20	0.108	0.138	0.295
45	30	0.124	0.169	0.360
45	100	0.166	0.206	0.573
46	10	0.074	0.099	-0.935
46	20	0.125	0.165	-0.558
46	30	0.142	0.181	0.025
46	100	0.184	0.246	0.308
47	10	0.046	0.063	0.212
47	20	0.083	0.107	0.434
47	30	0.109	0.140	0.437
47	100	0.159	0.210	0.486
48	10	0.038	0.065	0.019
48	20	0.068	0.106	0.285
48	30	0.100	0.140	0.281
48	100	0.155	0.188	0.449
49	10	0.057	0.087	0.211
49	20	0.107	0.156	0.140
49	30	0.130	0.180	0.282
49	100	0.134	0.161	0.670

Table B.1: Performance of chunk learners in predicting  $F_\beta$  per chunk over 50 runs

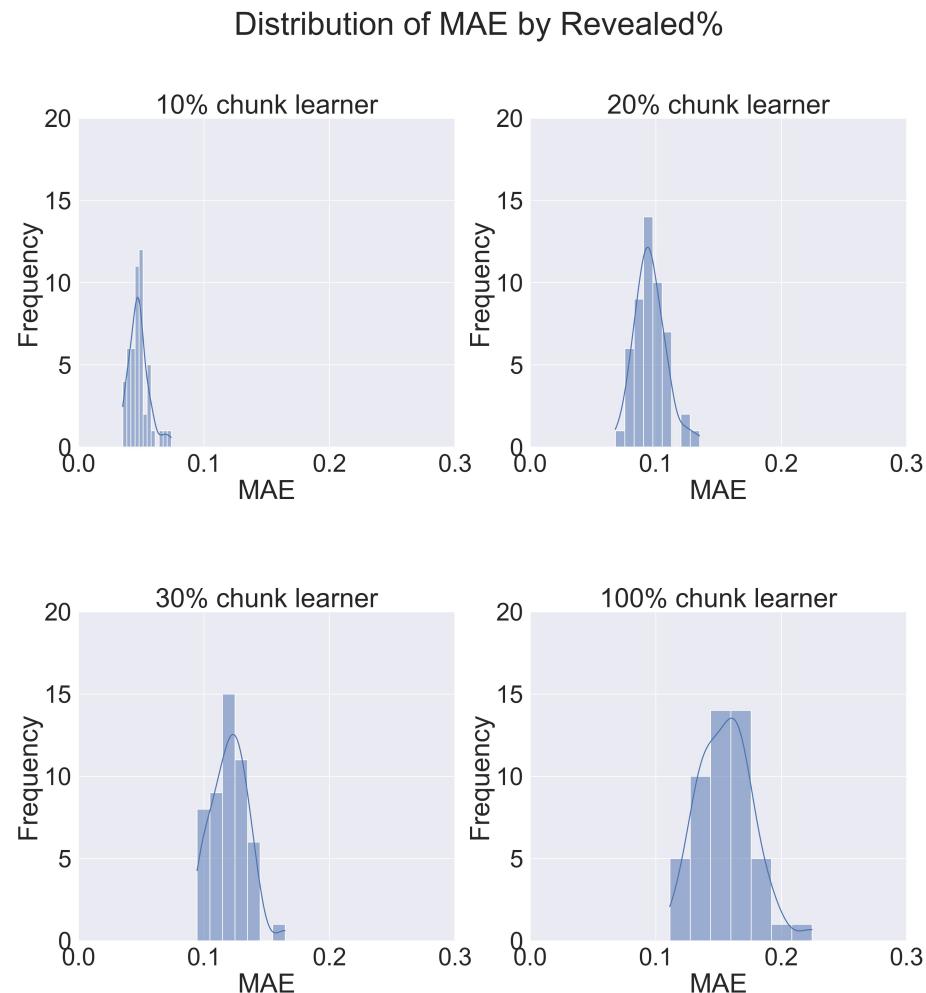


Figure B.1: Histograms for distribution of MAE per chunk over 50 runs

## B.2 Which features contribute the most to the F-Score Prediction ?

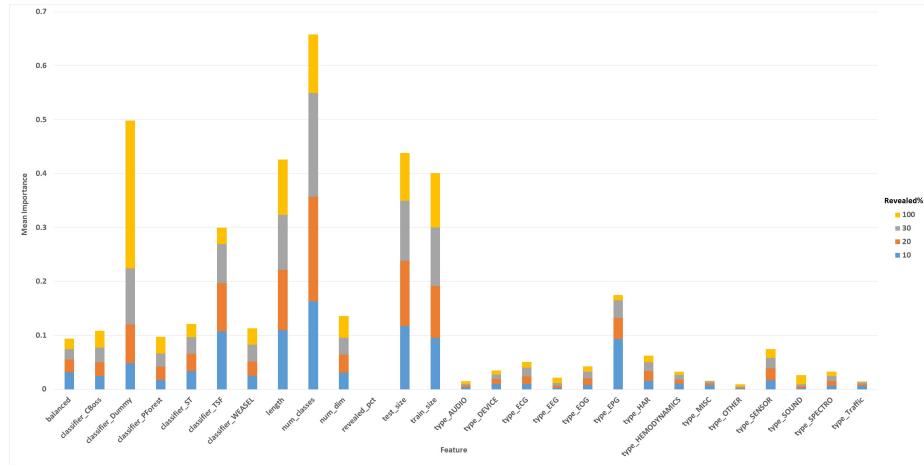


Figure B.2: Mean Feature Importance stacked by chunk

### B.3 Are the final recommendations of the framework consistent with actual results ?

#Seed	Revealed%	#Instances	TP	TN	FP	FN	Accuracy	Precision	Recall
0	10	81	47	15	19	0	0.77	0.71	1.00
0	20	81	56	16	9	0	0.89	0.86	1.00
0	30	80	60	15	5	0	0.94	0.92	1.00
0	100	80	65	15	0	0	1.00	1.00	1.00
1	10	85	44	16	24	1	0.71	0.65	0.98
1	20	85	54	15	16	0	0.81	0.77	1.00
1	30	84	54	15	15	0	0.82	0.78	1.00
1	100	83	59	15	9	0	0.89	0.87	1.00
2	10	81	46	16	18	1	0.77	0.72	0.98
2	20	81	56	15	10	0	0.88	0.85	1.00
2	30	81	57	15	9	0	0.89	0.86	1.00
2	100	81	65	15	1	0	0.99	0.98	1.00
3	10	83	48	15	18	2	0.76	0.73	0.96
3	20	83	52	15	16	0	0.81	0.76	1.00
3	30	82	54	16	9	3	0.85	0.86	0.95
3	100	82	63	15	4	0	0.95	0.94	1.00
4	10	83	50	16	17	0	0.80	0.75	1.00
4	20	83	52	15	16	0	0.81	0.76	1.00
4	30	83	58	15	10	0	0.88	0.85	1.00
4	100	82	66	15	1	0	0.99	0.99	1.00
5	10	80	43	17	20	0	0.75	0.68	1.00
5	20	80	52	16	12	0	0.85	0.81	1.00
5	30	79	57	15	7	0	0.91	0.89	1.00
5	100	78	62	15	1	0	0.99	0.98	1.00

#Seed	Revealed%	#Instances	TP	TN	FP	FN	Accuracy	Precision	Recall
6	10	78	39	15	24	0	0.69	0.62	1.00
6	20	78	52	15	11	0	0.86	0.83	1.00
6	30	78	53	15	10	0	0.87	0.84	1.00
6	100	77	57	15	5	0	0.94	0.92	1.00
7	10	84	41	20	16	7	0.73	0.72	0.85
7	20	84	62	15	7	0	0.92	0.90	1.00
7	30	84	61	15	8	0	0.90	0.88	1.00
7	100	83	67	15	1	0	0.99	0.99	1.00
8	10	79	34	17	27	1	0.65	0.56	0.97
8	20	79	44	15	20	0	0.75	0.69	1.00
8	30	79	48	16	15	0	0.81	0.76	1.00
8	100	79	55	15	9	0	0.89	0.86	1.00
9	10	80	39	15	26	0	0.68	0.60	1.00
9	20	79	48	15	16	0	0.80	0.75	1.00
9	30	78	53	15	10	0	0.87	0.84	1.00
9	100	77	56	15	6	0	0.92	0.90	1.00
10	10	84	53	16	14	1	0.82	0.79	0.98
10	20	84	56	15	13	0	0.85	0.81	1.00
10	30	83	61	15	7	0	0.92	0.90	1.00
10	100	82	62	15	5	0	0.94	0.93	1.00
11	10	81	42	21	18	0	0.78	0.70	1.00
11	20	81	52	15	14	0	0.83	0.79	1.00
11	30	80	56	15	9	0	0.89	0.86	1.00
11	100	80	65	15	0	0	1.00	1.00	1.00
12	10	82	46	16	20	0	0.76	0.70	1.00
12	20	81	45	15	21	0	0.74	0.68	1.00
12	30	81	52	15	14	0	0.83	0.79	1.00
12	100	80	60	15	5	0	0.94	0.92	1.00
13	10	80	42	20	16	2	0.78	0.72	0.95
13	20	80	34	19	9	18	0.66	0.79	0.65
13	30	79	52	15	12	0	0.85	0.81	1.00
13	100	79	60	15	4	0	0.95	0.94	1.00
14	10	78	45	17	13	3	0.79	0.78	0.94
14	20	78	51	15	12	0	0.85	0.81	1.00
14	30	77	53	15	9	0	0.88	0.85	1.00
14	100	77	60	15	2	0	0.97	0.97	1.00
15	10	81	39	15	23	4	0.67	0.63	0.91
15	20	81	47	15	19	0	0.77	0.71	1.00
15	30	81	53	15	13	0	0.84	0.80	1.00
15	100	80	58	15	7	0	0.91	0.89	1.00
16	10	83	46	20	15	2	0.80	0.75	0.96
16	20	83	48	15	19	1	0.76	0.72	0.98
16	30	82	57	16	9	0	0.89	0.86	1.00
16	100	79	61	15	3	0	0.96	0.95	1.00
17	10	73	37	15	21	0	0.71	0.64	1.00
17	20	73	46	15	11	1	0.84	0.81	0.98
17	30	72	44	18	4	6	0.86	0.92	0.88

#Seed	Revealed%	#Instances	TP	TN	FP	FN	Accuracy	Precision	Recall
17	100	69	50	16	3	0	0.96	0.94	1.00
18	10	78	43	18	11	6	0.78	0.80	0.88
18	20	78	51	15	12	0	0.85	0.81	1.00
18	30	78	55	15	7	1	0.90	0.89	0.98
18	100	77	62	15	0	0	1.00	1.00	1.00
19	10	81	40	16	25	0	0.69	0.62	1.00
19	20	80	45	18	15	2	0.79	0.75	0.96
19	30	80	52	17	11	0	0.86	0.83	1.00
19	100	79	55	15	9	0	0.89	0.86	1.00
20	10	78	45	16	15	2	0.78	0.75	0.96
20	20	78	50	15	13	0	0.83	0.79	1.00
20	30	78	52	15	11	0	0.86	0.83	1.00
20	100	78	63	15	0	0	1.00	1.00	1.00
21	10	84	47	20	15	2	0.80	0.76	0.96
21	20	84	50	15	18	1	0.77	0.74	0.98
21	30	83	59	15	9	0	0.89	0.87	1.00
21	100	81	62	15	4	0	0.95	0.94	1.00
22	10	76	38	16	22	0	0.71	0.63	1.00
22	20	75	44	15	16	0	0.79	0.73	1.00
22	30	74	45	15	14	0	0.81	0.76	1.00
22	100	74	53	15	6	0	0.92	0.90	1.00
23	10	79	38	16	25	0	0.68	0.60	1.00
23	20	79	46	15	18	0	0.77	0.72	1.00
23	30	79	52	15	12	0	0.85	0.81	1.00
23	100	77	61	15	1	0	0.99	0.98	1.00
24	10	78	39	22	15	2	0.78	0.72	0.95
24	20	78	47	15	16	0	0.79	0.75	1.00
24	30	78	48	15	15	0	0.81	0.76	1.00
24	100	78	59	15	4	0	0.95	0.94	1.00
25	10	87	41	19	25	2	0.69	0.62	0.95
25	20	87	52	16	15	4	0.78	0.78	0.93
25	30	87	51	17	9	10	0.78	0.85	0.84
25	100	87	67	15	5	0	0.94	0.93	1.00
26	10	83	53	18	12	0	0.86	0.82	1.00
26	20	83	60	15	8	0	0.90	0.88	1.00
26	30	83	62	15	6	0	0.93	0.91	1.00
26	100	83	68	15	0	0	1.00	1.00	1.00
27	10	81	46	19	16	0	0.80	0.74	1.00
27	20	81	52	15	14	0	0.83	0.79	1.00
27	30	81	55	15	11	0	0.86	0.83	1.00
27	100	81	66	15	0	0	1.00	1.00	1.00
28	10	79	46	15	17	1	0.77	0.73	0.98
28	20	79	50	15	13	1	0.82	0.79	0.98
28	30	77	53	15	9	0	0.88	0.85	1.00
28	100	74	56	15	3	0	0.96	0.95	1.00
29	10	83	43	16	24	0	0.71	0.64	1.00
29	20	82	50	15	17	0	0.79	0.75	1.00

#Seed	Revealed%	#Instances	TP	TN	FP	FN	Accuracy	Precision	Recall
29	30	81	55	15	11	0	0.86	0.83	1.00
29	100	81	64	15	2	0	0.98	0.97	1.00
30	10	82	46	16	19	1	0.76	0.71	0.98
30	20	82	55	15	12	0	0.85	0.82	1.00
30	30	81	54	15	12	0	0.85	0.82	1.00
30	100	79	59	15	5	0	0.94	0.92	1.00
31	10	82	42	17	23	0	0.72	0.65	1.00
31	20	82	46	15	21	0	0.74	0.69	1.00
31	30	82	49	15	18	0	0.78	0.73	1.00
31	100	81	57	15	9	0	0.89	0.86	1.00
32	10	83	17	29	1	36	0.55	0.94	0.32
32	20	82	47	17	9	9	0.78	0.84	0.84
32	30	81	55	15	11	0	0.86	0.83	1.00
32	100	80	64	15	1	0	0.99	0.98	1.00
33	10	81	40	18	21	2	0.72	0.66	0.95
33	20	81	48	15	18	0	0.78	0.73	1.00
33	30	81	52	15	14	0	0.83	0.79	1.00
33	100	81	66	15	0	0	1.00	1.00	1.00
34	10	77	39	17	19	2	0.73	0.67	0.95
34	20	77	44	16	11	6	0.78	0.80	0.88
34	30	77	55	15	7	0	0.91	0.89	1.00
34	100	77	61	15	1	0	0.99	0.98	1.00
35	10	82	45	18	18	1	0.77	0.71	0.98
35	20	82	57	15	10	0	0.88	0.85	1.00
35	30	82	59	15	8	0	0.90	0.88	1.00
35	100	82	67	15	0	0	1.00	1.00	1.00
36	10	75	34	18	21	2	0.69	0.62	0.94
36	20	75	31	23	5	16	0.72	0.86	0.66
36	30	75	48	15	12	0	0.84	0.80	1.00
36	100	73	53	15	5	0	0.93	0.91	1.00
37	10	84	45	16	22	1	0.73	0.67	0.98
37	20	84	50	15	19	0	0.77	0.72	1.00
37	30	84	53	17	13	1	0.83	0.80	0.98
37	100	84	69	15	0	0	1.00	1.00	1.00
38	10	84	40	17	27	0	0.68	0.60	1.00
38	20	83	50	15	18	0	0.78	0.74	1.00
38	30	81	54	15	12	0	0.85	0.82	1.00
38	100	81	60	15	6	0	0.93	0.91	1.00
39	10	77	40	17	18	2	0.74	0.69	0.95
39	20	77	46	15	15	1	0.79	0.75	0.98
39	30	77	51	15	11	0	0.86	0.82	1.00
39	100	75	55	15	5	0	0.93	0.92	1.00
40	10	81	42	19	20	0	0.75	0.68	1.00
40	20	80	52	15	13	0	0.84	0.80	1.00
40	30	80	56	15	9	0	0.89	0.86	1.00
40	100	80	63	15	2	0	0.98	0.97	1.00
41	10	76	36	16	22	2	0.68	0.62	0.95

#Seed	Revealed%	#Instances	TP	TN	FP	FN	Accuracy	Precision	Recall
41	20	76	38	18	16	4	0.74	0.70	0.90
41	30	76	46	15	15	0	0.80	0.75	1.00
41	100	76	56	15	5	0	0.93	0.92	1.00
42	10	86	41	16	29	0	0.66	0.59	1.00
42	20	86	51	15	20	0	0.77	0.72	1.00
42	30	86	54	15	17	0	0.80	0.76	1.00
42	100	85	65	15	5	0	0.94	0.93	1.00
43	10	85	51	17	17	0	0.80	0.75	1.00
43	20	85	54	18	8	5	0.85	0.87	0.92
43	30	84	62	15	7	0	0.92	0.90	1.00
43	100	83	67	15	1	0	0.99	0.99	1.00
44	10	77	39	16	22	0	0.71	0.64	1.00
44	20	77	17	25	10	25	0.55	0.63	0.40
44	30	77	47	15	15	0	0.81	0.76	1.00
44	100	73	55	15	3	0	0.96	0.95	1.00
45	10	78	39	20	10	9	0.76	0.80	0.81
45	20	78	52	15	11	0	0.86	0.83	1.00
45	30	77	56	15	6	0	0.92	0.90	1.00
45	100	77	58	15	4	0	0.95	0.94	1.00
46	10	83	41	19	19	4	0.72	0.68	0.91
46	20	83	51	15	15	2	0.80	0.77	0.96
46	30	82	54	15	13	0	0.84	0.81	1.00
46	100	81	59	15	7	0	0.91	0.89	1.00
47	10	85	46	15	24	0	0.72	0.66	1.00
47	20	85	54	15	15	1	0.81	0.78	0.98
47	30	85	58	15	12	0	0.86	0.83	1.00
47	100	85	65	15	5	0	0.94	0.93	1.00
48	10	84	41	22	15	6	0.75	0.73	0.87
48	20	84	55	15	14	0	0.83	0.80	1.00
48	30	84	58	15	11	0	0.87	0.84	1.00
48	100	82	66	15	1	0	0.99	0.99	1.00
49	10	82	33	27	16	6	0.73	0.67	0.85
49	20	81	51	15	15	0	0.81	0.77	1.00
49	30	81	54	15	12	0	0.85	0.82	1.00
49	100	80	64	15	1	0	0.99	0.98	1.00

Table B.2: Confusion Matrix of chunk learners in predicting good and bad classifiers per chunk over 50 runs

# Bibliography

- [Alo+18] Md Zahangir Alom et al. “The history began from alexnet: A comprehensive survey on deep learning approaches”. In: *arXiv preprint arXiv:1803.01164* (2018).
- [AML19] 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.
- [Bag+12] 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.
- [Bag+15] 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.
- [Bag+17] 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.
- [Bag+18] Anthony Bagnall et al. “The UEA multivariate time series classification archive, 2018”. In: *arXiv preprint arXiv:1811.00075* (2018).
- [Bag+20] Anthony Bagnall et al. “A tale of two toolkits, report the third: on the usage and performance of HIVE-COTE v1. 0”. In: *arXiv e-prints* (2020), arXiv–2004.
- [BB17a] Aaron Bostrom and Anthony Bagnall. “A shapelet transform for multivariate time series classification”. In: *arXiv preprint arXiv:1712.06428* (2017).
- [BB17b] 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](https://doi.org/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).
- [BCM16] Alessio Benavoli, Giorgio Corani, and Francesca Mangili. “Should we really use post-hoc tests based on mean-ranks?” In: *The Journal of Machine Learning Research* 17.1 (2016), pp. 152–161.
  - [Ber+06] Laurent Bernaille et al. “Traffic classification on the fly”. In: *ACM SIGCOMM Computer Communication Review* 36.2 (2006), pp. 23–26.
  - [Bos18] Aaron Bostrom. “Shapelet Transforms for Univariate and Multivariate Time Series Classification”. PhD thesis. University of East Anglia, 2018.
  - [BR16] 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.
  - [Bre01] Leo Breiman. “Random forests”. In: *Machine learning* 45.1 (2001), pp. 5–32.
  - [BRT13] 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.
  - [CCC16] Zhicheng Cui, Wenlin Chen, and Yixin Chen. “Multi-scale convolutional neural networks for time series classification”. In: *arXiv preprint arXiv:1603.06995* (2016).
  - [CCP09] Jorge Caiado, Nuno Crato, and Daniel Peña. “Comparison of times series with unequal length in the frequency domain”. In: *Communications in Statistics—Simulation and Computation®* 38.3 (2009), pp. 527–540.
  - [Che+18] Weihao Cheng et al. “Predicting Complex Activities from Ongoing Multivariate Time Series.” In: *IJCAI*. 2018, pp. 3322–3328.
  - [CN04] 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.
  - [CPB18] Raphael Couronné, Philipp Probst, and Anne-Laure Boulesteix. “Random forest versus logistic regression: a large-scale benchmark experiment”. In: *BMC bioinformatics* 19.1 (2018), pp. 1–14.
  - [Dau+18] Hoang Anh Dau et al. *The UCR Time Series Classification Archive*. [https://www.cs.ucr.edu/~eamonn/time\\_series\\_data\\_2018/](https://www.cs.ucr.edu/~eamonn/time_series_data_2018/). Oct. 2018.
  - [Dem06] Janez Demšar. “Statistical comparisons of classifiers over multiple data sets”. In: *The Journal of Machine Learning Research* 7 (2006), pp. 1–30.

- [Den+13] Houtao Deng et al. “A time series forest for classification and feature extraction”. In: *Information Sciences* 239 (2013), pp. 142–153.
- [DGM97] 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.
- [Din+08] 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.
- [Faw+19a] 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.
- [Faw+19b] 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.
- [Faw+20] Hassan Ismail Fawaz et al. “Inceptiontime: Finding alexnet for time series classification”. In: *Data Mining and Knowledge Discovery* 34.6 (2020), pp. 1936–1962.
- [FJ20] Johann Faouzi and Hicham Janati. “pyts: A Python Package for Time Series Classification”. In: *Journal of Machine Learning Research* 21.46 (2020), pp. 1–6. URL: <http://jmlr.org/papers/v21/19-763.html>.
- [FRM94] Christos Faloutsos, Mudumbai Ranganathan, and Yannis Manolopoulos. “Fast subsequence matching in time-series databases”. In: *Acm Sigmod Record* 23.2 (1994), pp. 419–429.
- [Gao+14] Jingkun Gao et al. “Plaid: a public dataset of high-resolution electrical appliance measurements for load identification research: demo abstract”. In: *proceedings of the 1st ACM Conference on Embedded Systems for Energy-Efficient Buildings*. 2014, pp. 198–199.
- [GB10] Xavier Glorot and Yoshua Bengio. “Understanding the difficulty of training deep feedforward neural networks”. In: *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings. 2010, pp. 249–256.
- [GL13] Tomasz Górecki and Maciej Luczak. “Using derivatives in time series classification”. In: *Data Mining and Knowledge Discovery* 26.2 (2013), pp. 310–331.
- [GM01] 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.
- [GO12] 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.

- [Gra+14] 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.
- [GRO14] 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.
- [Gua+19] Chaoyu Guan et al. “Towards a deep and unified understanding of deep neural models in nlp”. In: *International conference on machine learning*. PMLR. 2019, pp. 2454–2463.
- [Gup+20] Ashish Gupta et al. “A Fault-Tolerant Early Classification Approach for Human Activities using Multivariate Time Series”. In: *IEEE Transactions on Mobile Computing* (2020).
- [He+15] Guoliang He et al. “Early classification on multivariate time series”. In: *Neurocomputing* 149 (2015), pp. 777–787.
- [He+16] 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.
- [Hil+14] Jon Hills et al. “Classification of time series by shapelet transformation”. In: *Data Mining and Knowledge Discovery* 28.4 (2014), pp. 851–881.
- [HLT18] Huai-Shuo Huang, Chien-Liang Liu, and Vincent S Tseng. “Multivariate time series early classification using multi-domain deep neural network”. In: *2018 IEEE 5th International Conference on Data Science and Advanced Analytics (DSAA)*. IEEE. 2018, pp. 90–98.
- [Ita75] 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.
- [JJO11] 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.
- [Jog+17] Amod Jog et al. “Random forest regression for magnetic resonance image synthesis”. In: *Medical image analysis* 35 (2017), pp. 475–488.
- [JZ14] Zbigniew Jerzak and Holger Ziekow. “The DEBS 2014 Grand Challenge”. In: *Proceedings of the 8th ACM International Conference on Distributed Event-Based Systems*. DEBS ’14. Mumbai, India: Association for Computing Machinery, 2014, pp. 266–269. ISBN: 9781450327374. DOI: [10.1145/2611286.2611333](https://doi.org/10.1145/2611286.2611333). URL: <https://doi.org/10.1145/2611286.2611333>.

- [Kat16] 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.
- [KB14] Diederik P Kingma and Jimmy Ba. “Adam: A method for stochastic optimization”. In: *arXiv preprint arXiv:1412.6980* (2014).
- [KP01] 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.
- [KSH12] 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.
- [Lar+18] James Large et al. *From BOP to BOSS and Beyond: Time Series Classification with Dictionary Based Classifiers*. 2018. arXiv: [1809.06751 \[cs.LG\]](#).
- [LB15] 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.
- [LBH15] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. “Deep learning”. In: *nature* 521.7553 (2015), pp. 436–444.
- [LGI17] Thach Le Nguyen, Severin Gsponer, and Georgiana Ifrim. “Time series classification by sequence learning in all-subsequence space”. In: *2017 IEEE 33rd international conference on data engineering (ICDE)*. IEEE. 2017, pp. 947–958.
- [Lin+07] 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.
- [Lin+12] Jason Lines et al. “A shapelet transform for time series classification”. In: *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2012, pp. 289–297.
- [Lin+15] 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.
- [LKL12] 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.
- [LLF14] Kang Li, Sheng Li, and Yun Fu. “Early classification of ongoing observation”. In: *2014 IEEE International Conference on Data Mining*. IEEE. 2014, pp. 310–319.

- [LLY17] 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.
- [LMT16] 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.
- [Lön+19] Markus Löning et al. “sktime: A unified interface for machine learning with time series”. In: *arXiv preprint arXiv:1909.07872* (2019).
- [Low14] Richard Lowry. “Concepts and applications of inferential statistics”. In: (2014).
- [LTB18] 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).
- [Luc+19] 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.
- [Lv+19] Junwei Lv et al. “An effective confidence-based early classification of time series”. In: *IEEE Access* 7 (2019), pp. 96113–96124.
- [Mar08] 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.
- [MKY11] 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.
- [Mor+17a] 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.
- [Mor+17b] 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.
- [Mor+19] Usue Mori et al. “Early classification of time series using multi-objective optimization techniques”. In: *Information Sciences* 492 (2019), pp. 204–218.
- [MVB19] 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.

- [Par+13] Nathan Parrish et al. “Classifying with confidence from incomplete information”. In: *The Journal of Machine Learning Research* 14.1 (2013), pp. 3561–3589.
- [Ped+11] F. Pedregosa et al. “Scikit-learn: Machine Learning in Python”. In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.
- [Pet+16] 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.
- [Pro+06] Pavlos Protopapas et al. “Finding outlier light curves in catalogues of periodic variable stars”. In: *Monthly Notices of the Royal Astronomical Society* 369.2 (2006), pp. 677–696.
- [RA04] 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.
- [RK04] 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.
- [RK05] Chotirat Ann Ratanamahatana and Eamonn Keogh. “Three myths about dynamic time warping data mining”. In: *Proceedings of the 2005 SIAM international conference on data mining*. SIAM. 2005, pp. 506–510.
- [RK13] 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.
- [Rui+20] 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.
- [Ruß+19] Marc Rußwurm et al. “End-to-end learning for early classification of time series”. In: *arXiv preprint arXiv:1901.10681* (2019).
- [SAD12] 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.
- [SC78] 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.
- [Sch15] 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.

- [Sch16] Patrick Schäfer. “Scalable time series classification”. In: *Data Mining and Knowledge Discovery* 30.5 (2016), pp. 1273–1298.
- [SH12] Patrick Schäfer and Mikael Höglqvist. “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.
- [Shi+20] 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.
- [SK16] Tiago Santos and Roman Kern. “A Literature Survey of Early Time Series Classification and Deep Learning.” In: *Sami@ iknow*. 2016.
- [SL17a] 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.
- [SL17b] Patrick Schäfer and Ulf Leser. “Multivariate time series classification with WEASEL+ MUSE”. In: *arXiv preprint arXiv:1711.11343* (2017).
- [SL20] Patrick Schäfer and Ulf Leser. “TEASER: early and accurate time series classification”. In: *Data Mining and Knowledge Discovery* 34.5 (2020), pp. 1336–1362.
- [SM13] Pavel Senin and Sergey Malinchik. “Sax-vsm: Interpretable time series classification using sax and vector space model”. In: *2013 IEEE 13th international conference on data mining*. IEEE. 2013, pp. 1175–1180.
- [SPK18] Joan Serrà, Santiago Pascual, and Alexandros Karatzoglou. “Towards a Universal Neural Network Encoder for Time Series.” In: *CCIA*. 2018, pp. 120–129.
- [SSS17] Balraj Singh, Parveen Sihag, and Karan Singh. “Modelling of impact of water quality on infiltration rate of soil by random forest regression”. In: *Modeling Earth Systems and Environment* 3.3 (2017), pp. 999–1004.
- [SW17] Simone Scardapane and Dianhui Wang. “Randomness in neural networks: an overview”. In: *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 7.2 (2017), e1200.
- [Sze+15] Christian Szegedy et al. “Going deeper with convolutions”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 1–9.
- [Sze+17] 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.

- [Tan+19] Chang Wei Tan et al. “Time series classification for varying length series”. In: *arXiv preprint arXiv:1910.04341* (2019).
- [Tav+20] Romain Tavenard et al. “Tslearn, A Machine Learning Toolkit for Time Series Data”. In: *Journal of Machine Learning Research* 21.118 (2020), pp. 1–6. URL: <http://jmlr.org/papers/v21/20-091.html>.
- [TH16] 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.
- [TM16] 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.
- [TPW20] 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.
- [Wan+13] 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.
- [WYO17] 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.
- [Xin+11] Zhengzheng Xing et al. “Extracting interpretable features for early classification on time series”. In: *Proceedings of the 2011 SIAM International Conference on Data Mining*. SIAM. 2011, pp. 247–258.
- [XPK10] Zhengzheng Xing, Jian Pei, and Eamonn Keogh. “A brief survey on sequence classification”. In: *ACM Sigkdd Explorations Newsletter* 12.1 (2010), pp. 40–48.
- [XPP09] 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.
- [XPP12] Zhengzheng Xing, Jian Pei, and S Yu Philip. “Early classification on time series”. In: *Knowledge and information systems* 31.1 (2012), pp. 105–127.
- [YD19] Omolbanin Yazdanbakhsh and Scott Dick. “Multivariate Time Series Classification using Dilated Convolutional Neural Network”. In: *arXiv preprint arXiv:1905.01697* (2019).

- [YK09] 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.
- [Zei12] Matthew D Zeiler. “Adadelta: an adaptive learning rate method”. In: *arXiv preprint arXiv:1212.5701* (2012).
- [Zha+17] Bendong Zhao et al. “Convolutional neural networks for time series classification”. In: *Journal of Systems Engineering and Electronics* 28.1 (2017), pp. 162–169.
- [Zhe+14] 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.
- [Zhe+16] 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.

# Declaration of Authorship

I hereby declare that I have written this thesis "A comparison of Time Series Classification Algorithms based on their ability to learn in an Early Classification Context" 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, 25.05.2021