

UTRECHT UNIVERSITY

MASTER'S DEGREE THESIS

---

# Temporal Segmentation using Support Vector Machines in the context of Human Activity Recognition

---

*Author:*

R.Q. VLASVELD, UU

*Supervisors:*

Dr. G.A.W. VREESWIJK, UU

Anne VAN ROSSUM, Almende B.V.

*A thesis submitted in fulfilment of the requirements  
for the degree of Master of Science*

*in the*

Faculty of Science

Department of Information and Computing Sciences

January 2014

# **Declaration of Authorship**

I, R.Q. VLASVELD, declare that this thesis titled, 'Temporal Segmentation using Support Vector Machines in the context of Human Activity Recognition' and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

---

Date:

*“All truly great thoughts are conceived by walking.”*

Friedrich Nietzsche

UTRECHT UNIVERSITY

## *Abstract*

Faculty of Science  
Department of Information and Computing Sciences

Master of Science

**Temporal Segmentation using Support Vector Machines in the context of  
Human Activity Recognition**

by R.Q. VLASVELD

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centered vertically so can expand into the blank space above the title too...

**\*\*\* *TODO: Write abstract* \*\*\***

## *Acknowledgements*

The acknowledgements and the people to thank go here, don't forget to include your project advisor...

# Contents

<b>Declaration of Authorship</b>	i
<b>Abstract</b>	iii
<b>Acknowledgements</b>	iv
<b>List of Figures</b>	viii
<b>List of Tables</b>	ix
<b>Acronyms</b>	x
<b>1 Introduction</b>	2
1.1 Scope . . . . .	3
1.2 One-Class Support Vector Machines . . . . .	4
1.3 Data sets . . . . .	5
1.4 Contributions . . . . .	5
1.5 Thesis structure . . . . .	6
<b>2 Literature review</b>	7
2.1 Outline . . . . .	7
2.2 Statistical framework . . . . .	7
2.3 Overview . . . . .	8
2.4 Statistical models . . . . .	9
2.5 Temporal Segmentation . . . . .	10
2.5.1 Segmentation . . . . .	11
2.5.2 CUSUM . . . . .	14
2.5.3 Change detection . . . . .	15
2.5.4 Novelty and Outlier detection . . . . .	15
2.6 Change Detection methods . . . . .	16
2.6.1 CUSUM . . . . .	16
2.6.2 Model fitting . . . . .	16
2.6.3 Density-Ratio . . . . .	16
2.7 CUSUM . . . . .	16
2.8 Change-detection by Density-Ratio Estimation . . . . .	17
2.9 Support Vector Machines in Change Detection . . . . .	19

<b>3 Change detection by Support Vector Machines</b>	<b>21</b>
3.1 Problem Formulation Change Detection . . . . .	21
3.2 One-Class Classification . . . . .	22
3.2.1 Problem formulation . . . . .	23
3.2.2 One-Class Classification methods . . . . .	24
3.3 One-Class Support Vector Machine . . . . .	26
3.3.1 Support Vector Machine . . . . .	26
3.3.2 Kernels . . . . .	30
3.3.3 $\nu$ -Support Vector Machine . . . . .	31
3.3.4 Support Vector Data Description . . . . .	33
3.3.5 SVM model parameters . . . . .	35
<b>4 Application to Human Activity Time Series Data</b>	<b>37</b>
4.1 Algorithm overview . . . . .	38
4.2 Data Gathering . . . . .	40
4.2.1 Artificial data . . . . .	40
4.2.2 Real-world data . . . . .	40
4.3 Model Construction: Incremental SVDD . . . . .	41
4.4 Model Features . . . . .	43
4.5 Change Detection . . . . .	44
4.5.1 CUSUM based methods . . . . .	46
4.5.2 Ratio-thresholding . . . . .	46
4.5.3 Post-processing . . . . .	47
<b>5 Artificial data results</b>	<b>48</b>
5.1 Artificial data generation . . . . .	48
5.2 Quality metrics . . . . .	49
5.3 Results . . . . .	51
<b>6 Real-world results</b>	<b>54</b>
6.1 Data Sets . . . . .	54
6.2 Data Gathering . . . . .	56
6.2.1 Outdoor straight lines . . . . .	56
6.2.2 Outdoor free . . . . .	57
6.2.3 Indoor stairs . . . . .	58
6.3 Results . . . . .	58
6.3.1 Objective metrics . . . . .	60
6.3.2 Subjective observations and remarks . . . . .	61
6.4 Possible improvements . . . . .	64
<b>7 Conclusion</b>	<b>65</b>
7.1 Outline . . . . .	65
7.2 What is done . . . . .	66
7.3 Main findings . . . . .	68
7.4 Contributions . . . . .	68
7.5 Other research . . . . .	68
7.6 Limitations . . . . .	68
7.7 Future research . . . . .	68

7.8 Final section . . . . .	69
-----------------------------	----

<b>Bibliography</b>	<b>70</b>
---------------------	-----------

# List of Figures

1.1	Thesis goal	4
3.1	Difference between two and one-class classification	24
3.2	One-Class Classification (OCC) methods	25
3.3	Kernel mapping	27
3.4	Mapping spaces in Support Vector Machine (SVM)	27
3.5	SVM and the separating hyperplane	28
3.6	$\nu$ -Support Vector Machine ( $\nu$ -SVM)	32
3.7	Difference $\nu$ -SVM and Support Vector Data Description (SVDD)	33
3.8	SVDD boundary	35
3.9	SVDD boundary	36
4.1	Method setup	38
4.2	Expected radius behavior	43
4.3	Thresholding	47
5.1	Artificial data sets	50
5.2	Box plot results artificial data sets	52
5.3	Artificial data set results	53
6.1	WISDM Excerpt	55
6.2	UCI HAR Excerpt	55
6.3	Plots run 1	57
6.4	Plots run 5	58
6.5	Stills subject 2	59
6.6	Stills subject 3	59
6.7	Plots subject 3	60
6.8	Box plot results real-world runs	61

# List of Tables

4.1	Measured metrics . . . . .	41
4.2	Proposed algorithm . . . . .	45
5.1	Results artificial data sets . . . . .	51
6.1	Results real world runs . . . . .	61

# Acronyms

**AIC** Akaike Information Criterion. [13](#)

**AR** Autoregressive. [40, 45, 48](#)

**BIC** Bayesian Information Criterion. [13](#)

**changeFinder** is a unifying framework proposed by Takeuchi and Yamanishi [\[63\]](#) which combines outlier detection and change detection. [21, 38](#)

**CUSUM** Cumulative Sum. [9–12, 14–17, 39, 45, 46](#)

**DR** Dimensionality Reduction. [8](#)

**FAR** False Alarm Rate. [49, 58, 62, 63](#)

**GMM** Gaussian Mixture Model. [24](#)

**HAR** Human Activity Recognition. [2–4, 6, 11, 18, 19, 37](#)

**HMM** Hidden Markov Model. [13](#)

**ICSS** Iterated Cumulative Sums of Squares. [14, 15, 17](#)

**i.i.d.** independently and identically distributed. [22](#)

**I-SVDD** Incremental Support Vector Data Description. [39, 41–44, 67](#)

**KKT** Karush-Kuhn-Tucker. [42](#)

**KL** Kullback-Leibler. [19, 22, 23](#)

**KLIEP** Kullback-Leibler Importance Estimation Procedure. [12, 17, 18](#)

**MLE** Maximum Likelihood Estimates. [15](#)

**MOSUM** Moving Sum. 11, 14, 16

**OCC** One-Class Classification. viii, 2, 3, 8, 16, 21–26, 43, 66

**OC-SVM** One-Class Support Vector Machine. 3–6, 8, 19, 20, 22, 26, 31, 35–37, 41–43, 48, 67, 69

**PCA** Principal Component Analysis. 12, 13, 25

**PDF** Probability Density Function. 9, 23

**PLR** Piecewise Linear Representation. 11

**QP** Quadratic Programming. 29, 31

**RBФ** Radial Base Function. 30, 33–35, 39, 43, 50, 51, 60

**RKHS** Reproducing Kernel Hilbert Space. 30

**RT** Ratio-Thresholding. 39, 48, 50, 51

**SV** Support Vector. 32–35

**SVC** Support Vector Regression. 15, 29

**SVCPD** Support Vector based Change Point Detection. 6, 19, 21, 22, 37, 38, 41, 44, 45, 67

**SVDD** Support Vector Data Description. vi, viii, 2, 3, 5, 19, 21, 25, 26, 33–37, 39, 41–43, 50, 51, 55, 66, 67

**SVM** Support Vector Machine. viii, 2, 3, 6, 8, 9, 11, 18, 19, 21, 22, 25–33, 35, 38, 42

**$\nu$ -SVM**  $\nu$ -Support Vector Machine. vi, viii, 15, 18, 19, 21, 25, 26, 31–33, 35

**SWAB** Sliding Window And Bottom-up. 10–12

**UCI HAR** Human Activity Recognition Using Smartphones Data Set. 4, 54

**WISDM** Wireless Sensor Data Mining. 4, 54

– Overall todo's –

- \*\*\* *TODO: The term emperical data/distribution is nowhere used in the thesis. or: Empirical cumulative distribution. Especially Chapter 2.* \*\*\*
- \*\*\* *TODO: Check for “tang-constructions”: saying to much in a single sentence.* \*\*\*
- \*\*\* *TODO: Overall: meer focus op de waarom vraag/antwoord bij keuzes. Bouwt aan motivatie en context van het geheel.* \*\*\*
- \*\*\* *TODO: Consistent ‘data point’ of ‘data object’ gebruiken.* \*\*\*

# Chapter 1

## Introduction

With the wide availability of smartphones and built-in inertial sensors, more and more applications of Human Activity Recognition (HAR) are introduced [7, 21, 28]. Many approaches to recognize the performed activity rely on classification of recorded data [6, 8, 12, 22, 34, 74, 77], which can be done online or in batches after the performance [23]. Often a time-windowed based approach is used, which feeds short consecutive segments of data through a model construction and recognition phase. The constructed model is used to determine which (earlier learned) activity is performed. Besides the explicit classification of the data, an implicit obtained result is a segmentation of performed activities over time.

In this thesis the goal is to find the temporal segmentation of time series data explicitly, prior to the classification of the activities. This is done under the assumption that it could be beneficial for classification methods to possess this explicit segmentation, since the model construction phase can use more data than the time-windowed approach would allow. To find a temporal segmentation, we rely on change detection in time series data. In the context of this research, this time series data consists of recordings from inertial sensors found in smartphones, such as the accelerometer signal, during the performance of human activities. We are interested in both in- and outdoor activities in an uncontrolled environment and in a continuous manner, which can be regarded as real-world data. Although the number of performed activities will be limited, we assume the temporal segmentation must be more robust than in the case of data from a controlled environment.

For the change detection algorithm we have used a special form of Support Vector Machine (SVM), the One-Class Classification (OCC) based Support Vector Data Description (SVDD) as introduced by Tax [64]. This method models the data under consideration in the shape of a high-dimensional hypersphere. From this hypersphere the

radius is extracted, which can be used as an indication of change and be processed by one-dimensional change detection algorithms. Since a change in data characteristics, resulting from a change in performed activity, alters the data distribution, the model and thus the radius should reflect this change as well.

## 1.1 Scope

The scope of this thesis is to find a temporal segmentation of time series data from inertial sensors embedded in smartphones, obtained during the performance of human activities. As mentioned above, the segmentation is often obtained as a by-product from direct windowed classification. Other approaches have a very rigid form of segmentation, e.g. by limiting the number of possible segment types [15, 36], or have applied it to different type of time series data [25, 29, 46]. Others have extracted features from windows of data which are used for the segmentation [30]. These activities include sitting, standing, walking, running, and ascending and descending staircases. The problem of finding a temporal segmentation can be seen as a sub-problem to the problem of Activity Classification. This relation is illustrated in Figure 1.1. The figure shows that we can apply classification over windows of the time series data. In many approaches this is constructed with a (partially overlapping) sliding window. In the context of HAR, the window length commonly consists of around 2 seconds of data points. In our setup, we create a temporal segmentation instead of a fixed window length. Each segment should consist of homogeneous data, in the sense that the underlying data generating activity should not change within a segment. Likewise, the previous and next segments (if present) should be generated by a different performed activity.

To find these segments of homogeneous data, we employ a model construction phase. The One-Class Support Vector Machine (OC-SVM) model transforms the data points to a high-dimensional hyperspherical shape from which we extract the radius as characterizing metric. In this phase the data is processed by a fixed window length and a change in data characteristics should reflect in a (significant) change of the hyper-sphere's radius. Such a sudden change in the radius thus indicates a change in the time series data and indicates the end of the current and beginning of a new segment. In Section 2.5 a more detailed overview of temporal segmentation methods is discussed. The information about segments can be used in the classification phase of HAR, which applies classification to windows (or in this case: segments) of data. The final result, which is outside the scope of this thesis, will be a full classification of activities over the time series data.

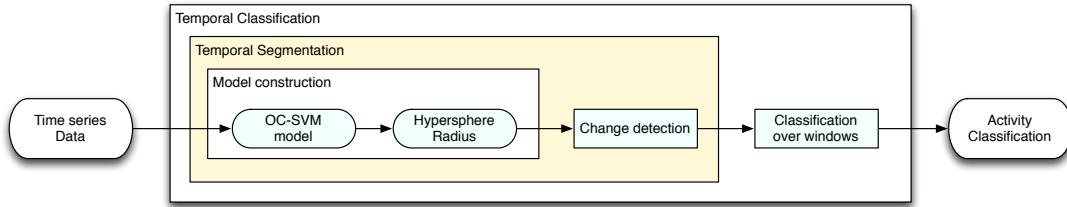


FIGURE 1.1: The goal of this thesis. The scope of this thesis is Temporal Segmentation, which can be useful in the context of Temporal Classification. Given a data set, we construct a OC-SVM high-dimensional spherical model, from which we extract the radius. This metric is then applied to direct change detection algorithms. The detected change points can support the classification of homogeneous segments of data.

## 1.2 One-Class Support Vector Machines

In the setup, as described in the previous section, we employ a model construction phase. For this model we will use an One-Class Support Vector Machine implementation: the Support Vector Data Description algorithm by Tax [64]. In earlier research multi-class SVMs are used for direct classification of inertial sensor data [6, 34, 53]. Others have used OC-SVM classifiers for novelty, or outlier, detection [14, 47, 50, 59, 65]. The research by Yin *et al.* [78] creates a OC-SVM classifier to filter normal activity traces, using training data. However, to our knowledge, OC-SVM methods have not yet been applied to inertial sensor time series data directly to find changes in behavior and create a temporal segmentation. This algorithm transforms the time series data, which can be of any high dimension, to a higher dimensional feature space. When the data is assumed to be of the same class, it can be characterized by an enclosing boundary, in the shape of a hypersphere. Using this OCC method, we are able to test new data points of being a member of this class. When a new data point, after transformation to the feature space, lies within the boundary, it is a member of the class, and otherwise it is regarded as an outlier.

In case the data distribution changes, the data points will have a larger variance which results in a larger hypersphere in the feature space. This means that with continuous model updating, a change in the underlying process of a time series data means the hypersphere will have a larger radius. Thus the OC-SVM method allows us to generate a one-dimensional time series data, reduced from a complex high-dimensional original signal. The final radius time series data can be inspected by simpler change detection methods, in a number of ways. Chapter 3 will discuss the OC-SVM methods in detail and Chapter 4 elaborates on the overall change detection method.

### 1.3 Data sets

In recent previous research two data sets for the context of HAR have been made public. In order to benchmark our proposed method we tried to employ it to these data sets. However, it turned out the data sets were not appropriate for continuous temporal segmentation. The first set, the Wireless Sensor Data Mining (WISDM) from [44], consists of non-continuous recordings of activities. Obviously, this voids the usefulness of the data set of our research. The second set, the Human Activity Recognition Using Smartphones Data Set (UCI HAR) from [6], also lacks continuous recordings and furthermore exhibits patterns and labeling which seems to be incorrect. Besides these problems, since there are no visual images of the performed activities, there is no ground truth which can be consulted in case of ambiguity. To overcome these problems, we have created our own real-world data sets.

Our own data sets are recorded by placing a smartphone with inertial sensors in the front right pants pocket of a subject. The subject was asked to perform common and simple activities in an uncontrolled environment, both in- and outdoor situated. During the performance the subjects were also recorded with a video camera from the third-person point of view. This allows us to create a ground truth for the temporal segmentation of performed activities.

Besides the recorded data sets we have used artificial data sets to obtain an objective performance measure of the method. These artificial data sets are modeled to the sets used in [14, 63]. In Chapter 5 we discuss the artificial data sets and results in detail. The real-world recordings and results are discussed in Chapter 6.

### 1.4 Contributions

As a result from our proposed method and application to recorded real-world data sets, this thesis will provide the following contributions:

1. **Continuous recorded human activity data:** the proposed method is applied to our own recorded human activity data. This data set can be grounded by using the video recordings of the subject while performing the activities. The data set, consisting of raw inertial sensor recordings, manual change point and activity labeling, and video recordings, are used for testing the method and are made publicly available [73].
2. **Application of OC-SVM methods to inertial sensor data:** to our best knowledge, this thesis is the first application of OC-SVM methods, especially the

SVDD method, to inertial sensor signals in the context of creating a temporal segmentation. Furthermore, this thesis contributes to the small amount of research that has a temporal segmentation of inertial sensor data as its primary objective.

3. **Adaptation of SVCPD:** we have adopted and slightly adapted the SVCPD method by Camci [14] in order to find change points in time series data. In our simplification we have not encountered noticeably decrease in performance.

## 1.5 Thesis structure

The structure of this thesis is as follows. In Chapter 2 a literature review is provided. We will look at the different interpretations and implementations for change, novelty, and outlier detection. Previous work in the field of HAR is discussed and we end with existing applications of SVMs to change detection. The next chapter further analyses SVMs and two implementations of OC-SVMs. It relates the properties of the constructed OC-SVM models to change detection in time series data. That relation is applied to our problem formulation in Chapter 4, where we construct our change detection method **\*\*\***  
***TODO: add name \*\*\****, based on the work of Camci [14]. The two following chapters apply the method to both artificial and real-world data. Chapter 5 show the objective performance of the method to artificial constructed data. In Chapter 6 we apply the method to our self-recorded real-world data. In that chapter we show the ability of detecting change in HAR data, recorded by inertial sensors. This thesis is concluded in Chapter 7, in which we reflect on the performed research and state possibilities for feature research.

# Chapter 2

## Literature review

### 2.1 Outline

*Not intented for the reader.*

- Literature review about Temporal Segmentation (previous draft was more about classification)
- Consider methods for the context of filter-methods for classification
- Take a loot at 3-4 different kind of methods for change detection:
  - Introduction with a lot of techniques
  - Explain why look at a few
  - CUSUM - or other more traditional methods
  - Density-ratio estimation
  - Support Vector Machines (?) - if there are more sources about this
  - (Dimensionality reduction) –*i* probably not
  - *Not to much attention to all techniques, focus is on SVM*
- With each method, shortly look at characteristics, strengths and weaknesses and consider applicability to accelerometer sensor data

### 2.2 Statistical framework

- Nieuwe indeling Chapter 2 –

**Intro/outline** Introduction and overview of this chapter. Put in the context of temporal classification.

**Temporal segmentation** Give overview of used temporal segmentation methods, move towards a split in model building/dimensionality reduction and change detection methods

**Dimensionality reduction** Explain that Dimensionality Reduction (DR) can be useful in data processing. Name a few techniques, make bridge to OC-SVM.

**Change detection** Describe methods for *direct* change detection. End with notion

**Model fitting**

**Cusum**

**Density ratio**

## 2.3 Overview

As described in Chapter 1, the focus of this research is the application of temporal segmentation to accelerometer data. This is with the assumption that explicit temporal segmentation can aid the process of temporal classification.

**\*\*\* TODO: Give examples of temporal classification techniques \*\*\***

In recent years a lot of research is performed in the classification of human activities, recorded by on-body sensors. Especially since smartphones with accelerometer sensors became more available, research has focused on recordings from these devices.

**\*\*\* TODO: Include highlights from old-classification text \*\*\***

In the following of this chapter we will start with an overview of statistical models used in this field of study. That section will act as a starting point to look at the different types, applications and methods of temporal segmentation. The following section, Section 2.5 follows the distinction made and looks at a collection of temporal segmentation and change detection methods. It will also relate the concept of OCC to temporal segmentation and change detection. The final section of this chapter, Section 2.9, will discuss a specific class of algorithms for OCC, namely the methods that use SVMs for the classification task.

## 2.4 Statistical models

**\*\*\* TODO: Introduce the different models and thus different approaches \*\*\***

Many applications require the detection of time points at which the underlying properties of a system change. This problem has received a lot of attention in the fields of data mining, etc... **\*\*\* TODO: list and refs \*\*\***. Often this problem is formulated in a statistical framework, by inspecting the underlying data generating Probability Density Function (PDF) of the time series data. A change point is then defined as a significant change in the properties of the PDF, such as the mean and variance.

The widely used Cumulative Sum (CUSUM) method by Basseville *et al.* [11] follows this approach. It originates from control methods for detection from benchmarks. This method and some derivatiges are discussed and analyzed in Section 2.7.

Many methods rely on pre-specified parametric model assumptions and considers the data be independent over time, which makes it less flexible to real-world applications. The methods proposed by Kawahara *et al.* [41] and Lui *et al.* [48] try to overcome these problems by estimating the *ratio* between the PDFs, instead of estimating each PDF. This approach is discussed and analyzed in Section 2.8.

The density-estimation methods rely on the log-likelihood ratio between PDFs. The method of Camci [14] follows an other approach within the statistical framework, by using a SVM. One problem it tries to overcome is the (claimed) weakness of many methods to detect a decrease in variance. The method represents the distribution over the data points as a hyper-sphere in a higher dimension using the kernel trick. A change in the PDF is represented by a change in the radius of this sphere. Section 2.9 discusses the SVM-methods. **\*\*\* TODO: Put emphasis that this is the source of inspiration for the chosen method? \*\*\***

– NEW –

In the field of temporal segmentation of time series data we can distinguish different approaches and common frameworks in which they perform. The problem is commonly formulated in the statistical framework and regarding the properties of the underlying data generating PDF fo the time series data. In that sense, a change point is then defined as a significant change in the properties of the PDF, such as the mean or variance.

In the search for the change in properties, temporal segmentation and change point detection methods can roughly be categorized in three methods in the way the data is processed, as discussed by Avci *et al.* [7]:

**Support Vector Machine based methods** **Top-Down** methods iteratively divide the signal in segments by splitting at the best location. The algorithm starts with two segments and completes when a certain condition is met, such as when an error value or number of segments  $k$  is reached.

- **Bottom-Up** methods are the natural complement to top-down methods. They start with creating  $n/2$  segments and iteratively join adjacent segments while the value of a cost function for that operation is below a certain value.
- **Sliding-window** methods are simple and intuitive for online segmenting purposes. It starts with a small initial subsequence of the time series. New data points are joined in the segment until the fit-error is above a threshold. Since the data is only processed very locally, these methods can yield in poor results [42].
- **Sliding Window And Bottom-up**, as introduced by Keogh *et al.* [42], combines the ability of the sliding window mechanism to process time series online and the bottom-up approach to create superior segments in terms of fit-error. The algorithm processes the data in two stages. The first stage is to join new data points in the current segment created by a sliding window. The second stage processes the data using Bottom-Up and return the first segment as the final result. Since this second stage retains some (semi-)global view of the data, the results are comparative with normal Bottom-Up.

It is clear that for the application of this research sliding-window and preferably Sliding Window And Bottom-up (SWAB)-based algorithms should be considered. In the following we will discuss classes of algorithms grouped by the type of decision function, assuming a SWAB-based data processing order.

A widely used class of methods is based on the CUSUM approach by Basseville *et al.* [11]. This approach originates from control methods for detection and faults and benchmarking. In Section 2.5.2 a further overview of this class of methods is discussed.

- Density
- Outlier based

## 2.5 Temporal Segmentation

*\*\*\* TODO: Is there any research of temporal segmentation on accelerometer/inertia signals? Note if there is/isn't \*\*\**

Give overview of used temporal segmentation methods. Problem is that it is not really well-defined.

- Given overview of segmentation techniques, for times series data
- Use different “point-of-views”, or terms
- “Segmentation”
- “Change detection”
- “Novelty detection”
- Specific view on SVMs

This section gives an overview of the literature on temporal segmentation in the context of HAR. It takes a look on different implementations and methodologies. A wide range of terms and subtle differences are used in the field, such as ‘segmentation’, ‘change detection’, ‘novelty detection’ and ‘outlier detection’. These will be the categorical terms for which we discuss the literature.

Aspects we will consider when discussing the literature are applicability to HAR, parameter tuning and ... \*\*\* ***TODO: What more aspects?*** \*\*\*

### 2.5.1 Segmentation

Methods that apply temporal segmentation on time series data can be roughly categorized in three different methods, as discussed by Avci *et al.* [7]. In that survey the distinction between *Top-Down*, *Bottom-Up*, and *Sliding-Window* approaches is made, based on the way data is processed by the algorithm. Since we are interested in online change detection, the literature discussed in this section will mainly be forms of sliding-window algorithms.

In Keogh *et al.* [42] a comparison of the aforementioned approaches is made. The research also introduces the SWAB method, which combines the simple and intuitive approach of the sliding-window approach with the (semi-)global view over the data from the bottom-up method. This proposed method is dependent on user provided parameters, setting the maximum error when performing the algorithm. Each segment is approximated by using a Piecewise Linear Representation (PLR). The provided error threshold controls the granularity and number of segments.

Other methods have been proposed, such as an adaptive threshold based on the signal energy by Guenterberg *et al.* [29], the adaptive CUSUM-based test by Alippi *et al.* [4] and

the Moving Sum (MOSUM) by Hsu [38] in order to eliminate this user-dependency. The latter of these methods is able to process the accelerometer values directly, although better results are obtained when features of the signal are processed, as done in the former method. Here the signal energy, mean and standard deviation are used to segment activities and by adding all the axial time series together, the Signal-To-Noise ration is increased, resulting in a robuster method.

The method of Guenterberg *et al.* [29] extracts features from the raw sensor signal to base the segmentation on other properties than the pure values. The method of Bernecker *et al.* [12] uses other statistical properties, namely autocorrelation, to distinguish periodic from non-periodic segments. Using the SWAB method the self-similarity of a one-dimensional signal is obtained. The authors claim that only a slight modification is needed to perform the method on multi-dimensional data. After the segmentation phase, the method of Bernecker *et al.* extracts other statistical features which are used in the classification phase.

The proposal of Guo *et al.* [30] dynamically determines which features should be used for the segmentation and simultaneously determines the best model to fit the segment. For each of the three dimensions features such as the mean, variance, covariance, correlation, energy and entropy are calculated. By extending the SWAB method, for every frame a feature set is selected, using an enhanced version of Principal Component Analysis (PCA). The research also considered the (Stepwise) Feasable Space Window as introduced by [49], but since it results in a higher error rate than SWAB the latter was chosen to extend. Whereas the before mentioned algorithms use a linear representation, this methods considers linear, quadratic and cubical representations for each segment. This differs from other methods where the model is fixed for the whole time series, such as [25], which is stated to perform inferior on non-stationary time series such as daily life.

The time series data from a sensor can be considered as being drawn from a certain stochastic process. Probabilistic models can be constructed on that signal, yielding in probabilistic and Bayesian based segmentation methods. The CUSUM-methods follows a statistical approach and relies on the log-likelihood ratio [32] to measure the difference between two distributions. To calculate the ratio, the probability density functions need to be calculated. The method of Kawahara *et al.* [41] proposes to estimate the ratio of probability densities (known as the *importance*), based on the log likelihood of test samples, directly, without explicit estimation of the densities. The method by Liu *et al.* [48] uses a comparable dissimilarity measure using the Kullback-Leibler Importance Estimation Procedure (KLIEP) algorithm. They claim this results in a robuster approach for real-world scenarios. Although this is a model-based method, no strong assumptions (parameter settings) are made on the models.

The method of Adams and MacKay [1] builds a probabilistic model on the segment run length, given the observed data so far. Instead of modeling the values of the data points as a probabilistic distribution, the length of segments as a function of time is modeled by calculating the posterior probability. It uses a prior estimate for the run length and a predictive distribution for newly-observed data, given the data since the last change point. This method contrasts with the approach of Guralnik and Srivastava [31] in which change points are detected by a change in the (parameters of an) underlying, observed, model. For each new data point, the likelihoods of being a change point and part of the current segment are calculated, without a prior model (and thus is a non-Bayesian approach). It is observed that when no change point is detected for a long period of time, the computational complexity increases significantly.

Another application of PCA is to characterize the data by determining the dimensionality of a sequence of data points. The proposed method of Berbić *et al.* [9] determines the number of dimensions (features) needed to approximate a sequence within a specified error. With the observation that more dimensions are needed to keep the error below the threshold when transitions between actions occur, cut-points can be located and segments will be created. The superior extension of their approach uses a Probabilistic PCA algorithm to incorporate the dimensions outside the selected set as noise.

In the method by Himberg *et al.* [36] a cost function is defined over segments of data which is to be minimized. The cost functions thereby searches for internally homogeneous segments of data, reflecting states in which the devices and the user are. The cost function can be any arbitrary function and in the implementation the sum of variances over the segments is used. Both in a local and global iterative replacement procedure (as an alternative for the computationally hard dynamic programming algorithm) the best breakpoint locations  $c_i$  for a pre-defined number of segments  $1 \leq i \leq k$  are optimized.

Many methods obtain an implicit segmentation as a result of classification over a sliding window \*\*\* ***TODO: add refs*** \*\*\*. The method of Yang *et al.* [76] explicitly performs segmentation and classification simultaneously. It argues that the classification of a pre-segmented test-sequences becomes straightforward with many classical algorithms to choose from. The algorithm matches test examples with the *sparsest* linear representation of mixture subspace models of training examples, searching over different temporal resolutions.

The method of Chamroukhi *et al.* [15] is based on a Hidden Markov Model (HMM) and logistic regression. It assumes a  $K$ -state hidden process with a (hidden) state sequence, each state providing the parameters (amongst which the order) for a polynomial. The order of the model segment is determined by model selecting, often using the Bayesian

Information Criterion (BIC) or the similar Akaike Information Criterion (AIC) [3], as in [34].

Field of computer vision: [79], [46].

“Segmentation and Recognition of Motion Streams by Similarity Search” [46]. 29, 2007

“Aligned Cluster Analysis for Temporal Segmentation of Human Motion” [79]. 63, 2008

## 2.5.2 CUSUM

**\*\*\* TODO: Link this subsection with previous and next \*\*\***

*Notes:* Non-Bayesian change detection algorithm (i.e. no prior distribution believe available for the change time). The CUSUM method is developed by Page [56] for the application of statistical quality control (it is also known as a control chart). Primary for detection of mean shift. The MOSUM of squares test for monitoring variance changes [38]. Use of Cumulative Sums of Squares for Retrospective Detection of Changes of Variance [39]

An often used approach in the statistical framework of change detection is the CUSUM as introduced by Page [56]. Originally used for quality control in production environments, its main function is to detect change in the mean of measurements and has been applied to this problem [11]. It is a non-Bayesian method and thus makes no assumptions, in the form of prior belief distributions, for the change points. Many extensions to this method have been proposed. Some focus on the change in mean, such the method of Alippi and Roveri [4]. Others apply the method to problems in which the change of variance is under consideration. Among others are there the centered version of the cumulative sums, introduced by Brown, Durbin and Evans [13] and the MOSUM of squares by [38].

The method of Inclán and Tiao [39] builds on the centered version of CUSUM [13] to detect changes in variance. Using the batch Iterated Cumulative Sums of Squares (ICSS) algorithm they are able to find multiple change points, offline while post-processing the data. Let  $C_k = \sum_{i=1}^k \alpha_t^2$  be the cumulative sum of squares for a series of uncorrelated random variables  $\{\alpha_t\}$  of length  $T$ . The centered (and normalized) sum of squares is defined as

$$D_k = \frac{C_k}{C_T} - \frac{k}{T}, \quad k = 1, \dots, T, \quad \text{with } D_0 = D_T = 0. \quad (2.1)$$

For a series with homogeneous variance, the value of  $D_k$  will oscillate around 0. In case of a sudden change, the value will increase and exceed some predefined boundary with

high probability. By using an iterative algorithm, the method is able to minimize the masking effect of successive change points.

One of the motivations for the ICSS algorithm was the heavy computational burden involved with Bayesian methods, which need to calculate the posterior odds for the log-likelihood ratio testing. The ICSS algorithm avoids applying a function at all possible locations, due to the iterative search. The authors recommend the algorithm for analysis of long sequences.

### 2.5.3 Change detection

**\*\*\* TODO: change order of this and temporal segmentation sections? So first change, then segmentation? \*\*\***

Whereas the above mentioned researches focus on *segmentation*, many have focused on *change detection*. Although these techniques are closely related, there is a subtle difference. In the case of *change detection* to goal is to find, possibly unrelated, sudden change points in a signal [63]. In contrast, the goal of *temporal segmentation* is to find homogeneous segments of data, which can be the result of multiple detected changes.

The ICSS by Inclán and Tiao [39] is a statistical method which obtains results (when applied to stock data) comparable to Maximum Likelihood Estimates (MLE) and Bayesian **\*\*\* TODO: Bayesian What? \*\*\***. Whereas CUSUM can be applied to search for a change in mean, the ICSS is adapted to find changes in variance. It obtains a *likelihood ratio* for testing the hypothesis of one change against no change in the variance. Using an iterative approach, all possible change points are considered. The proposal of [18] extends on the CUSUM-based methods to find change points in mean and variance, by creating a more efficient and accurate algorithm.

Section 3.1 discusses the problem of change detection in time series further, and gives a formal problem definition.

**\*\*\* TODO: move CUSUM based techniques to this subsection \*\*\***

### 2.5.4 Novelty and Outlier detection

In the case of time series data, many approaches rely on the detection of outlier, or novel, data objects. An increase in outlier data objects indicates a higher probability of change, as discussed by Takeuchi and Yamanischi [63]. The concepts of novelty detection and outlier detection are closely related and often used interchangeable.

The algorithm by Ma and Perkins [51] uses Support Vector Regression (SVC) to model a series of past observations. Using the regression, a matching function  $V(t_0)$  is constructed which uses the residual value of  $t_0$  to create an outlier (out-of-class) confidence value. The method defines novel *events* as a series of observations for which the outlier confidence is high enough. In an alternative approach by the same authors [50] a classification method, in contrast with regression, is used to detect outlier observations. The latter approach uses the  $\nu$ -Support Vector Machine ( $\nu$ -SVM) algorithm by Schölkopf *et al.* [59], which is further discussed in Section 2.9

A extensive overview of novelty detection based on statistical models is given by Markou and Singh [51]. In the second part of their review [52], novelty detection by a variety of neural networks is discussed. The survey by Hodge and Austin [37] distinguishes between three types of outlier detection: 1) unsupervised, 2) supervised, and 3) semi-supervised learning. As we will see in Section 3.2, in this research we are interested in methods from Type 3, from which OCC is a member. More theoretical background and foundation on outliers in statistical data can be found in [10].

## 2.6 Change Detection methods

Describe methods for *direct* change detection. End with notion that all work on low/single dimensional time series data, so dimensionality reduction is needed.

### 2.6.1 CUSUM

### 2.6.2 Model fitting

### 2.6.3 Density-Ratio

– OLD –

## 2.7 CUSUM

*Notes:* Non-Bayesian change detection algorithm (i.e. no prior distribution believe available for the change time). The CUSUM method is developed by Page [56] for the application of statistical quality control (it is also known as a control chart). Primary for detection of mean shift. The MOSUM of squares test for monitoring variance changes

[38]. Use of Cumulative Sums of Squares for Retrospective Detection of Changes of Variance [39]

An often used approach in the statistical framework of change detection is the CUSUM as introduced by Page [56]. Originally used for quality control in production environments, its main function is to detect change in the mean of measurements and has been applied to this problem [11]. It is a non-Bayesian method and thus makes no assumptions, in the form of prior belief distributions, for the change points. Many extensions to this method have been proposed. Some focus on the change in mean, such the method of Alippi and Roveri [4]. Others apply the method to the problems in which the change of variance is under consideration. Among others are there the centered version of the cumulative sums, introduced by Brown, Durbin and Evans [13] and the MOSUM of squares by [38].

The method of Inclán and Tiao [39] builds on the centered version of CUSUM [13] to detect changes in variance. Using the batch ICSS algorithm they are able to find multiple change points, offline while post-processing the data. Let  $C_k = \sum_{i=1}^k \alpha_t^2$  be the cumulative sum of squares for a series of uncorrelated random variables  $\{\alpha_t\}$  of length  $T$ . The centered (and normalized) sum of squares is defined as

$$D_k = \frac{C_k}{C_T} - \frac{k}{T}, \quad k = 1, \dots, T, \quad \text{with } D_0 = D_T = 0. \quad (2.2)$$

For a series with homogeneous variance, the value of  $D_k$  will oscillate around 0. In case of a sudden change, the value will increase and exceed some predefined boundary with high probability. By using an iterative algorithm, the method is able to minimize the masking effect of successive change points.

One of the motivations for the ICSS algorithm was the heavy computational burden involved with Bayesian methods, which need to calculate the posterior odds for the log-likelihood ratio testing. The ICSS algorithm avoids applying a function at all possible locations, due to the iterative search. The authors recommend the algorithm for analysis of long sequences.

## 2.8 Change-detection by Density-Ratio Estimation

Many approaches to detect change points monitor the logarithm of the likelihood ratio between two consecutive intervals. A change point is regarded to be the moment in time when the underlying probabilistic generation function changes. Some methods which rely on this are novelty detection, maximum-likelihood estimation and online learning of autoregressive models [41]. A limitation of these methods is that they rely on pre-specified parametric models. Non-parametric models for density estimation have

been proposed, but it is said to be a hard problem [33, 62]. A solutions to this is to estimate the *ratio* of probabilities instead of the probabilities themselves. One of the recent methods to achieve this is the KLIEP by Sugiyama *et al.* [61].

The method proposed by Kawahara and Sugiyama [41] is composed of an online version of the KLIEP algorithm. The method also considers *sequences* of samples (rather than samples directly) because the time series samples are generally not independent over time. An advantage over other non-parametric approaches, such as sequential one-class support vector machines, is that the model has an natural cross-validation procedure. This makes that the value of tuning parameters, such as the kernel bandwidth, can be objectively obtained.

In their formulation change is detected by monitoring the logarithm of the likelihood ratio between the reference (past) and test (current) time intervals

$$S = \sum_{i=1}^{n_{\text{te}}} \ln \frac{p_{\text{te}}(\mathbf{Y}_{\text{te}}(i))}{p_{\text{rf}}(\mathbf{Y}_{\text{te}}(i))} \quad (2.3)$$

Where  $\mathbf{Y}_{\text{te}}(i)$  is a sequence of samples from the test interval. A change is detected when  $S > \mu$ , for some predetermined threshold  $\mu$ . The question is then how to calculate the density ratio

$$w(\mathbf{Y}) := \frac{p_{\text{te}}(\mathbf{Y})}{p_{\text{rf}}(\mathbf{Y})} \quad (2.4)$$

because this ratio is unknown and should be estimated. The naive approach is to estimate the ratio by taking the ratio of the estimated densities. Since this is known to be a hard problem and sensitive for errors, the solution would be to estimate the ratio directly.

The procedure of the method proposed by Kawahara and Sugiyama [41] is to first apply the batch KLIEP algorithm with model selection for initial parameter  $\alpha$  and kernel width calculation. Then for every new sample the reference and test intervals are shifted and the calculated parameters  $\alpha$  are updated. Finally the logarithm of the likelihood ratio is evaluated. If it is above the threshold  $\mu$  the current time is reported as a change point.

Improvements in this line of research by Liu *et al.* [48] has led the application of improved density-ratio estimation methods to the problem of change detection. Such an improvement is the Unconstrained Least-Squares Importance Fitting (uLSIF) method [40] and an extension which possesses a superior non-parametric convergence property: Relative uLSIF (RuLSIF) [75].

## 2.9 Support Vector Machines in Change Detection

Many proposals in the field of remote sensing and HAR make use of SVMs as a (supervised) model learning method. An elaborate overview of applications is given in [53]. In this section we review methods of change detection based on the applications of SVMs as model construction method. A number of proposals have been made using this method. Schölkopf *et al.* [59] applies Vapnik's principle, to never solve a problem which is more general than the one that one is actually interested in, to novelty detection. In the case of novelty detection, they argue there is no need for density estimation of the distribution. Simple algorithms estimate the density by considering how many data points fall in a region of interest. The  $\nu$ -SVM method instead starts with the number of data points that should be within the region and estimates a region with that desired property. It builds on the method of Vapnik and Vladimir [71], which characterizes a set of data by separating it from the origin. The  $\nu$ -SVM method adds the kernel method, allowing non-linear decision functions, and incorporates 'softness' by the  $\nu$ -parameter. Whereas the method in [71] focuses on two-class problems, the  $\nu$ -SVM method introduces the method the one-class problems.

The method introduced by Ma and Perkins [50] creates a *projected phase* of a time series data, which is intuitively equal to applying a high-pass filter to the time series. The projected phase of the time series combines a history or data points to a vector, which are then classified by the  $\nu$ -SVM method. The method is applied to a simple synthetic sinusoidal signal with small additional noise and a small segment with large additional noise. The algorithm is able to detect that segment, without false alarms.

The algorithms of SVMs has been applied to the problem of HAR, as by Anguita *et al.* [6]. In that research a multi-class classification problem is solved using SVMs and the One-Vs-All method. The method exploits fixed-point arithmetic to create a hardware-friendly algorithm, which can be executed on a smartphone.<sup>1</sup>

With the same concepts as  $\nu$ -SVM, the SVDD method by Tax and Duin [65, 67] uses a separating hypersphere (in contrast with a hyperplane) to characterize the data. The data points that lie outside of the created hypersphere are considered to be outliers, which number is a pre-determined fraction of the total number of data points.

The method by Yin *et al.* [78] uses the SVDD method in the first phase of a two-phase algorithm to filter commonly available normal activities. The Support Vector based Change Point Detection (SVCPD) method of Camci [14] uses SVDD applies it to

---

<sup>1</sup>Smartphones have limited resources and thus require energy-efficient algorithms. In [5] a Hardware-Friendly SVM is introduced. It uses less memory, processor time and power consumption, with a loss of precision.

time-series data. By using an OC-SVM the method is able to detect changes in mean variance. Especially the detection of variance *decrease* is an improvement over other methods that rely in the Kullback-Leibler (KL) divergence, since these are unable to detect decreases in variance [63]. This main research subject of this thesis is to apply the method of Camci [14] to sensor data (such as accelerometer time series) obtained by on-body smartphones.

The inner workings of OC-SVMs are discussed in detail in Section 3.3.

**\*\*\* TODO: Add notion of ‘supervised’ vs ‘unsupervised’ learning. OC-SVM is unsupervised. \*\*\***

# Chapter 3

## Change detection by Support Vector Machines

This chapter discusses the concepts and algorithms that will be used as a basis for the proposed method, as introduced in Chapter 4. The first section formulates the problem of change detection and relates it to outlier and novelty detection. It transforms the traditional problem of outlier detection into change detection for times series data. In Section 3.2 the problem of One-Class Classification is explained and discusses a number of implementations. In the final sections of this chapter, Section 3.3, two SVM-based OCC methods,  $\nu$ -SVM and SVDD, and the influence of model parameters are further discussed in detail.

### 3.1 Problem Formulation Change Detection

The problems of outlier and novelty detection, segmentation, and change detection in time series data are closely related. The terminology depends on the field of application, but there are subtle differences. The problem of outlier detection is concerned with finding data objects in a data set which have small resemblance with the majority of the data objects. These objects can be regarded as erroneous measurements. In the case of novelty detection these objects are considered to be member of a new class of objects. The unifying framework of Takeuchi and Yamanishi [63], changeFinder, creates a two stage process expressing change detection in terms of outlier detection. The first stage determines the outliers in a time series by giving a score based on the deviation from a learned model, and thereby creates a new time series. The second stage runs on that new created time series and calculates a average over a window of the outlier scores. The problem of change detection is then reduced to outlier detection over that

average-scored time series. The implementation by Camci [14], SVCPD, implements outlier detection with the SVDD algorithm to detect changes in mean and variance.

The problem of change point detection can be formulated using different type of models, as discussed in 2.5.1. The methods by Takeuchi and Yamanishi [63] and Camci [14] use the following formulation for change detection, which we will also use for our problem formulation. The algorithm searches for sudden changes in the time series data. In other words, slowly changing properties in the data are not considered to be changes. Considered a time series  $x_1x_1\dots$ , which is drawn from a stochastic process  $p$ . Each  $x_t$  ( $t = 1, 2, \dots$ ) is a  $d$ -dimensional real valued vector and  $p$  a probability density function of the sequence  $x_1x_2\dots$ . Assume  $p$  can be decomposed in two different independently and identically distributed (i.i.d.) stationary stochastic processes  $p^1$  and  $p^2$  and are one-dimensional Gaussian density functions. For a time point  $a$  data points for which  $t < a$  are drawn from  $p^1 = N(\mu_1, \sigma_1^2)$  and for  $t \geq a$  from  $p^2 = N(\mu_2, \sigma_a^2)$ . If  $p^1$  and  $p^2$  are different, then the time point  $t = a$  is a *change point*. In [63] the similarity between the stochastic processes are expressed by the KL divergence  $D(p^2||p^1)$ . It is observed that this measure is not able to detect a change by decrease in variance [14, 63]. This problem is the motivation for Camci [14] to create the SVCPD algorithm.

The definition of change point being sudden changes in the time series data is in line with the search of changes in activities. Since we are only interested in different activities (which are represented by sudden changes), slight changes within an activity are not of interest.

## 3.2 One-Class Classification

As discussed in the previous section, change detection in time series can be implemented using outlier, or novelty, detection. To regard a data point as an outlier, there must be a model of the normal time series data and a (dis)similarity measure defined over the model and data objects. When a data point differs enough from the created model, it can be labeled as an outlier. The class of OCC algorithms is especially designed for that purpose. The algorithms build up a model of the data, assuming only normal data objects are available (or a very limited amount of example outlier data objects). This is also known as novelty detection or semi-supervised detection and is of Type 3 in the system by Hodge and Austin [37]. This differs from classical classification algorithms, which commonly rely of both positive and negative examples.

In Section 3.2.1 the problem formulation of OCC methods is explained. In section 3.2.2 an overview of OCC methods is given. The following section will discuss one specific set of methods, the OC-SVM which use SVMs to model the normal training data.

### 3.2.1 Problem formulation

The problem of OCC is closely related to the (traditional) two-class classification situation<sup>1</sup>. In the case of traditional classification algorithms, the problem is to classify an unknown object to one of the pre-defined categories. Every object  $i$  is represented as a vector  $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d}), x_{i,j} \in \mathbb{R}$  with  $d$  real-valued measurements. Using this notation, an object  $\mathbf{x}_i$  thus represents one point in a feature space:  $\mathbf{x}_i \in \mathcal{X} = \mathbb{R}^d$ . The two classes of objects,  $\omega_1$  and  $\omega_2$ , are labeled  $-1$  and  $+1$  respectively. The objects with a label  $y_i \in \{-1, +1\}$  attached are in the training set  $\mathcal{X}^{tr}$  (note that it can be both positive and negative example objects). This problem is solved by determining a decision boundary in the feature space of the data objects and label the new data object based on the location relative to this boundary. This is expressed by the decision function  $y = f(\mathbf{x})$ :

$$f : \mathbb{R}^d \rightarrow \{-1, +1\} \quad (3.1)$$

In case of the OCC problem, only one class (often referred as the target class, or positive examples) of training data is used to create a decision boundary. The goal is to determine whether a new data object belongs to the target class or to the unknown class and thus is an outlier. One could argue that this problem is equal to the traditional two-class problem by considering all other classes as negative examples, although there are important differences. In pure OCC problems there are no negative example objects available. This could be because the acquisition of these examples is very expensive, or because there are only examples of the ‘normal’ state of the system and the goal is to detect ‘abnormal’ states. Since the algorithm’s goal is to differentiate between normal and abnormal objects (relative to the training objects), OCC is often called outlier, novelty or anomaly detection, depending on the origin of the problem to which the algorithm is applied<sup>2</sup>. The difference between two and one-class classification and the consequence for outlier objects is illustrated in Figure 3.1. In the two-class classification problem the object **o** will be member of the **-1** class whilst the OCC problem will label it as an outlier. In [64] a more detailed analysis of the OCC is given.

The OCC algorithms have been applied to a wide range of applications. The first is, obviously, outlier detection of objects which do not resemble the bulk of the training

---

<sup>1</sup>Two-class problems are considered as to be the basic problem, since multi-class classification problems can be decomposed into multiple two-class problems [26].

<sup>2</sup>The term One-Class Classification originates from Moya *et al.* [54].

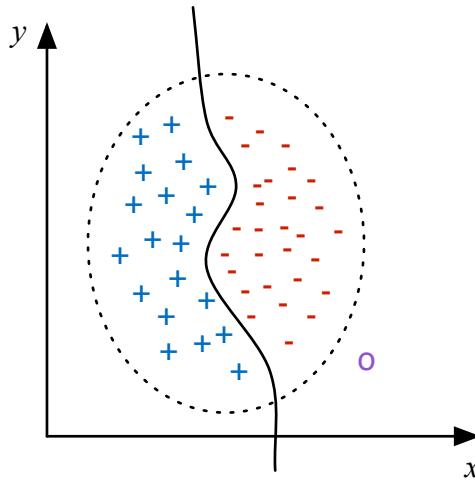


FIGURE 3.1: This plot shows the difference between two and one-class classification. The solid line indicates a possible decision boundary between the  $+1$  and  $-1$  example objects. The dashed circle indicates the closed boundary around all the data objects. In the first type the object  $\text{o}$  is considered to be member of the  $-1$ -class, whilst in the latter (OCC) formulation it is an outlier.

data. It can be a goal by itself and can be used as a filtering mechanism for other data processing methods. Often methods that rely on data-characteristics are sensitive for remote regions in the data set. Using OCC these remote regions can be removed from the data set. A second application is for the problem as described above, in which only data from a single target class is available. When this problems originates from e.g. a monitoring process, the OCC is able to recognize abnormal system states, without the need to create (or simulate) faulty states beforehand. The final possible application given by Tax [64] is the comparison of two data sets. By constructing a OCC-classifier using a training data set, one can compare that set to a new data set. This will result in a similarity-measure expressed in the properties of outliers. This measure is related to other similarity measures, such as PDF estimation and the KL divergence as discussed in Section 3.1.

### 3.2.2 One-Class Classification methods

The methods and algorithms used for the OCC-problem can be organized into three categories [55, 64], visually represented in Figure 3.2. The first category consists of methods that estimate the density of the training data and set a threshold on this density. Among those are Gaussian models, Gaussian Mixture Models (GMMs) and Parzen density estimators. In order to get good generalization results with these methods, the dimensionality of the data and the complexity of the density need to be restricted. This can cause a large bias on the data. When a good probability model is assumed, these

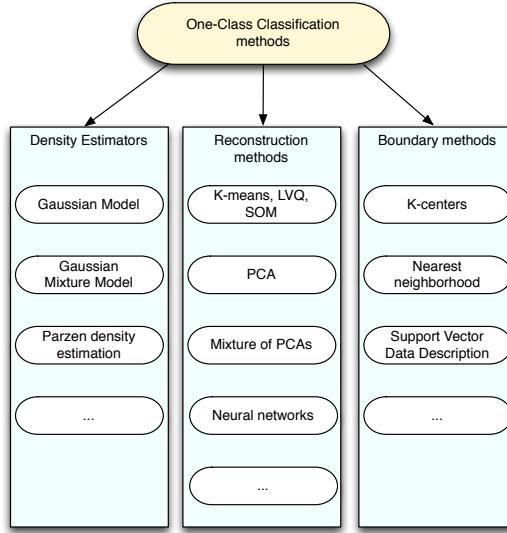


FIGURE 3.2: Overview of OCC methods categorized in Density Estimators, Reconstruction methods and Boundary methods. This categorization follows the definition of Tax in [64].

methods work very well, since when one threshold is optimized, a minimum volume is automatically found for the given probability density model [64].

Boundary methods are based on Vapnik's principle<sup>3</sup> which imply in this case that estimating the complete data density for a OCC may be too complex, in case one is only interested in the closed boundary. Examples of methods that focus on the boundary (a direct threshold) of the training data distribution are K-centers, Nearest-neighborhood and SVDD or a combination of those methods [35]. Especially the latter has a strong bias towards minimal volume solutions. These type of methods are sensitive to scaling of features, since they rely on a well-defined distance measure. The number of objects that is required is smaller than in the case of density methods. The boundary method SVDD, constructed by Tax and which has shown good performance [43], will be further discussed in Section 3.3.4.

Reconstruction methods take a different approach. Instead of focusing on classification of the data and thus on the discriminating features of data objects, they model the data. This results in a compact representation of the target data and for any object a reconstruction error can be defined. Outliers are data objects with a high reconstruction error, since they are worse represented by the constructed model. Examples of reconstruction methods are K-means, PCA and different kind of neural network implementations.

In [43, 55] an overview of applications for OCC-algorithms, and explicitly for SVM-based methods (such as SVDD and  $\nu$ -SVM), is given. It shows successful applications

<sup>3</sup>With a limited amount of data available, one should avoid solving a more general problem as an intermediate step to solve the original problem [72].

for, amongst others, problems in the field of Handwritten Digit Recognition, Face Recognition Applications, Spam Detection and Anomaly Detection [47, 57]. As discussed in Section 3.1, this can be used for change detection in time series data.

In this Section we have discussed the problem of OCC and different kind of implementations. An often used implementation is the SVM-based method [55], since it shows good performance in comparative researches [43, 60]. In the following section (3.3) two implementations of OC-SVM will be discussed, the SVDD method of Tax and Duin [65] and the  $\nu$ -SVM-algoritm by Schölkopf.

### 3.3 One-Class Support Vector Machine

In this section we will discuss the details of a SVM and OC-SVM implementations. The classical implementation of a SVM is to classify a dataset in two distinct classes. This is a common usecase, although sometimes there is no training data for both classes available. Still, one would like to classify new data points as regular, in-class, or out-of-class, e.g. in the case of a novelty detection. With that problem only data examples from one class is available and the objective is to recognize new data points that are not part of that class. This unsupervised learning problem is closely related to density estimation. In that context, the problem can be the following. Given an underlying probability distribution  $P$  the goal is to find a subset  $S$  for which the probability that a data point from  $P$  lies outside  $S$  equals some predetermined value  $\nu$  between 0 and 1 [59].

In the following of this section we will start with a discussion of traditional SVMs. In Section 3.3.2 we will show how kernels allow for non-linear decision functions. That is followed by two different implementations of OC-SVM:  $\nu$ -SVM in Section 3.3.3 by Schölkopf *et al.* [59], which closely follows the above problem statement regarding density estimation, and SVDD by Tax and Duin [65] in Section 3.3.4. The final part of this section discusses the influence of model parameters on the performance of SVMs

#### 3.3.1 Support Vector Machine

We will first discuss the traditional two-class SVM before we consider the one-class variant, as introduced by Cortes and Vapnik in [20]. Consider a data set  $\Omega = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ ; points  $x_i \in \mathbb{R}^d$  in a (for instance two-dimensional) space where  $x_i$  is the  $i$ -th input data point and  $y_i \in \{-1, 1\}$  is the  $i$ -th output pattern, indicating the class membership.

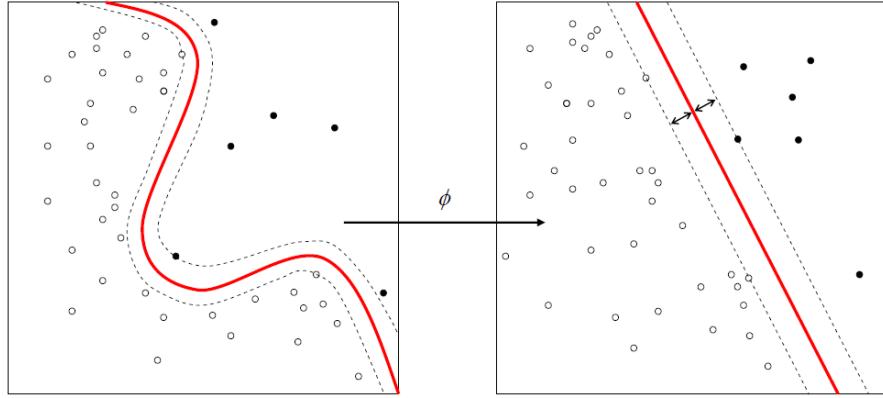


FIGURE 3.3: The non-linear boundary in the input space  $\mathcal{I}$  (left) is transformed to a linear boundary in the feature space  $\mathcal{F}$  (right) by mapping the data points with the function  $\phi$ . The kernel trick uses a function  $K$  which performs an implicit mapping.

Image from Wikipedia.org

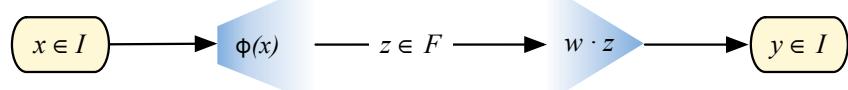


FIGURE 3.4: The input vector  $x$  from input space  $\mathcal{I}$  is mapped by a non-linear function  $\phi(x)$  to a feature vector  $z$  in the high dimensional feature space  $\mathcal{F}$ . The weights of  $w$  create a linear separating hyperplane, which maps the high dimensional vector to the predicted outcome  $y$ .

A SVM can create a boundary between linear-separable data points, making it a non-probabilistic binary linear classifier. More flexible non-linear boundaries can be obtained by the use of a non-linear function  $\phi(x)$ , as illustrated in Figure 3.3. This function maps the input data from space  $\mathcal{I}$  to a higher dimensional space  $\mathcal{F}$ . The SVM can create a linear separating hyperplane in the space  $\mathcal{F}$  that separates the data points from the two classes. When the hyperplane is projected to the (lower) original input space  $\mathcal{I}$  it creates a non-linear separating curve. A schematic overview of this process is illustrated in Figure 3.4. The original data points  $x \in \mathcal{I}$  are through  $\phi(x)$  projected to the points  $z \in \mathcal{F}$ . In that space the linear hyperplane is constructed by weights  $w$ . This results in a non-linear hyperplane  $y$  in input space  $\mathcal{I}$ . The mapping and projection of data points can be efficient (and implicit) performed by using the kernel trick, which is discussed in Section 3.3.2.

The separating hyperplane is represented by

$$w^T x + b = 0, \quad (3.2)$$

with  $w \in \mathcal{F}$  and  $b \in \mathbb{R}$ . The hyperplane that is created determines the *margin* between the classes; the minimal distance from one of the data points to the hyperplane. In

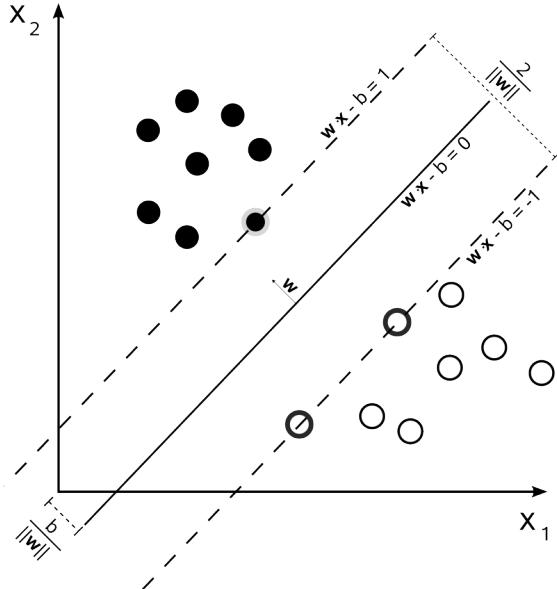


FIGURE 3.5: Illustration of the separating hyperplane of a SVM. Here  $w$  is the normal vector for the separating hyperplane and the distance between the two margins is  $\frac{2}{\|w\|}$ .  
Image from Wikipedia.org

geometric sense,  $w$  is the normal vector indicating the direction of the hyperplane and  $\frac{b}{\|w\|}$  determines the offset of the hyperplane to the origin. Since the distance between the two margins is equal to  $\frac{2}{\|w\|}$ , the maximum-margin hyperplane is found by minimizing  $\|w\|$ . The data points which lie on the margin are the *support vectors*. This geometrical interpretation is illustrated in Figure 3.5. All data points for which  $y_i = -1$  are on one side of the hyperplane and all other data points (for which  $y_i = 1$ ) are on the other side. The minimal distance from a data point to the hyperplane is for both classes equal. This results in a *maximal margin* between the two classes. Thus, the SVM searches for a maximal separating hyperplane.

With every classification method there is a risk of overfitting. In that case the random error or noise of the data set is described instead of the underlying data. The SVM classifier can use a *soft margin* by allowing some data points to lie within the margin, instead of on the margin or farther away from the hyperplane. For this it introduces *slack variables*  $\xi_i$  for each data point and the constant  $C > 0$  determines the trade-off between maximizing the margin and the number of data points within that margin (and thus the training errors). The slack variables minimize the *sum of deviations* rather than the *number* of incorrect data points [19]. The objective function for a SVM is the following minimization function:

$$\min_{w, b, \xi_i} \frac{\|w\|^2}{2} + C \sum_{i=1}^n \xi_i \quad (3.3)$$

subject to:

$$\begin{aligned} y_i(w^T \phi(x_i) + b) &\geq 1 - \xi_i && \text{for all } i = 1, \dots, n \\ \xi_i &\geq 0 && \text{for all } i = 1, \dots, n \end{aligned} \tag{3.4}$$

This minimization problem can be solved (using Quadratic Programming (QP)) and transformed to its (Lagrange) dual formulation. In the dual formulation the problem scales with the number of training examples  $n$  instead of the dimensionality  $d$  of the samples. Solving this problem directly in the high dimensional feature space  $\mathcal{F}$  makes it untractable. The linear approximation function corresponds to the kernel function in the dual formulation. Solving this dual formulation is equivalent to solving the primal formulation [19]. In the dual formulation the Lagrange multipliers  $a_i \geq 0$  are introduced and the decision function becomes:

$$f(x) = \operatorname{sgn}\left(\sum_{i=1}^n a_i y_i K(x, x_i) + b\right), \tag{3.5}$$

where  $K(x, x_i) = \phi(x)^T \phi(x_i)$  is the dot product of data objects in feature space  $\mathcal{F}$  (which is further discussed in Section 3.3.2). Here every data point in  $\mathcal{I}$  for which  $a_i > 0$  is weighted in the decision function and thus “supports” the classification machine: hence the name “Support Vector Machine”. Since it is shown that under certain circumstances SVMs show an equality to sparse representations [27, 60], there will often be relatively few Lagrange multipliers with a non-zero value.

Using this formulation two important properties arise [24]:

1. Searching for the maximum margin decision boundary is equivalent to searching for the support vectors; they are the training examples with non-zero Lagrange multipliers.
2. The optimization problem is entirely defined by pairwise dot products between training examples: the entries of the kernel matrix  $K$ .

An effect of the first property, combined with the equality to sparse representations, is that SVMs often have good results, even in the case of high dimensional data (also known as the curse of dimensionality) or limited training examples [19]. The second property is what enables powerful adaption for SVMs to learn non-linear decision boundaries. The workings of the kernel matrix  $K$  and the non-linear boundaries is discussed in the following section.

### 3.3.2 Kernels

In the previous section, the mapping function  $\phi(x)$  and the kernel function  $K$  were briefly mentioned. The decision function in Equation (3.5) only relies on the dot products of mapped data points in the feature space  $\mathcal{F}$  (i.e. all pairwise distances between the data points in that space). It shows [24] that as long as any function has the same result, without an explicit mapping to the higher dimension  $\mathcal{F}$ , the dot products can be substituted by the kernel function  $K$ , as introduced by Aizerman *et al.* [2] and applied to Support Vector Classifiers (SVCs) by Vapnik [72]:

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{z} \cdot \mathbf{z}') = \phi(\mathbf{x}) \cdot \phi(\mathbf{x}') \quad (3.6)$$

where vectors  $\mathbf{z}$  and  $\mathbf{z}'$  are projections of data objects  $\mathbf{x}$  and  $\mathbf{x}'$  through  $\phi(\mathbf{x})$  on the features space  $\mathcal{F}$  and the dot product kernel  $K$  is determined by the sum

$$K(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^m g_j(\mathbf{x})g_j(\mathbf{x}') \quad (3.7)$$

of basis functions  $g_j(\mathbf{x})$ . Note that the evaluation of dot products in the feature space  $\mathcal{F}$  between vectors is performed indirectly via the evaluation of the kernel  $K$  between (support) vectors in the input space  $\mathcal{I}$ . This is known as the *kernel trick* and gives the SVM the ability to create non-linear decision function without high computational complexity.

The kernel function  $K$  can have different forms, such as linear, polynomial and sigmoidal but the most used (and flexible) form is the Gaussian Radial Base Function (RBF). The basis functions have the form

$$g(x) = \text{sign} \left( \sum_{i=1}^n \alpha_i \exp \left\{ \frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{\sigma^2} \right\} \right), \quad (3.8)$$

where  $\sigma$  defines the width and  $\|\mathbf{x} - \mathbf{x}'\|^2$  is the dissimilarity measure expressed as Euclidian distance. The inner product kernel  $K$  then becomes

$$K(\mathbf{x}, \mathbf{x}') = \exp \left\{ -\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\sigma^2} \right\}. \quad (3.9)$$

The kernel  $K$  maps input space  $\mathcal{I}$  to the feature space  $\mathcal{F}$  which is a Reproducing Kernel Hilbert Space (RKHS) of (theoretically) infinite dimensions. As Smola *et al.* [60] state, this Gaussian kernel yields good performance, especially when no assumptions can be made about the data. As an explanation, they show a correspondence between learning SVMs with RBF kernels and good regularization operators. This may give insights in

why SVMs have been found to exhibit high generalization ability (by learning with few training objects).

The number of basis functions, the center parameters that correspond with support vectors, and the weights in the output layer are all automatically determined via the optimal hyperplane in features space  $\mathcal{F}$  [19]. The width parameter  $\sigma$  is equal for all basis functions and is set a priori and determines the flexibility and complexity of the boundary. In Section 3.3.5 this (hyper)parameter for a SVM is further discussed.

The mapping from input space  $\mathcal{I}$  to  $\mathcal{F}$  via  $\phi(x)$  is subject to the continuity assumptions. This general assumption in pattern recognition state that two near objects in feature space should also resemble each other in real life [64]. Thus, objects which are close in feature space should be close in the original input space. When this assumption would not hold, the example objects would be scatter through the feature space and finding a decision function becomes very complex.

### 3.3.3 $\nu$ -Support Vector Machine

The first of the OC-SVM methods we will discuss is often referred to as  $\nu$ -SVM and introduced by Schölkopf *et al.* [59]. Instead of estimating the density function of an distribution  $P$ , it focuses on an easier problem: the algorithm find regions in the input where the “probability density lives”. This results in a function such that most of the data is in the region where the function is nonzero.

The constructed decision function  $\mathcal{F}$  resembles the function discussed in Section 3.3.1. It returns the value  $+1$  in a (possibly small) region capturing most of the data points, and  $-1$  elsewhere. The method maps the data points from input space  $\mathcal{I}$  to a feature space  $\mathcal{F}$  (following classical SVMs). In that space  $\mathcal{F}$  it separates the data points with maximal margin from the origin, with a hyperplane. For a new data points  $x$ , the function value  $f(x)$  determines whether the data point is part of the distribution (i.e. the value is  $+1$ ) or a novelty (i.e. the value is  $-1$ ). The hyperplane is represented by  $g(x) = w \cdot \phi(x) + \rho = 0$  and the decision function is  $f(x) = \text{sgn}(g(x))$ . This hyperplane and the separation from the origin is illustrated in Figure 3.6.

The objective function to find the separating hyperplane is the following minimization function, which can be solved using QP:

$$\min_{w, \xi_i, \rho} \frac{\|w\|^2}{2} + \frac{1}{\nu n} \sum_{i=1}^n \xi_i - \rho \quad (3.10)$$

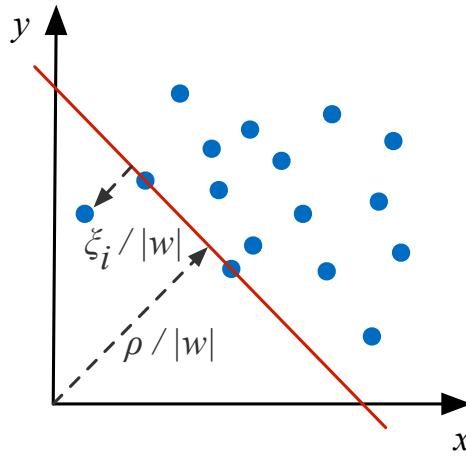


FIGURE 3.6: Graphical representation of  $\nu$ -SVM. The separating hyperplane  $w \cdot \phi(x_i) + \rho = 0$  creates a maximal margin, in the feature space, between the data points and the origin. Slack variables  $\xi_i$  are used to create a soft margin.

subject to:

$$\begin{aligned} (w \cdot \phi(x_i)) &\geq \rho - \xi_i && \text{for all } i = 1, \dots, n \\ \xi_i &\geq 0 && \text{for all } i = 1, \dots, n \end{aligned} \tag{3.11}$$

The decision function in the dual formulation with Lagrange multipliers is denoted as:

$$f(x) = \text{sgn}((w \cdot \phi(x_i)) - \rho) = \text{sgn}\left(\sum_{i=1}^n \alpha_i K(x, x_i) - \rho\right) \tag{3.12}$$

In the classical SVM objective function, as denoted in Equation (3.3), the parameter  $C$  decided the smoothness of the boundary, with respect to the slack variables  $\xi_i$ . In the formulation of  $\nu$ -SVM the equivalent parameter is  $\nu \in (0, 1)$  (hence the name). It characterizes the solution in two ways:

1.  $\nu$  is an upper bound on the fraction of outliers, i.e. training examples regarded as out-of-class.
2.  $\nu$  is a lower bound on the fraction of Support Vectors (SVs), i.e. training examples with a nonzero Lagrange multiplier  $\alpha_i$ .

When  $\nu$  approaches 0, the penalty factor for nonzero Lagrange multipliers ( $\frac{1}{\nu n}$ ) becomes infinite, and thus the solution resembles a *hard margin* boundary.

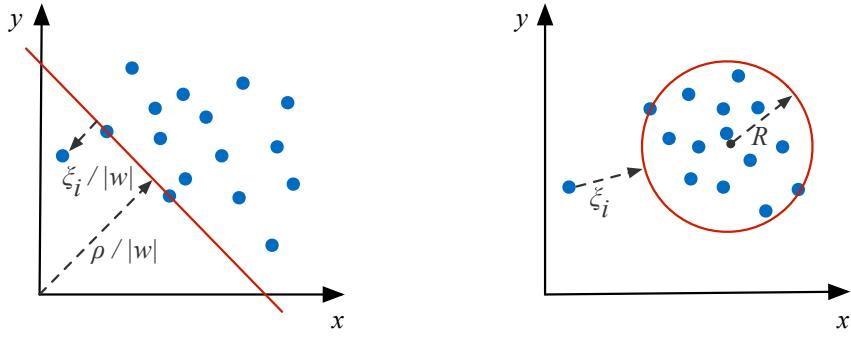


FIGURE 3.7: Graphical representation of the difference between  $\nu$ -SVM (left) and SVDD (right). Note that for the sake of simplicity the kernel functions are not applied.

This method creates a *hyperplane*, characterized by  $w$  and  $\rho$ , that separates the data with maximal margin from the origin in the feature space  $\mathcal{F}$ . In the following section we will discuss an alternative method, which uses an circumscribing *hypersphere* to characterize the data (and the region of distribution  $P$  where the density (or: support) lives).

### 3.3.4 Support Vector Data Description

The method introduced by Tax and Duin [65], known as Support Vector Data Description, follows a spherical instead of planar approach. The boundary, created in feature space  $\mathcal{F}$ , forms a hypersphere around the (high density region of the) data. The volume of this hypersphere is minimized to get the smallest enclosing boundary. The chance of accepting outlier objects is thereby also minimized [68]. By allowing outliers using slacks variables, in the same manner as classical SVM and  $\nu$ -SVM, a soft margin is constructed.

The constructed hypersphere is characterized by a center  $\mathbf{a}$  and a radius  $R > 0$  as distance from the center to (any data point that is a SV on) the boundary, for which the volume, and thus the radius  $R$ , will be minimized. The center  $\mathbf{a}$  is a linear combination of the support vectors. Like the classical SVM and SVDD it can be required that all the distances from the data points  $x_i$  to the center  $\mathbf{a}$  are strict less than  $R$  (or equivalent measure). A soft margin can be allowed by using slack variables  $\xi_i$ . In that case, the penalty is determined by  $C$  and the minimization is expressed as Equation (3.13). This principle is illustrated in the right image or Figure 3.7. Instead of a separating hyperplane, constructed by  $\nu$ -SVM and illustrated on the left of the Figure, the SVDD creates a hypersphere (in the illustration a circle) around the data points. By using kernel functions (e.g. the RBF) the hyperspheres in the high dimensional feature space  $\mathcal{F}$  corresponds to a flexible and tight enclosing boundary in input space  $\mathcal{I}$ . Possible resulting closed boundaries are illustrated in Figure 3.8. This enclosing boundary is

obtained by minimizing the following error function  $L$  which contains the volume of the hypersphere and the distance from the boundary to the outlier objects:

$$L(R, \mathbf{a}, \boldsymbol{\xi}) = R^2 + C \sum_{i=1}^n \xi_i \quad (3.13)$$

subject to:

$$\begin{aligned} \|x_i - \mathbf{a}\|^2 &\leq R^2 + \xi_i && \text{for all } i = 1, \dots, n \\ \xi_i &\geq 0 && \text{for all } i = 1, \dots, n \end{aligned} \quad (3.14)$$

In the dual Lagrangian formulation of this error function  $L$  the multipliers  $\boldsymbol{\alpha}$  are maximized:

$$L = \sum_i \alpha_i (x_i \cdot x_i) - \sum_{i,j} \alpha_i \alpha_j (x_i \cdot x_j) \quad (3.15)$$

subject to:

$$0 \leq \alpha_i \leq C, \sum_i \alpha_i = 1 \quad (3.16)$$

In the maximization of Equation (3.15) a large fraction of the multipliers  $\alpha_i$  become zero and for a small fraction  $\alpha_i > 0$ . This small fraction, for which  $\alpha_i$  is non-zero, are called the SVs and these objects lie on the boundary of the description. The center of the hypersphere only depends on this small number of SVs and the objects for which  $\alpha_i = 0$  can be disregarded from the solution. Testing the membership of a (new) object  $\mathbf{z}$  is done by determining if the distance to the center  $\mathbf{a}$  of the sphere is equal or smaller to the radius  $R$ :

$$\|\mathbf{z} - \mathbf{a}\|^2 = (\mathbf{z} \cdot \mathbf{z}) - 2 \sum_i \alpha_i (\mathbf{z} \cdot \mathbf{x}_i) + \sum_{i,j} (\mathbf{x}_i \cdot \mathbf{x}_j) \leq R^2 \quad (3.17)$$

As with Equation (3.5), the solution of this equation only relies on dot products between the data points in  $\mathbf{x}$  and  $\mathbf{z}$ . This means that the kernel projection and trick, as discussed in Section 3.3.2, can be applied to SVDD as well [65, 66].

Because the Gaussian RBF often yields good (i.e. tight) boundaries, this set of kernels fuctions is commonly used:

$$(\mathbf{x} \cdot \mathbf{y}) \rightarrow K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{\sigma^2}\right) \quad (3.18)$$

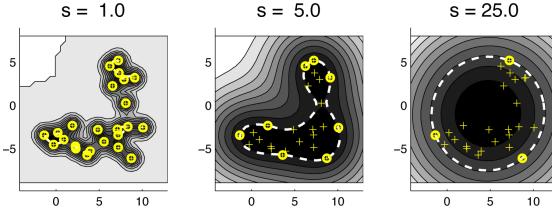


FIGURE 3.8: The SVDD method trained on a banana-shaped data set with different sigma-values for the RBF kernel. Solid circles are support vectors, the dashed line is the boundary. Image by Tax [64].

Using this kernel function, the Lagrangian error function  $L$  of Equation (3.15) changes to:

$$L = 1 - \sum_i \alpha_i^2 - \sum_{i \neq j} \alpha_i \alpha_j K(x_i, x_j) \quad (3.19)$$

Using Equation (3.17), the following kernel formulation needs to hold for a new object  $\mathbf{z}$  to lie within the hypersphere:

$$\sum_i \alpha_i K(\mathbf{z}, x_i) \leq \frac{1}{2} \left( 1 - R + \sum_{i,j} \alpha_i \alpha_j K(x_i, x_j) \right) \quad (3.20)$$

When the Gaussian RBF kernel is applied, and in case the data is preprocessed to have unit length (for the  $\nu$ -SVM solution), the two different OC-SVM implementations  $\nu$ -SVM and SVDD are shown to have identical solutions [58, 66]

### 3.3.5 SVM model parameters

SVM-model selectiong and tuning depends on two type of parameters [19]:

1. Parameters controlling the ‘margin’ size,
2. Model parameterization, e.g. the kernel type and complexity parameters. For the RBF kernel the width parameter determines the model complexity.

In case of a RBF kernel, the width parameter  $\sigma$  determines the flexibility and complexity of the boundary. The value of this parameter greatly determines the outcomes of the algorithm (e.g. SVDD) as illustrated in Figure 3.8. With a small value for the kernel width  $\sigma$ , each data point will tend to be used as a support vector (for almost all  $\alpha_i > 0$ ) and the SVDD solution resembles a Parzen density estimation. For large values of  $\sigma$ , the solution will resemble the original hypersphere solution is obtained (in contrast with a tight boundary around the data). With a large value for the width  $\sigma$ , the boudary approximates the spherical boundary. The influence of the  $\sigma$  parameter on the SVDD solution is illustrated in Figure 3.9.

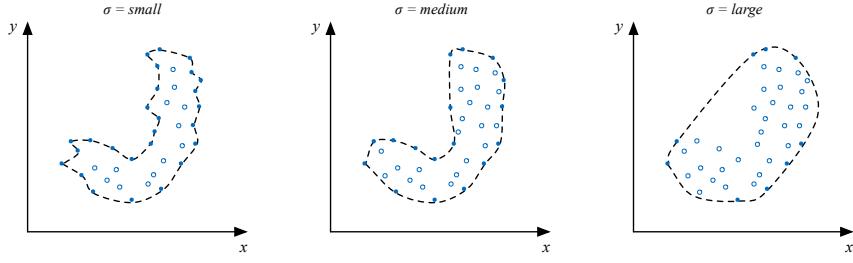


FIGURE 3.9: The SVDD method trained on a banana-shaped data set with different  $\sigma$ -values for the RBF kernel. Solid circles are support vectors, the dashed line is the boundary.

As discussed in Section 3.3.3, the SVM parameter  $C$  (or  $\nu$  in case of  $\nu$ -SVM) is of high influence on the “smoothness” of the decision function. It acts as an upper bound to the fraction of outliers and as a lower bound to the fraction of SVs. A more detailed discussion of the influence of the SVM model parameters can be found in Section 9.8 of [19] and Section 7.4 from [24]. A detailed discussion of the  $\nu$  and kernel parameters can be found in [58].

The following chapter will discuss our proposed method, which incorporates the SVDD algorithm. It relates the OC-SVM model construction to outlier detection and eventually change detection, leading to finding a temporal segmentation of time series data.

## Chapter 4

# Application to Human Activity Time Series Data

\*\*\* *TODO: Decide on name* \*\*\*

- Human Activity Support Vector Change Detection
- Inertial Sensor Time Series Svm Temporal Segmentation
- Support Vector bases Human Activity Temporal Segmentation
- Human Activity Temporal Segmentation by Support Vector Machines: HATS-SVM
- HATS with One-Class SVM: HATS-OCS, HATS-OS, HATSOS
- OCS-HATS
- HATS-SVDD

This chapter will introduce our method and setup for the temporal segmentation of human activities, using an OC-SVM approach. It shows our method to apply a OC-SVM based temporal segmentation method to inertial sensor time series data, often used in the context of Human Activity Recognition. To our best knowledge it is the first in its kind, especially on the application of SVDD and SVCPD to real-world HAR data. The setup follows the structure as illustrated in Figure 4.1 and further discussed in Section 4.1. It starts with the data collection phase, in which we both use artificial and real-world data. The data sets are discussed in Section 6.2. After (optional) pre-processing, the SVDD algorithm is used to create a model representation, which is iteratively updated. In Section 4.3 this construction phase is discussed and shows the windowed approach.

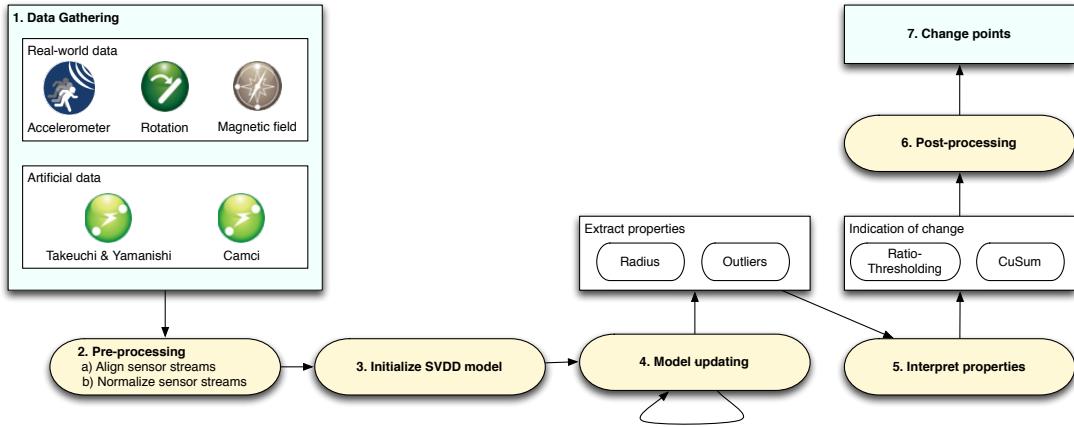


FIGURE 4.1: Schematic overview of the change detection method. The first step is the data gathering, described in Section 4.2. After the pre-processing, the data is used to construct a SVDD model, as described in Section 4.3. Section 4.4 describes which features of the model are used for the final change indication algorithm, discussed in Section 4.5.

From this constructed model we can extract some features. In our case, we use the approximated radius  $R$  of the hypersphere to detect changes. Section 4.4 goes into detail about the rationale of using that feature. Finally, in Section 2.6 we discuss two different, but simple, methods to detect changes in the time series data, based on the extracted radius  $R$ .

## 4.1 Algorithm overview

The proposed method of this thesis follows the unifying framework as introduced by Takeuchi and Yamanishi [63] and an similar implementation by Camci [14] with SVMs. The unifying framework relates the detection of outliers with change points and divides the whole process in two stages. The first stage determines the outliers in a time series by giving a score based on the deviation from a learned model, and thereby creates a new time series. The second stage runs on that new created time series and calculates a average over a window of the outlier scores. The problem of change detection is then reduced to outlier detection over that average-scored time series. This method is named changeFinder by the authors. The implementation by Camci, which uses SVMs to detect changes is named Support Vector based Change Point Detection.

Whereas changeFinder uses double probability estimation algorithm, our approach follows SVCPD by constructing a SVM over a sliding window. The SVCPD algorithm uses the location of new data points in the feature space  $\mathcal{F}$  with respect to the hypersphere and the hypersphere's radius  $R$  to determine whether the new data point represents a change point. Our approach is a slight simplification of the algorithm: we only use

the updated model properties (i.e. the radius  $R$ ) to detect a change. This difference is further discussed and justified in Section 4.5.

This section gives a description of the method used for the experiments and change detection mechanism. First described is the method to process the gathered sensor data. A schematic overview is given in Figure 4.1 and shows the steps of the method. A more detailed explanation of the “Update model” step follows.

As graphically represented in Figure 4.1, the change detection method starts by processing the data from sensor, such as the accelerometer, magnetic orientation and rotation metrics.

The first step is to process the raw streams of data originating from a multiple of sensors. The two processes applied are alignment and normalization. Due to noisy sampling, not all the timestamps in the data streams are sensed at the same timestamp. Since the SVDD method requires all the data stream at every timestamp and can not handle missing data on one of the timestamps, all the unique timestamps are filtered out. Whilst this results in an overall filtering effect, in practice between 1% and 5% of each data stream is disregarded. The effect of this filtering is not significant and the data is not modified.

Due to the nature of the sensor signals, a normalization step is required in order to set the weight for all the data streams equal. The range of the accelerometer signal typically spans  $-20$  to  $20$ , the magnetic field from  $-60$  to  $60$  and the rotations range is from  $-1$  to  $1$ . This means that a relative small change in the accelerometer stream could have a much larger impact on the model than the same (absolute) change in the rotation stream, whilst the latter has a larger relative impact. The normalization step ensures that all data is weighted equally and changes in the data are all proportional.

In step 3 the SVDD model is initialized. The first full window over the data stream is used to construct an initial model. During the initialization the parameters for the SVDD are provided, begin the kernel type (radial), RBF with  $\sigma$  and the outlier-fraction  $C$ .

Step 4 is executed for every step-size  $s$  data points in the stream. Every update the oldest  $s$  data points are removed from and  $s$  new data points are added to the SVDD model. The Incremental Support Vector Data Description (I-SVDD) algorithm by Tax and Duin [66] is used for efficient incremental model updating. The model is (partially) reconstructed and new model properties, such as the radius of the hypersphere, is the result of this step <sup>1</sup>.

---

<sup>1</sup>Other measures are also possible, for instance the distance from all the outliers to the boundary of the hypersphere or the number of outliers

This final step of this method, step 5, is the interpretation of the model properties. Many algorithms can be used for this process, all which take a one-dimensional time series as input and determine where change has occurred. In our setup we used the Ratio-Thresholding (RT) and CUSUM methods, to show the modularity of this step. This final step is further discussed in Section 2.6.

## 4.2 Data Gathering

In this section we briefly discuss the different data gathering methods used for the change detection algorithms and experiments. Section 4.2.1 discusses the artificial data sets we will use. In Section 4.2.2 an overview of the real-world data sets used is provided. Both sections refer to Chapters 5 and 6 for more details, respectively.

### 4.2.1 Artificial data

In order to provide an objective comparison to other methods, we will use artificial data sets which are also used in the researches on which our method is based. These are the data sets used by Takeuchi and Yamanishi [63] and Camci [14]. Both construct a collection of one-dimensional time series data according to a second order Autoregressive (AR) model:

$$x_t = a_1 x_{t-1} + a_2 x_{t-2} + \epsilon_t. \quad (4.1)$$

Over the different data series the mean and variance of the Gaussian random variable  $\epsilon_t$  differs and changes at pre-determined change points. Using this data set an objective quality measure over the change detection methods can be obtained and compared. All the used data sets are listed and analyzed in Chapter 5.

### 4.2.2 Real-world data

In the second type of data sets we apply our method of change detection and temporal segmentation to real-world data sets. For our setup we record the activities of humans performed both in- and outdoor in an uncontrolled environment. Activities performed include sitting, standing, walking, running in a straight and curved line, and walking up- and downstairs. Our method uses the signals from the accelerometer, magnetic field, and rotation sensors. These time series data are used to detect change points. A video recording from the performed activity is used to annotate the time series with real change points. The discovered change point are compared with these annotated change

Metric	Description	Units of measure	Typical range
Accelerometer	Acceleration force along each axis (including gravity).	$m/s^2$	-20 – 20
Gravity	Force of gravity along each axis.	$m/s^2$	-10 – 10
Gyroscope	Rate of rotationg around each axis.	$rad/s$	-15 – 15
Light	Light sensitive sensor at the front of the phone.		0 – 10000
Linear acceleration	Acceleration force along each axis (excluding gravity).	$m/s^2$	-20 – 20
Magnetic field	Geomagnetic field strength along each axis.	$\mu T$	-60 – 60
Orientation	Degrees of rotation around the three physical axis.	Degrees	-100 – 360
Rotation	Measre of rotation around the device's rotation axis.	Unitless	-1 – 1

TABLE 4.1: Measured metrics. The set of axis is always the triple (x, y, z) direction.

points to give a subjective quality measure. In Chapter 6 we give a detailed analysis of the performed activities and the recorded data sets.

For the experiments we used a smartphone with inertial sensors as recording device. The activities were recorded using a free ANDROID application [45]. This application was chosen for its convenient data format of the sensor recording and its regularity of the sampling interval. Table 4.1 lists all the recorded metrics. For our expiments we used the data for the accelerometer, magnatief field and rotation.

We implemented our algorithm in MATLAB. For the SVDD model construction phase we used the DD\_TOOLS library by Tax [69], which is dependend on the PR\_TOOLS package, originating from [70]. In an earlier attempt we used the a LIBSVM [16] add-on for the SVDD algorithm [17], but we had troubles making that implementation work with our data.

### 4.3 Model Construction: Incremental SVDD

After the data is collected and pre-processed (or in the case of the artificial data sets: generated), we construct an online incremental sliding window algorithm. We follow the

SVCPD method by Camci [14] and the I-SVDD implementation by Tax and Laskov [68]. The latter algorithm combines the techniques of online, unsupervised, and incremental learning methods with the earlier introduced OC-SVM algorithm SVDD. The method is initialized with a fixed window length and then in every step a new data object is added to and the oldest data object is removed from the working set. The model construction phase can be regarded as the first stage of the unifying framework by Takeuchi and Yamanishi [63].

Using the following abstract form of the SVM optimization problem, the extension of the incremental SVM to the SVDD can be carried out:

$$\max_{\mu} \min_{\substack{0 \leq x \leq C \\ \mathbf{a}^T \mathbf{x} + b = 0}} : W = -\mathbf{c}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T K \mathbf{x} + \mu(\mathbf{a}^T \mathbf{x} + b), \quad (4.2)$$

where  $\mathbf{c}$  and  $\mathbf{a}$  are  $n \times 1$  vectors,  $K$  is a  $n \times n$  matrix and  $b$  is a scalar. The SVDD implementation of this abstract form is set by the parameters  $\mathbf{c} = \text{diag}(K)$ ,  $\mathbf{a} = \mathbf{y}$  and  $b = 1$ . The procedure for the incremental version has two operations for a data object  $k$ : add to and remove from the working set. When a data object  $k$  added, its weight  $x_k$  is initially set to 0. In case of an object removal, the weight is forced to be  $x_k = 0$ . Both the operations conclude with the recalculation of  $\mu$  and the weights  $\mathbf{x}$  for all the objects, in order to obtain the optimal solution for the enlarged or reduced data set. These two operations are the basis for I-SVDD. With respect to the timed index, new data objects are added to and old data objects are removed from the working set.

The size of the initial window of data objects has a lower bound determined by the hyperparameter  $C$  (Equation (3.13)). Because of the equality constraint  $\sum_{i=1}^n a_i x_i = 1$  and the box constraint  $0 \leq x_i \leq C$ , the number of objects in the working set must be at least  $\lceil \frac{1}{C} \rceil$ . Thus the algorithm is initialized by selection the first  $\lceil \frac{1}{C} \rceil$  objects for the working set. In every loop of the algorithm the number of elements added are at least enough to keep  $\lceil \frac{1}{C} \rceil$  objects in the window, after removing the oldest objects. The optimality of the algorithm is proved by analysis of the Karush-Kuhn-Tucker (KKT) conditions [68].

From experiments it shows that the (online) I-SVDD method results in less false alarms than the static SVDD. An explanation for this is that I-SVDD follows the changing data distribution, i.e. small changes over time. Such changes are a drift in mean or increase in frequency. The I-SVDD algorithm continuously re-models the SVM representation. This re-modeling of the SVM representation is beneficial for our purpose, since we are only interested in sudden changes between activities. The following section will discuss the features extracted from the constructed OC-SVM model at every loop-step of the algorithm, and Section 4.5 describes how these features are used to detect change.

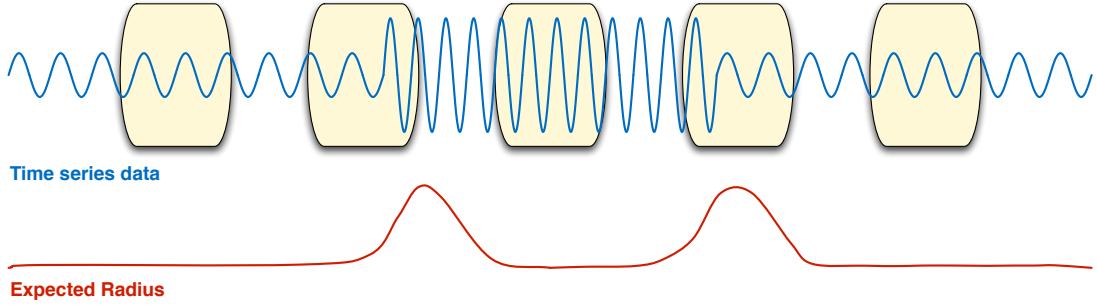


FIGURE 4.2: Expected behavior of the radius  $R$  of the hypersphere. The upper part shows a typical sinusoidal time series signal. The lower graph visualizes an abstract expectation of the values of  $R$ . Five possible time windows are illustrated. The left, middle, and right windows cover an area of homogeneous signal. The expected value of  $R$  is low. The other two windows cover a change entering and leaving the window, respectively. At these locations the radius  $R$  is expected to increase temporarily.

## 4.4 Model Features

In the previous section we have discussed the I-SVDD method, which creates a OC-SVM representation of a working set of data objects at every step of the algorithms loop. This section shows how we interpret the constructed model and extract features to obtain a measure which can be used for an indication of change points. The next section discusses how this obtained measure is used to indicate change.

The I-SVDD algorithm creates a spherical OC-SVM representation of the working set at every step of the algorithm. This model is obtained by the minimization of Equation (3.13), which incorporates the radius  $R$  of the sphere and the distances  $\xi$  from the outliers to the boundary. We will use the radius  $R$  of the hypersphere as an indication of change.

In [66] Tax and Duin provide an analysis of the error of the SVDD algorithm. This error is based on 1) the fraction  $f_{T-}$  of target objects that is rejected, and 2) the fraction  $f_{O+}$  of outliers that is accepted. Since in OCC situations typically there are (almost) no examples of outlier objects, Tax and Duin construct a method to generate outliers based on the assumption that the outliers are uniformly distributed around the target set. To minimize the error, calculated by the fractions  $f_{T-}$  and  $f_{O+}$ , we should minimize the volume of the target data description (i.e. the boundary of SVDD). That is because the fraction of accepted outliers  $f_{O+}$  is an estimate of the volume of the target data description, with respect to the volume of the outlier distribution. Tax and Duin provide a method to optimize the parameters of the SVDD method, i.e. the trade-off parameter  $C$  and the RBF kernel width  $\sigma$ . This optimization will result in the modification of the radius  $R$  of Equation (3.13) and affects the Lagrangian inequality (3.20).

Since in our method the parameters  $C$  and  $\sigma$  are kept constant, and thereby also the fraction of target objects being rejected, the only free parameter is the radius  $R$ . During the I-SVDD algorithm the volume of the hypersphere is being minimized. This means that only the radius of the sphere will change during the analysis over the time series data. We will interpret changes in the radius  $R$  with the same continuity assumptions as mentioned in Section 3.3.2. This means that when two objects (in this case radius sizes) are close to each other, their original objects are also close. In other words, when the underlying distribution of time series is changed, and thus the characteristics of the data objects changes, the corresponding radius  $R$  will also change. Since in that case the data objects in the working set become more heterogeneous, the radius of the hypersphere will increase. When, instead, the data objects change from a heterogeneous set to a more homogeneous, we expect the radius to decrease in value, since the data objects are closer to each other in the feature space. This relation between the signal and expected radius  $R$  is illustrated in Figure 4.2. It shows that in the two time windows that cover a heterogeneous set of data objects, the radius of the hypersphere is expected to be relatively high.

With the I-SVDD algorithm we have effectively implemented a form of dimensionality reduction by feature extraction, using the radius  $R$ . The following section discusses the algorithms which can be applied to the extracted radius  $R$  as a volume estimate. We thereby follow the setup of the unifying framework by Takeuchi and Yamanishi [63], of which this section described the first stage.

## 4.5 Change Detection

In the previous sections we have discussed what data is used for the model construction, and how that model is interpreted in the context of change detection. We have shown how the multi-dimensional signal is reduced to a single dimension time series data, based on the extracted radius  $R$  from the constructed hypersphere. This metric is used to discover changes in the underlying data generation process. We argued that in the case of the data objects set becomes more heterogeneous, we expect an increase in the radius  $R$ . In the same manner we expect that when the change is in the oldest data objects of the working set the radius  $R$  will decrease, since the data becomes more homogeneous. An abstract representation of these expectations is illustrated in Figure 4.2.

In this section we discuss the methods applicable to the extracted radius  $R$  of the hypersphere, in order to discover a change in the underlying distribution properties. This employs the second stage of our algorithm, analogue to the second stage of the

Step	Action
1	Start with $n$ data objects as working set and construct hypersphere
2	Add next data object $x_t$ and drop first one
3	Identify new hyper-sphere and its <i>approximate radius</i>
4	calculate radius average of hyperspheres since last change point
5	calculate radius ratio $\bar{h}$ . If lower than $th_{low}$ or greater than $th_{high}$ then mark $t$ as change point
6	continue from step 2, until no more data objects available
7	Apply post-processing: merge close change points

TABLE 4.2: Algorithm of the proposed change detection method, using the Ratio-Thresholding method of 4.5.2 for change detection.

unifying framework by Takeuchi and Yamanishi [63]. It is also based on the SVCPD method by Camci [14], although we have slightly simplified the algorithm.

In the SVCPD the first step for a new data point  $z$  is to check whether the data point lies within the enclosing (hyperspherical) boundary of the constructed model. This will label the new data point as an in- or outlier, relative to the current data set. Depending on that outcome, the model is updated to include the new data point. The second step is to inspect the (changed) model properties, i.e. the radius  $R$ . In our method we only employ the second step: the inspection of the radius  $R$ . We regard the first step (determine if  $z$  is an outlier) to be superfluous: in case it is an outlier, the model update step will change the radius  $R$  in order to let the boundary include the new data point  $z$ . Would the new data object  $z$  be already an inlier, then the updated model would not have a significant change in radius  $R$ . Our algorithm is outlined in Table 4.2.

The problem of change detection in complex multi-dimensional time series is now reduced to finding change in a one-dimensional time series. Besides the dimensionality reduction, the form of the signal is also simplified. Whereas the original signal was represented as a sinusoidal or second order AR model, the new time series is of much simpler form. In our method we leave open the precise method which is used to find change points in the obtained change indication values. Below, we will give two examples of methods which can be applied. The first is a family of methods based on the CUSUM method. The second example is a ratio-based thresholding mechanism used by Camci [14]. We conclude this section by discussing some post-processing techniques we have used, to decrease the number of false positives.

### 4.5.1 CUSUM based methods

Many of the simple change detection methods are based on CUSUM, originally introduced by Page [56]. Other variations and extensions of the method are proposed and used ([4, 38, 39]). Here we will discuss the simple form, which only detects changes as an increase in value of the examples, in this case the radius  $R$ . The cumulative sum for all the values is calculated:

$$S_n = \sum_{k=1}^n r_k \quad (4.3)$$

and a change is detected when the value of the cumulative sums minus the minimum encountered exceeds a predetermined threshold  $h$ :

$$S_n - \min_{0 \leq i < n} S_i \geq h. \quad (4.4)$$

Extensions of the CUSUM method have been proposed. Amongst others, extensions which extend its use to two-threshold methods, and methods appropriate for changes in mean or variance [39]. Whilst adaptive methods have been proposed [4], most methods require manual determination of the threshold parameter  $h$ .

### 4.5.2 Ratio-thresholding

The second example of an algorithm that interprets the change indication time series and transforms it to change detection is the radius ratio thresholding method used by Camci [14]. The method shares characteristics with CUSUM, since it relies on summation of historic data and thresholds for change detection. The method calculates the radius ratio  $\hbar_t$  by taking the ratio of the current radius  $r_t$  at time  $t$  to the average radii since the last change point  $y$ :

$$\hbar_t = \frac{r_t}{\text{mean}(r_{y:t-1})}, \quad (4.5)$$

where  $\text{mean}(r_{y:t-1})$  is the average of the previous approximated radii. By using only historic data, this method can be incorporated in an online change detection method. The ratio  $\hbar_t$  is compared with two thresholds: the low threshold  $th_{\text{low}}$  and the high threshold  $th_{\text{high}}$ . When the data objects in the working set become more homogeneous at time  $t$ , the radius  $R$  of the hypersphere and thereby the ratio  $\hbar_t$  will decrease. When the new value of  $\hbar_r$  is lower than the low threshold  $th_{\text{low}}$ , a change detection is observed for time  $t$ . The same holds for a more heterogeneous set of data and an increase of  $\hbar_r$  higher than  $th_{\text{high}}$ . This is illustrated in Figure 4.3.

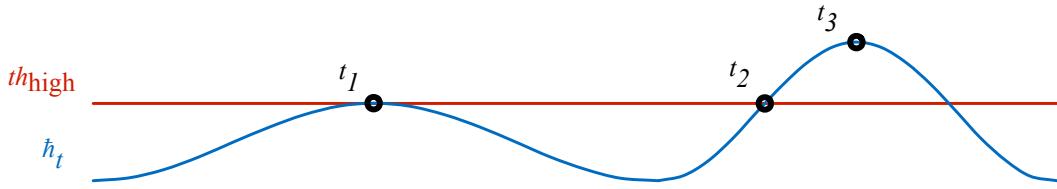


FIGURE 4.3: Abstract illustration of the radius ratio thresholding. The peak at  $t_1$  generates a single change point. The region between  $t_2$  and  $t_3$  generates multiple change points, since the value of  $h_t$  keeps increasing.

The appropriate values for the thresholds are highly depended on the characteristics of the data and the trade-off parameter  $C$  (as discussed in Section 3.3.4). This trade-off parameters regulates the fraction of data objects that will be rejected from the constructed model. Since a high value of  $C$  assigns a high penalty to outliers, data objects are more likely to be incorporated into the hypersphere by increasing the radius. For low values of  $C$  the cost of rejecting data objects is, compared to the benefits of a smaller hypersphere, relatively low. This gives a relation between the parameters  $C$ ,  $th_{low}$ ,  $th_{high}$ , and the sensitivity of the radius ratio based thresholding method. A high value of  $C$ , or values close to 1 for  $th_{low}$  and  $th_{high}$  result in a sensitive change detection procedure. The reverse results in less sensitive methods. The correct values for the intended sensitivity need to be empirically determined.

#### 4.5.3 Post-processing

In our method we use the ratio based thresholding, since it is also used by Camci [14]. From preliminary experiments we observed that a single change in the data can cause many detected change points, using the ratio based method. This phenomenon can be explained by the nature of the thresholding method. Since a single (high or low) threshold is set for ratio of the radii, high values are considered to be a change point before the highest point of the peak. As illustrated in Figure 4.3, the peak at  $t_1$  represents a single change point. The region from  $t_2$  to  $t_3$  is strictly increasing. At  $t_2$  a change point is detected, since  $h_t$  exceeds the threshold. But since it is increasing up to  $t_3$ , all values of  $h_t$  in that period will trigger the ratio based method to indicate a change point.

To overcome this problem, we apply a post-processing method on the generated change points. All the change points that are within a time period  $\delta$  of each other are merged together.

In the following two chapters we will take a detailed look at the used data sets and apply the method as described in this method to those sets.

# Chapter 5

## Artificial data results

In this chapter we will apply the proposed algorithm as stated in the previous chapter to artificial data sets. For comparison, we use some of the data sets that are used by Camci [14] and Takeuchi and Yamanishi [63]. These data sets are generated by Gaussian noise distributions. In Section 5.1 the precise generation method and characteristics of the data sets are discussed. For our final change detection mechanism we use the RT-method, also applied by Camci [14]. This is a simple and fast method and allows us to mainly focus on the OC-SVM model properties that should reflect changes in the data. The results of these two methods are discussed in Section 5.3.

### 5.1 Artificial data generation

In this section we present four data sets as used by [14, 63] to provide for a objective performance comparison. All the data sets are modeled by a second-order AR-model as in Equation (5.1), where  $\epsilon_t$  is a Gaussian distribution modeling the noise with mean 0 and variance  $\sigma^2 = 1$ ,  $a_1 = 0.6$  and  $a_2 = -0.5$ .

$$x_t = a_1 x_{t-1} + a_2 x_{t-2} + \epsilon_t, \quad (5.1)$$

The length of  $x_t$  is 10000 and change points are generated at each  $y \times 10000^{\text{th}}$  data point, with  $y = (1, 2, 3, \dots, 9)$ .

Using this general AR-model, we create four different data sets, all visualized in Figure 5.1. Each data set is characterized by a change in mean, variance, or both of the Gaussian term  $\epsilon_t$ :

1. **Fixed increasing mean:** at each change point  $y$  the mean is increased with  $\Delta(y) = 5$ . This set originates as the first data set of Camci [14]. The main characteristic is that the difference between the means of the Gaussian term  $\epsilon_t$  is fixed and thus the relative distance becomes smaller. The data is plotted in Figure 5.1a.
2. **Reduced increasing mean:** at each change point  $y$  the mean is increased with  $\Delta(y) = 10 - y$ . This set originates as the first data set of Takeuchi and Yamanishi [63]. In contrast with the previous set, in this set the absolute difference in mean is reduced over time. The data is plotted in Figure 5.1b.
3. **Reduced increasing mean, increasing variance:** at each change point  $y$  the mean is increased with  $\Delta(y) = 10 - y$  and the standard deviation is of  $\epsilon_t$  is  $0.1/(0.01 + (10000 - x)/1000)$ . This set originates as the second set of Camci [14]. This set combines the reduced increasing mean property of the previous set with an increasing variance, causing the change points to become more vague over time. The data is plotted in Figure 5.1c.
4. **Alternating variance:** at each change point  $y$  the variance  $\sigma^2$  is set to 9.0 when  $y$  is odd, and to 1.0 otherwise. This set originates as the third set from Camci [14] and Takeuchi and Yamanishi [63]. This data set is used to measure the ability of detecting changes in decrease of variance, as Camci [14] claims is not trivial. The data is plotted in Figure 5.1d.

## 5.2 Quality metrics

For an objective measure of quality of the proposed method applied to the four aforementioned data sets, we employ two metrics, inspired by and slightly altered from Takeuchi and Yamanishi [63]. The first metric is the average\_delay, for which we first need to define the notion of delay for each change:

$$\text{delay}(y) = |t_y - t_y^*|, \quad (5.2)$$

where  $t_y$  is the time of the closest discovered change point and  $t_y^*$  is the true change time. We can then define the average\_delay over all the discovered change points  $Y$  as:

$$\text{average\_delay}(Y) = \frac{\sum_{y=1}^Y \text{delay}(y)}{Y}. \quad (5.3)$$

Besides the delay, which should be low to indicate discovered change points near the actual change point, we need to penalize discovered change points which do not relate

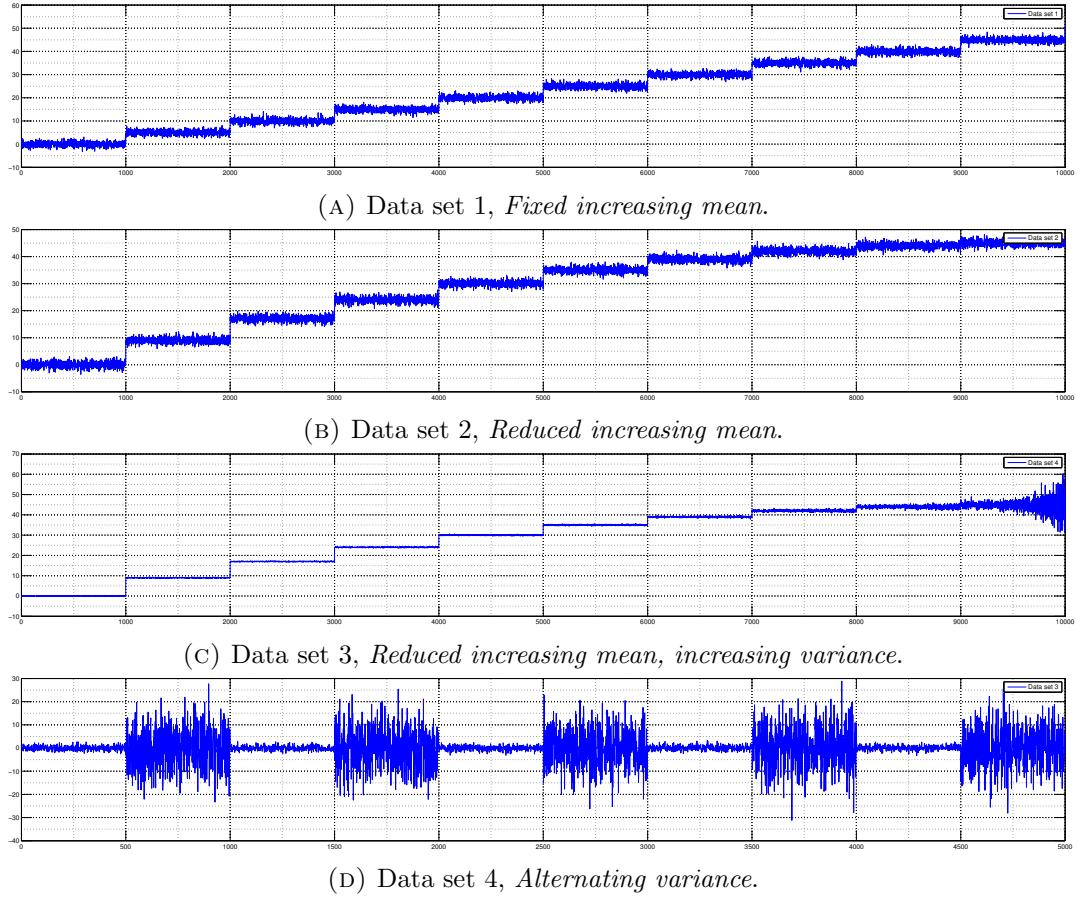


FIGURE 5.1: Plots of the artificial data sets used.

to a true change point. For this we use the False Alarm Rate, expressed as

$$\text{far}(Y) = \frac{Y_{\text{false}}}{Y}. \quad (5.4)$$

For a high quality method the average\_delay and the  $\text{far}(Y)$  should both be low. In our results we will also list the parameters used, since these need to be user-defined and are part of the (manual) optimization process. The parameters are:

- **Window length:** The window length that is processed by the SVDD algorithm. This is the number of data points that are continuously modeled in the shape in a hypersphere.
- **Sigma of RBF:** The value of  $\sigma$  for the RBF kernel. The effect of this parameter is discussed in Section 3.3.5.
- **High threshold:** The  $th_{\text{high}}$  threshold value for the RT method described in Section 4.5.2. A higher value indicates a more precise and robust setting.

	Set 1	Set 2	Set 3	Set 4
Window length	50	100	50	50
Sigma of RBF	13	13	15	13
High threshold	1.6	1.6	1.5	2.2
Low threshold	0.1	0.1	0.5	0.1
Closeness	10	10	10	50
far( $Y$ )	0	0.1	0.1	0.1
Average_delay	2.27	6.18	0.64	15.36
Standard deviation Delay	1.19	7.94	0.92	17.70

TABLE 5.1: Parameter settings and results of the artificial data sets.

- **Low threshold:** The  $th_{low}$  threshold value for the RT method described in Section 4.5.2. A lower value indicates a more precise and robust setting.
- **Closeness:** In the post-processing phase the detected change points which are within a certain *closeness* range will be merged together. This parameter indicates the minimal number of data points that change points must differ in order to be considered separate detections.

For all the sets we have used a value of  $C = 0.1$ , which indicated that 10% of the data points are used as outliers. This is a normal setting for the SVDD method.

### 5.3 Results

The results of all the data sets are provided in Table 5.1, together with the used parameters. For each data set it shows the high and low thresholds used in the RT based change detection method. During the post-processing all detected change points within a certain *closeness* range are merged, for which the value differs per data set. The difference between the actual and final detected change points are expressed as an offset. The data sets were processed with different window lengths and  $\sigma$  value for the RBF kernel.

To give a better impression of the spread of the accuracy of the detected change points, a box plot for each data set is displayed in Figure 5.2 and are the results are visualized in Figure 5.3. This last figure shows the number of data points between each actual and closest detected change point. Using the box plot notation, we can visualize not only the average but also the spread of the delay. With each data set embodying different characteristics, different observations can be drawn from the data sets. In the references figures, the blue lines represent the radius of the constructed hypersphere. The vertical red lines are the discovered change points. For each data set the following is observed:

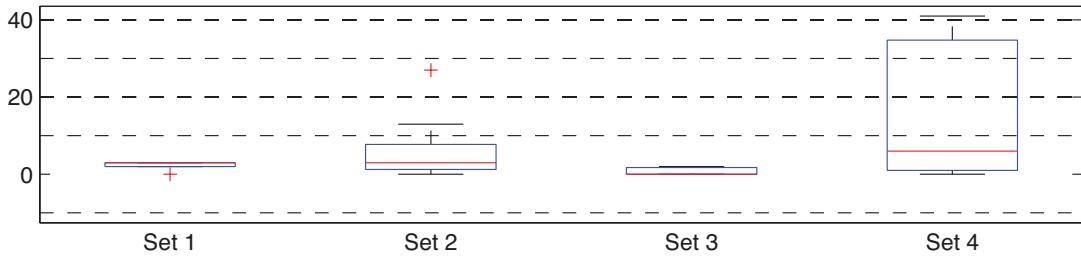


FIGURE 5.2: Box plot of the results for the artificial data sets, indicating the number of data points between the actual and closest detected change points. A lower and more compact box plot is better.

1. **Fixed increasing mean:** In Figure 5.3a we see that each change point is almost equally fast and accurate detected. This implies that the relative difference between the means of the Gaussian distribution has no effect on the detectability of the algorithm.
2. **Reduced increasing mean:** Figure 5.3b shows that the absolute difference between the means of the Gaussian distribution is a certain influence to the detectability. Although all change points are correctly detected (at the beginning of the series are two outliers), the change in radius becomes less significant at every change point.
3. **Reduced increasing mean, increasing variance:** In Figure 5.3d we see the same effect, in regard to the mean, as with the previous data set. Furthermore, the increasing variance makes change detection harder. This is visible around the 9000<sup>th</sup> data point and after, where the radius keeps increasing.
4. **Alternating variance:** The ability to detect decreases in variance of the Gaussian distribution is displayed in Figure 5.3c. Although some change points are discovered with a delay almost up to the length of the window processing the data, the decreases are all detected. The box plot of this set, number 4 in Figure 5.2, shows a wider spread for the results.

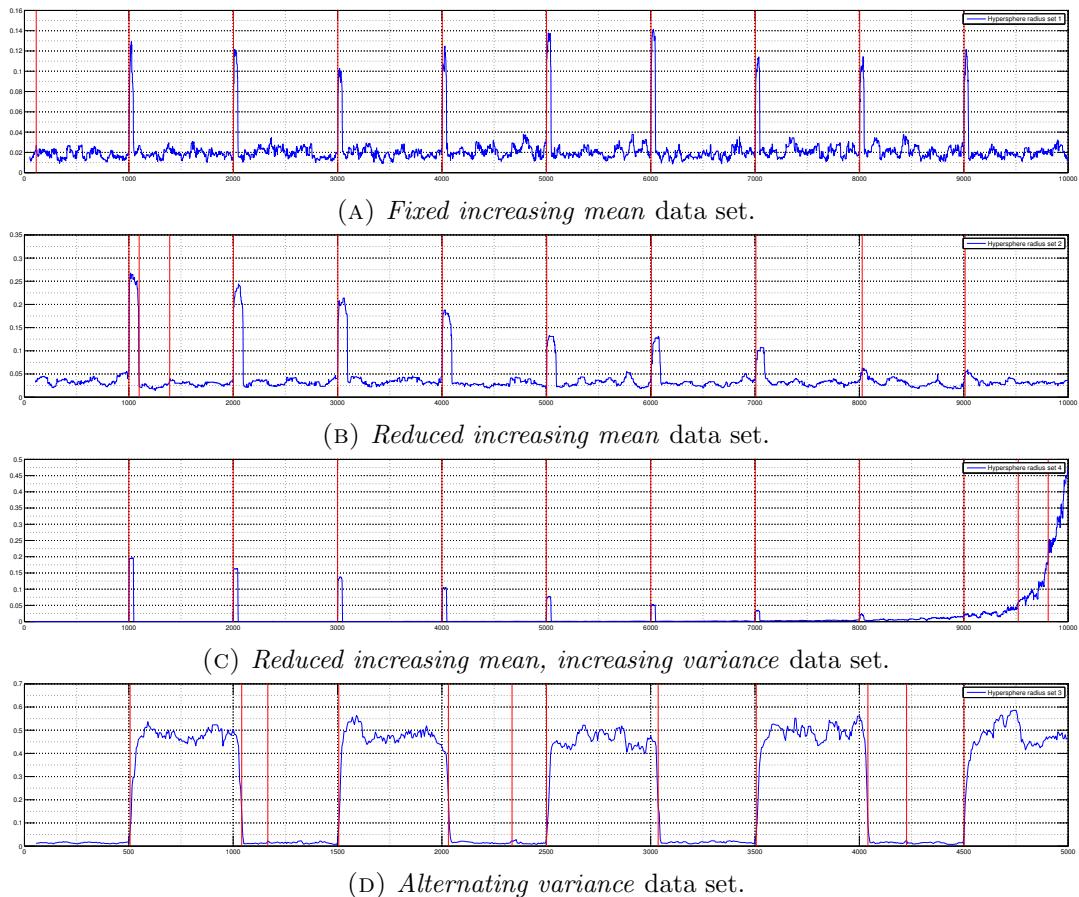


FIGURE 5.3: Hypersphere radius sizes for the four artificial data sets. The detected change points are displayed as red vertical bars.

# Chapter 6

## Real-world results

In this chapter we discuss the real-world data sets used for this research and the performance of our proposed algorithm on these data sets. In Section 6.1 we discuss common available data sets and argue that these are inappropriate for our research. It is followed by an overview of the setup we have created to record our own continuous data sets in an uncontrolled environment, both for in- and outdoor activities. The series of activities performed are discussed in Section 6.2. In that section we show typical plots of the inertial sensor data and stills from the video recordings. That will illustrate the setting in which the activities are performed and the characteristics of continuous data. Finally, the results of the algorithm are discussed in Section 6.3. We will use the objective measures, introduced in Section 5.2, and also give subjective claims on the performance of the algorithm applied to the real-world data. Since the data sets used are real-world recordings, we will put emphasis on the latter type of quality measures i.e. the subjective claims.

### 6.1 Data Sets

In the previous chapter we have applied our method to the artificial data sets from [14, 63]. In this chapter we apply the method to real-world data sets. The data sets come from inertial sensor data recorded during human activities, such as walking, running, ascending and descending stairs, and standing still. To compare our results with other research, we have tried to apply our method to two commonly used data sets. The first is the data from the WISDM lab [44]. As an excerpt from the data in Figure 6.1 shows, the activities are recorded non-continuous. Whilst this data set is convenient for activity classification, it does not allow for detection of change points. Even if the gap between the activities would be removed, the transition segments between activities

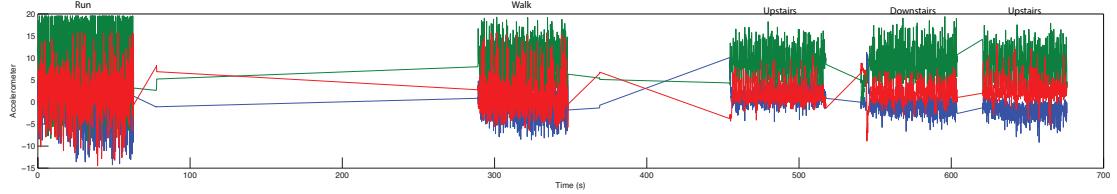


FIGURE 6.1: Excerpt from the WISDM data set [44]. First five activities for subject 33 are displayed. Due to the discontinuous nature of the recordings, this data set can not be used for change detection.

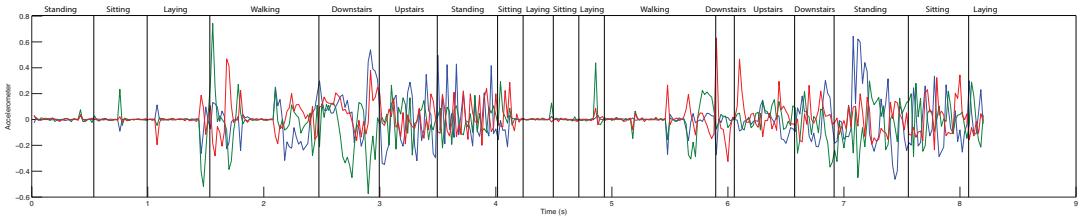


FIGURE 6.2: Excerpt from the UCI HAR data set [6]. The recorded training activities for subject 25 are displayed. Due to the imprecise (and possibly incorrect) labeling and the short duration of activities, the data set can not be used for the change detection method as described in this thesis.

would still be missing and thus the data set would not reflect continuous recordings. The other commonly used data set is the UCI HAR from the Machine Learning Repository and originates from [6]. As shown in an excerpt in Figure 6.2, the labeling to the data segments seems to be incorrect. Furthermore, it is not clear how the data is recorded and whether all activities are performed continuously.

Because of the above described shortcomings of commonly used data sets, we have decided to record activities by ourselves. The following characterizing requirements were set for our data sets:

1. The activities need to be performed and recorded in a continuous manner,
2. The environment should be natural and uncontrolled, both in- and outdoor,
3. The data should be clearly annotated, as objective as possible,
4. The required hardware (and software) should be widely available.

With these requirements in mind, we have created the following setup. All activities were recorded by a smartphone worn on body in the right front pants pocket of the subject. The used phones are the HTC SENSATION XE and HTC DESIRE, running the ANDROID smartphone operating system. During the activities the inertial sensor data is recorded using a free application, SENSOR LOGGER [45]. This application records the

sensor data to CSV files, which are then transformed to MATLAB compatible files. The constructed algorithm processes the data in an online manner (using only current and historic data). The MATLAB implementation uses the SVDD library DD\_TOOLS by Tax [69]. During the activities, the subject is recorded with a video camera to enable manual labeling and determination of the change points. As illustrated in the top row of Figure 6.5, the manual labeling is an ambiguous task. In Figure 6.5b a transition between walking and running is shown. It is impossible to determine the exact location of the change point and thus objective quality measures (such as the delay) are of less importance.

Using the video recordings, the data sets are manually annotated, as shown in the plots of Figures 6.3, 6.4 and 6.7. The graphs display each from a single recording the inertial sensor values (see Table 4.1 for explanation of the metrics), with manually annotated activities above the data. The manually determined change points are used to draw quality conclusion over the algorithm's performance, by comparing it to the discovered change points. All of the used data sets, including the video material and manual labeling, are made publicly available [73] for further research. **\*\*\* TODO: Publish web page with data sets \*\*\*** The following section discusses the performed activities during the recordings.

## 6.2 Data Gathering

For our data sets of real-world activities we have sampled both in- and outdoor activities. The subjects performing the activities are three males, of an age around 25. We have used four types of recordings. The first type (Run 1 and Run 2) consists of various outdoor recordings with segments of walking, running and standing. For the second type (Run 4, 5, and 6) some of the same activities are performed in a straight line, whilst others are performed around corners or in a wide circular direction. The third type (Run 8) consists mainly of walking up and down circular shaped stairs, with small horizontal parts between the segments. The last type is for Run 3, in which we walked in a straight line, took a 90° counter-clockwise turn (around the corner) and continued walking. A more detailed description of the performed activities follows.

### 6.2.1 Outdoor straight lines

The first series of outdoor performed activities consist of standing still, walking, and running segments in a straight line. Halfway the recording is a 180° turn. The plots of Figure 6.3 shows the recordings of Subject 1. The recorded metrics originate from

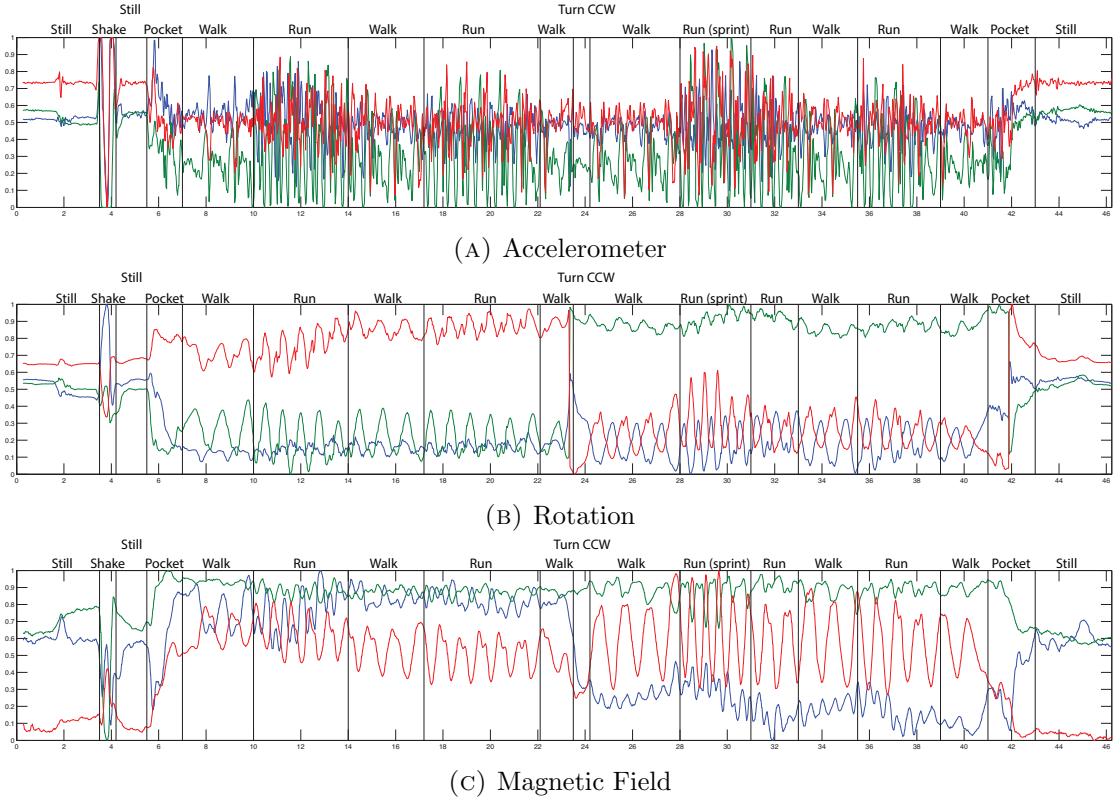


FIGURE 6.3: Annotated plots of a full recording performed by Subject 1. The vertical lines are manually determined change points. The label above each segment indicates the current performed activity.

the accelerometer, rotational and magnetometer sensors. The solid black lines indicate manually annotated change points. A still from the video recording, showing Subject 2 performing a (sprinting) run segment, is in Figure 6.5f.

### 6.2.2 Outdoor free

The second series of outdoor recordings consist of more free-form activities. The performed activities consist, as with the first set, of standing, walking, and running. Furthermore, the segments are performed in a circular manner. e.g. the running activities are performed around a fountain with an estimated diameter of 10 meters. The plots of the recordings of Subject 2, for the accelerometer, rotation, and magnetometer are shown in Figure 6.4. For Subject 2 the still of the video recordings are in the upper row of Figure 6.5, for Subject 1 in the lower row.

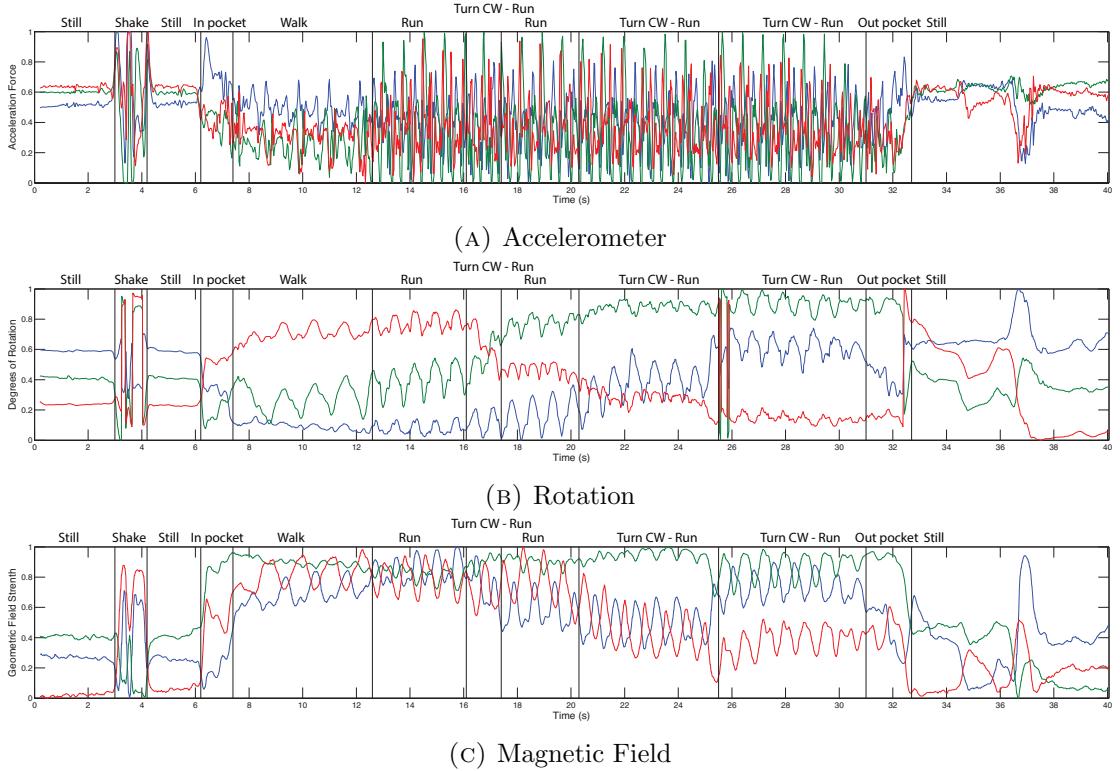


FIGURE 6.4: Annotated plots of a full recording performed by Subject 2. The vertical lines are manually determined change points. The label above each segment indicates the current performed activity.

### 6.2.3 Indoor stairs

The final series of activities recorded consist of walking up and down the stairs, indoor. The stairs are shaped in an semi-circular form, with a small horizontal segment on each floor. The still of Figure 6.6a illustrates the circular shape of the stairs. In Figure 6.6c the small flat segment is visible. Between the stairs the subject walked in a hallway with a  $180^\circ$  counter-clockwise turn, as shown in Figure 6.6b. The recorded sensor data of the performed activities are visible in the plots of Figures 6.7a to 6.7c.

## 6.3 Results

In this section we will provide the results of our methods applied to the real-world data sets and discuss the quality of the solution. As stated above, we will use the objective quality metrics as used in Section 5.2. That will be in the form of a tabular overview of the (manually set) best meta-parameters and results, expressed in False Alarm Rate (FAR) and Average\_delay (following Section 5.3). Furthermore, we will discuss the results from an empirical point of view, since the problem of finding change points in human activities is ambiguous. By stating our observations and remarks and providing



FIGURE 6.5: Stills of video recordings during performing of the outdoor activities of walking and running. Figure 6.5b shows the change point between the two activities. Figure 6.5d shows the running in clockwise motion.



FIGURE 6.6: Stills of video recordings during performing of the indoor activities of walking, ascending and descending the stairs.

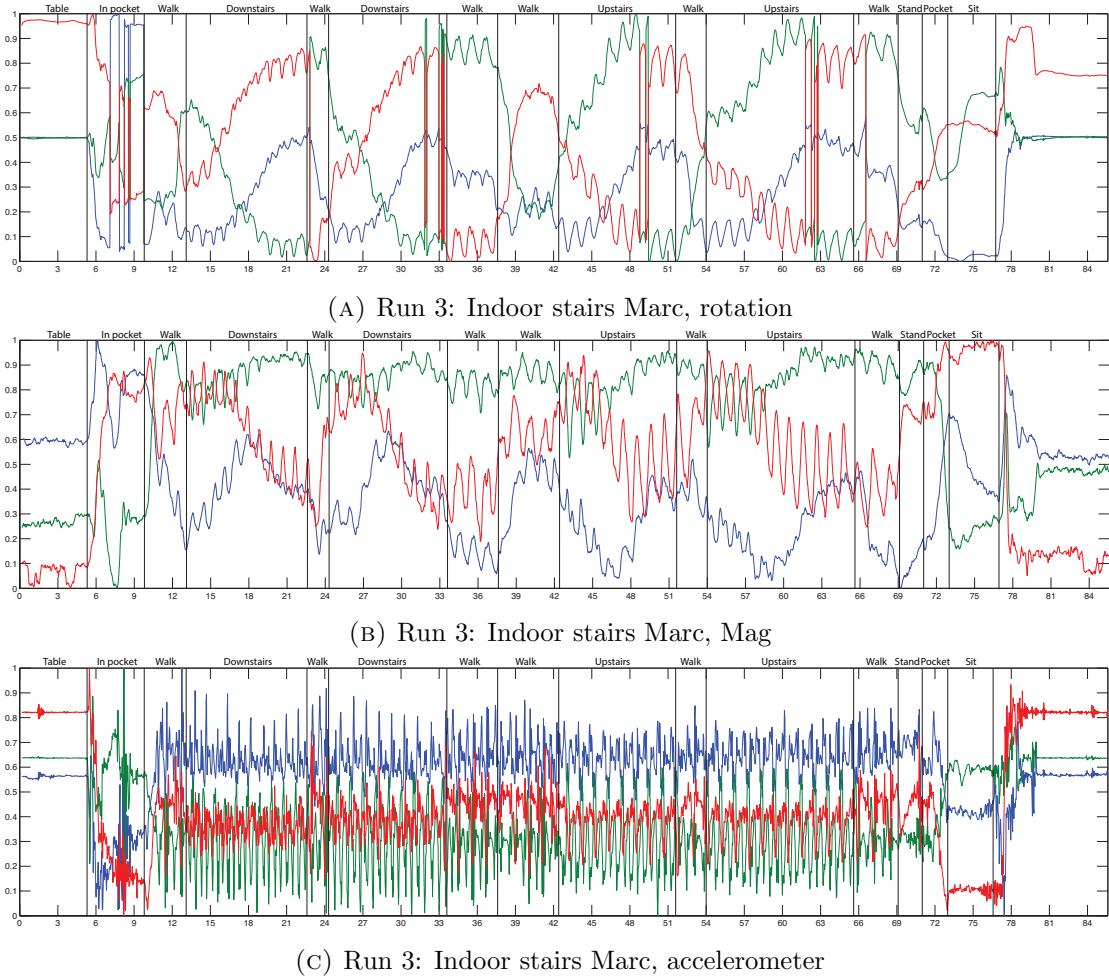


FIGURE 6.7: Annotated plots of a full recording performed by Subject 3 while walking the stairs. The vertical lines are manually determined change points. The label above each segment indicates the current performed activity.

fragments of visual results, we will discuss the characterizing strengths and weaknesses of our method.

### 6.3.1 Objective metrics

The employed two-stage algorithm requires the setting of user defined parameters. In Table 6.1 for each run the results are listed, together with the used parameters. We have used a window length of 50 frames, which represents around 1.16 seconds of data (the frame-rate of the recordings is not fixed, and we apply a slight filtering). The other parameter values are the result of an educated guess and manual optimization. In general, a high value for the upper threshold and low value for the lower threshold indicates a robust method, since it indicates that change points are significantly different from homogeneous segments. The closeness time frame should be low, since that indicates there are few false positives.

	Run 1	Run 2	Run 3	Run 4	Run 5	Run 6	Run 7
Window length	50	50	50	50	50	50	50
Sigma of RBF	13	13	13	13	13	13	4
High threshold	1.2	1.5	1.3	1.1	1.3	1.2	1.7
Low threshold	0.8	0.8	0.6	0.7	0.7	0.8	0.6
Closeness (s)	0.7	0.85	1	0.6	0.7	1	0.8
far( $Y$ )	0.1	0.05	0	0	0.2	0	0.48
Average_delay	0.52	0.83	0.47	0.92	1.26	1.05	0.92
STD Delay	0.29	0.55	0.52	0.75	1.67	0.78	0.92

TABLE 6.1: Parameter settings and results of the real-world data sets.

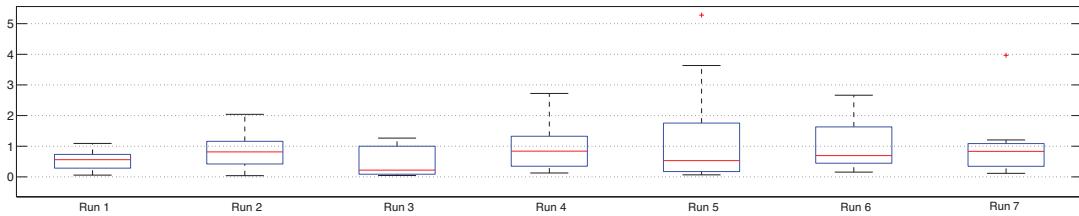


FIGURE 6.8: Box plot of the results for the real-world runs, indicating the number of data points between the actual and closest detected change points. A lower and more compact box plot is better.

To give a better impression on the quality of the Average\_delay, Figure 6.8 shows a box plot for each run. It shows the spread of the Average\_delay in seconds over all the change points. A lower and more compact box plot is better.

### 6.3.2 Subjective observations and remarks

Since the manual annotation of change points is ambiguous, the objective quality metrics can give a distorted view of the performance. To overcome this, we will state observations and remarks found after close (visual) inspection of the annotated and discovered change points. We will start with general observations applicable to all the runs and will continue with individual remarks for each run. Where possible, a graphical illustration will support our observations.

Considering the overall performance and results, we have the following observations:

- **Merging:** in our method we use a *closeness* time period  $t_c$  (usually up to 0.5 - 1.0 seconds) to merge discovered change points that have a very small distance between them. Multiple merging-strategies can be applied, and in our method we use the naive implementation by simply ignoring all the discovered change points that occur less than  $t_c$  seconds after the previous change point. This works well when there is a block of noisy data with a high amount of (falsely) discovered

change points. On the flip side, this method will ignore real new change points which occur during the noisy period. This problem on itself can be formulated as a change detection problem.

- **Masking:** Since the parameters for the detection method are set globally, it is difficult to discover all the change points (and only the change points) without a high FAR. If a change point is proceeded by a noisy block, then the real change point will be merged with the previous change points. In an other case, where a change point is very clear represented in the data but the next, close by, other change point is more subtle, the latter change point is *masked* by the former. This masking effect is also discussed in [39], where an iterated approach is applied.
- **False heterogeneity:** For our video-data synchronization, we started and ended each recording with a few seconds in which the recording smartphone was kept still in the air. During this period, the data variance, and thus the constructed hypersphere, becomes very small. As a results, even small movements are considered to be changes and in the case of the mid-air still smartphone a lot of false change points are detected. Due to the merging effect described above, the final real change point (often shaking the smartphone) is not discovered.
- **Incorrect weighting:** In our analyses we have used the data from the accelerometer, magnetic field, and rotation sensors. In the segments which embodied movement in a circular manner (such as walking and running around a fountain), the turn was not (or at a different time point) discovered. Alternatively, when we only used the magnetic field and rotation sensors, the turn was correctly discovered. Other transitions, such as from walking to running, are harder to discover without the (linear) accelerometer sensor data.

**\*\*\* TODO: Add positive remarks \*\*\***

When looking at the individual runs, we have the following observations for each run number:

### 1. Subject 2, straight walk and run

- The first transition from running to walking, around 14s is harder to discover than the transition for the same activities around 33s. This shows is that in real-world applications there is a diversity between the same transitions and activities.
- Around 8s, 9s, and 16s a few steps (while walking) are regarded as change points. During the running segments from 17s and 28s there is a lower probability of change for each step.

## 2. Subject 1, straight walk and run

- Following the video recordings, we have annotated a change point from running to walking around 37s. Our method discovers a change point almost a second before. In retrospect, we can see that the data distribution indeed changes from the discovered change point on. This shows us two important principles. The first is that the annotations are very subjective. The second is that between different activities the transition period is longer than we would think. Looking at the data, we can see that the body slows down, even before we visually notice it on the video recordings.

## 3. Subject 2, walk around corner

- The 90° counter-clockwise turn during the walking activity is hard to discover when the accelerometer sensor data is included. When only the magnetic field and rotation sensors are used, the turn requires a lower sensitivity. With only these two sensors all the other change points in this run are also successfully discovered.

## 4. Subject 2, walk and run around fountain

- Like in the first run, the walking segment from 24s results in a change point for each step. Further inspection of the data reveals that each step is indeed different from the other. Due to the global parameter settings, the sensitivity is too high for this segment to recognize it as one.
- During the circular run, from 12s till 24, there are two change points discovered. The difference for the rotational vectors need to accumulate to a certain value before they have enough influence to let the rotation be regarded as a change point. **\*\*\* TODO: Beter formuleren \*\*\*.**

## 5. Subject 1, walk and run around fountain

- As with the other runs, the accelerometer data makes it harder to detect turns. It requires a higher sensitivity, which results in a higher FAR.

## 6. Subject 2, walk and run fountain 2

- During the standing segment around 38s there are a lot of false positives. It seems to be the *false heterogeneity* problem described above.

## 7. Subject 3, indoor stairs

- The walking segment around 22s, between two segments of walking down-stairs, shows little difference in the data. To recognize it as a change point a high sensitivity and low closeness time period  $t_c$  is required.

- During some segments (downstairs from 24s, upstairs from 42s and 54s) the method discovers more change points than our annotation. A closer inspecting of the raw data reveals indeed changes in behavior. To exclude these (semi) false positives, a better tuning of parameters is required.
- Because of the circular shape of the stairs, the magnetic field sensors constantly differs. Although our method is build to exclude slow shifting changes (because we are only interested in sudden changes), with our used window width it still eventually results in change points.
- The difference between taking the stairs and walking is smaller than, e.g. , walking and running. The delay between these segments seems to be larger, as illustrated around 33s.

**\*\*\* TODO: Samenvatten observations en remarks, brugje naar possible improvements maken \*\*\***

## 6.4 Possible improvements

From the above section we can extract a few possible improvements for the algorithm. Most of the observations and remarks are due to a different required sensitivity for local segments of the run. Currently, our method employs a global setting for the thresholds. A possible improvement would be to use locally optimized parameters, although that would require an reflective data processing approach or more a priori about the data. With our approach, the lack of a priori knowledge about the data distribution is one of the design decisions.

Furthermore, we discovered that sometimes, especially with movements incorporating circular motions, the addition of (linear) accelerometer data makes the discovery of change points harder. In such cases, each of the inertial sensor data streams should be relatively weighted (e.g. using the covariance) in order to signal an overall change when it only occurs in one of the streams.

# Chapter 7

## Conclusion

### 7.1 Outline

*Not intended for the reader.*

- Tie together, integrate, synthesize all the issues from the chapters, reflecting on introductory thesis statements.
- Provide answers to research questions
- Highlight study limitations
- Directions for further research.
- Content: Summarize
  - What have I researched
  - Main arguments
  - How researched
  - What discovered
  - What pre-existing views were challenged (?)
- Overview of
  - New knowledge
  - Significance of research
  - Limitations (data, concepts)
  - Speculations on implications of research

- Further research and development, or: links with different fields, other methods applied to same data.
- Must: make clear statement of my original contribution
- Ideally:
  - Shows links across key ideas spread in chapters
  - Show commitment and enthusiasm for research
  - Positive impression
- Avoid:
  - claim findings that have not been proven
  - introducing new data/concepts
  - Be self-critical, but dont put weaknesses or limitation in research
  - Too long/too short

Paragraph: into on conclusion; what have I researched and what is my contribution.  
Break into sections further down.

====

In this research we have applied OCC methods to accelerometer data, recorded during the performance of various human activities, in order to create a temporal segmentation of those recordings. The goal was to find the change points between activities, assuming that it will aid better understanding of the sensor values. To do so, we used the SVDD algorithm by Tax [65], following the method of applying it to a sliding window of data, as done by Camci [14]. Instead of using both the information whether a new data point is an outlier and the radius of the hypersphere to indicate change, we only used the latter model properties. The increase of the radius of the constructed hypersphere, which encloses the data objects and increases with heterogeneous segments of data, is extracted and used for indication of change.

**\* \*\*\* TODO: Write overview of this chapter. In the following section we will... \*\*\***

## 7.2 What is done

The scope of this research was the temporal segmentation of recordings obtained from human activities. This was performed in the wider context of activity classification.

Currently, many algorithms obtain an implicit segmentation as a side-effect of direct activity classification. In this research the primary goal was finding a temporal segmentation, which is assumed to be able to aid the classification step. This problem is reduced to finding change points in temporal data sets originating from inertial sensors.

We have considered a range of segmentation methods, from which the method by Camci [14] showed potential. It is based on the detection of outliers in a time series data, assuming that an increase of outliers indicates a change in the underlying generating model. For the detection of outliers it uses the SVDD method by Tax [65]. It tests every new data point with the constructed model; whether the data point is in or outside the model, and whether the model increases or decreases in radius size. The constructed sliding-window algorithm, SVCPD, is applied to artificial Gaussian noise data.

In this research we have applied the method by Camci [14] to data from inertial sensors, recorded by smartphones. We have used the accelerometer data for measurement of speed, the gyroscope data for measurement of rotation, and the magnetometer to measure the direction and orientation of the performed activities. The proposed method reduces this 9-dimensional signal to a single property, obtained from the constructed model, which can be interpreted to indicate change.

The constructed model follows the implementation of Tax and Laskov [68] for an incremental version of SVDD, I-SVDD. This algorithm creates a OC-SVM model by processing the data over a sliding window. The model represents a hypersphere from which the radius size is extracted. That property is of interest since the assumption is that a heterogeneous window of data will show an increase of radius, in relation to a homogeneous segment of data. The homogeneous segment of data will have a relatively small radius since the data points will be close together, resulting from the continuity assumptions.

To test our approach, we have used the same data as used in the experiments from Camci [14] and Takeuchi and Yamanishi [63], for which the results are comparable with the original research (see Section 5.3). Furthermore, since this research is focused around human activities, we have recorded and manually annotated in- and outdoor activities. The discovered change points were compared to the annotated change points, obtained from video recordings. The results where discussed in Section 6.3. In this research we found that other public available and wide used common data sets, such as the WISDM [44] and UCI HAR [6], were not useful for our purpose, since the activities are non-continuously recorded.

**\*\*\* TODO:** *In the above, refer to the corresponding chapters/sections. \*\*\**

### 7.3 Main findings

What we have found, show links across chapters.

For each finding a different paragraph. How did we arrive at this finding, and it challenges/improves on previous research. Not just a summary, but the effects and implications.

While temporal segmentation of time series data is not uncommon, the application to inertial signals is currently not under much research attention. It is applied in other contexts, such as motion tracking and first person vision sensing, the segmentation obtained from inertial signals is often the implicit result of classification.

### 7.4 Contributions

### 7.5 Other research

Relate the finding with earlier research. Provide similarities and differences. Better/-worse results?

### 7.6 Limitations

Show what limitations were encountered or decided (data collection, scope). Note the focus of the limitations and how they worked together (to get positive note). This can create bridge to "future research".

### 7.7 Future research

The motivation for this research was, amongst others, the assumption that an explicit temporal segmentation of inertial signals can be used to obtain a better classification of the performed activities. There are two reasons for that. First, since the segmentation adds information about sequential data points (whether they belong to the same class), the a priori knowledge about the distribution changes. Second, the window of data used for model construction can be much larger (than the often used length of 1-3 seconds). This enables a better mathematical model to be constructed. A future research in this line of assumptions is to compare the performance of classification methods using this obtained temporal segmentation and those that do not.

Since the temporal segmentation of inertial signal is not yet a very active field of study, more studies can be performed to obtain it. In our approach we have used the raw data directly. In the field of activity recognition many methods make use of pre-processing steps, from which one step often is feature extraction from the raw data. That, and other pre-processing steps, could be incorporated in the explicit temporal segmentation of inertial signals. Future research could focus on extracted features that aid the detection of change points.

*\*\*\* TODO: Add future research on OC-SVM applications. \*\*\**

*\*\*\* TODO: Improvements of our method, Section 6.4 \*\*\**

What has not been covered? What is worthwhile. Show that I am thinking ahead and know the field of study.

## 7.8 Final section

Reminding of original contribution and significance. This section is the overall conclusion and thus should be very concise and precise.

*\*\*\* TODO: Formulate the research questions in the introduction of thesis, and in the chapters. \*\*\**

# Bibliography

- [1] Ryan Prescott Adams and David JC MacKay. Bayesian online changepoint detection. *arXiv preprint arXiv:0710.3742*, 2007.
- [2] A Aizerman, Emmanuel M Braverman, and LI Rozoner. Theoretical foundations of the potential function method in pattern recognition learning. *Automation and remote control*, 25:821–837, 1964.
- [3] Hirotugu Akaike. A new look at the statistical model identification. *Automatic Control, IEEE Transactions on*, 19(6):716–723, 1974.
- [4] Cesare Alippi and Manuel Roveri. An adaptive cusum-based test for signal change detection. In *Circuits and Systems, 2006. ISCAS 2006. Proceedings. 2006 IEEE International Symposium on*, pages 4–pp. IEEE, 2006.
- [5] Davide Anguita, Alessandro Ghio, Stefano Pischiutta, and Sandro Ridella. A hardware-friendly support vector machine for embedded automotive applications. In *Neural Networks, 2007. IJCNN 2007. International Joint Conference on*, pages 1360–1364. IEEE, 2007.
- [6] Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra, and Jorge L Reyes-Ortiz. Human activity recognition on smartphones using a multiclass hardware-friendly support vector machine. In *Ambient Assisted Living and Home Care*, pages 216–223. Springer, 2012.
- [7] Akin Avci, Stephan Bosch, Mihai Marin-Perianu, Raluca Marin-Perianu, and Paul Havinga. Activity recognition using inertial sensing for healthcare, wellbeing and sports applications: A survey. In *Architecture of Computing Systems (ARCS), 2010 23rd International Conference on*, pages 1–10. VDE, 2010.
- [8] Ling Bao and Stephen Intille. Activity recognition from user-annotated acceleration data. *Pervasive Computing*, pages 1–17, 2004.

- [9] J. Barbić, A. Safanova, J.Y. Pan, C. Faloutsos, J.K. Hodgins, and N.S. Pollard. Segmenting motion capture data into distinct behaviors. In *Proceedings of Graphics Interface 2004*, pages 185–194. Canadian Human-Computer Communications Society, 2004.
- [10] Vic Barnett and Toby Lewis. *Outliers in statistical data*, volume 3. Wiley New York, 1994.
- [11] M. Basseville, I.V. Nikiforov, et al. *Detection of abrupt changes: theory and application*, volume 104. Prentice Hall Englewood Cliffs, NJ, 1993.
- [12] Thomas Bernecker, Franz Graf, Hans-Peter Kriegel, Christian Moennig, Dieter Dill, and Christoph Tuermer. Activity recognition on 3d accelerometer data (technical report). 2012.
- [13] Robert L Brown, James Durbin, and John M Evans. Techniques for testing the constancy of regression relationships over time. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 149–192, 1975.
- [14] Fatih Camci. Change point detection in time series data using support vectors. *International Journal of Pattern Recognition and Artificial Intelligence*, 24(01):73–95, 2010.
- [15] F Chamroukhi, S Mohammed, D Trabelsi, L Oukhellou, and Y Amirat. Joint segmentation of multivariate time series with hidden process regression for human activity recognition. *Neurocomputing*, 2013.
- [16] Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2(3):27, 2011.
- [17] Wei-Cheng Chang, Ching-Pei Lee, and Chih-Jen Lin. A revisit to support vector data description (svdd).
- [18] Tsung-Lin Cheng. An efficient algorithm for estimating a change-point. *Statistics & Probability Letters*, 79(5):559–565, 2009.
- [19] Vladimir Cherkassky and Filip M Mulier. *Learning from data: concepts, theory, and methods*. Wiley. com, 2007.
- [20] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- [21] M.O. Derawi. Accelerometer-based gait analysis, a survey. *Norsk informasjonsikkerhetskonferanse (NISK)*, 2010.

- [22] R.W. DeVaul and S. Dunn. Real-time motion classification for wearable computing applications. 2001, project paper, <http://www.media.mit.edu/wearables/mithril realtime.pdf>, 2001.
- [23] Andrés Duque, Fco Javier Ordóñez, Paula de Toledo, and Araceli Sanchis. Offline and online activity recognition on mobile devices using accelerometer data. In *Proceedings of the 4th international conference on Ambient Assisted Living and Home Care*, pages 208–215. Springer-Verlag, 2012.
- [24] Peter Flach. *Machine learning: the art and science of algorithms that make sense of data*. Cambridge University Press, 2012.
- [25] Erich Fuchs, Thiemo Gruber, Jiri Nitschke, and Bernhard Sick. Online segmentation of time series based on polynomial least-squares approximations. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 32(12):2232–2245, 2010.
- [26] Keinosuke Fukunaga. *Introduction to statistical pattern recognition*. Access Online via Elsevier, 1990.
- [27] Federico Girosi. An equivalence between sparse approximation and support vector machines. *Neural computation*, 10(6):1455–1480, 1998.
- [28] E. Guenterberg, H. Ghasemzadeh, V. Loseu, and R. Jafari. Distributed continuous action recognition using a hidden markov model in body sensor networks. *Distributed Computing in Sensor Systems*, pages 145–158, 2009.
- [29] E. Guenterberg, S. Ostadabbas, H. Ghasemzadeh, and R. Jafari. An automatic segmentation technique in body sensor networks based on signal energy. In *Proceedings of the Fourth International Conference on Body Area Networks*, page 21. ICST (Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering), 2009.
- [30] Tian Guo, Zhixian Yan, and Karl Aberer. An adaptive approach for online segmentation of multi-dimensional mobile data. In *Proceedings of the Eleventh ACM International Workshop on Data Engineering for Wireless and Mobile Access*, pages 7–14. ACM, 2012.
- [31] V. Guralnik and J. Srivastava. Event detection from time series data. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 33–42. ACM, 1999.
- [32] Fredrik Gustafsson. The marginalized likelihood ratio test for detecting abrupt changes. *Automatic Control, IEEE Transactions on*, 41(1):66–78, 1996.

- [33] Wolfgang Härdle. *Nonparametric and semiparametric models*. Springer Verlag, 2004.
- [34] Zhen-Yu He and Lian-Wen Jin. Activity recognition from acceleration data using ar model representation and svm. In *Machine Learning and Cybernetics, 2008 International Conference on*, volume 4, pages 2245–2250. IEEE, 2008.
- [35] Kathryn Hempstalk, Eibe Frank, and Ian H Witten. One-class classification by combining density and class probability estimation. In *Machine Learning and Knowledge Discovery in Databases*, pages 505–519. Springer, 2008.
- [36] J. Himberg, K. Korpiaho, H. Mannila, J. Tikanmaki, and H.T.T. Toivonen. Time series segmentation for context recognition in mobile devices. In *Data Mining, 2001. ICDM 2001, Proceedings IEEE International Conference on*, pages 203–210. IEEE, 2001.
- [37] Victoria Hodge and Jim Austin. A survey of outlier detection methodologies. *Artificial Intelligence Review*, 22(2):85–126, 2004.
- [38] Chih-Chiang Hsu. The mosum of squares test for monitoring variance changes. *Finance Research Letters*, 4(4):254–260, 2007.
- [39] Carla Inclán and George C Tiao. Use of cumulative sums of squares for retrospective detection of changes of variance. *Journal of the American Statistical Association*, 89(427):913–923, 1994.
- [40] Takafumi Kanamori, Shohei Hido, and Masashi Sugiyama. A least-squares approach to direct importance estimation. *The Journal of Machine Learning Research*, 10: 1391–1445, 2009.
- [41] Y. Kawahara and M. Sugiyama. Change-point detection in time-series data by direct density-ratio estimation. In *Proceedings of 2009 SIAM International Conference on Data Mining (SDM2009)*, pages 389–400, 2009.
- [42] E. Keogh, S. Chu, D. Hart, and M. Pazzani. An online algorithm for segmenting time series. In *Data Mining, 2001. ICDM 2001, Proceedings IEEE International Conference on*, pages 289–296. IEEE, 2001.
- [43] Shehroz S Khan and Michael G Madden. A survey of recent trends in one class classification. In *Artificial Intelligence and Cognitive Science*, pages 188–197. Springer, 2010.
- [44] Jennifer R Kwapisz, Gary M Weiss, and Samuel A Moore. Activity recognition using cell phone accelerometers. *ACM SIGKDD Explorations Newsletter*, 12(2): 74–82, 2011.

- [45] KZ-S. Sensor logger. URL <https://play.google.com/store/apps/details?id=com.kzs6502.sensorlogger>.
- [46] C. Li, SQ Zheng, and B. Prabhakaran. Segmentation and recognition of motion streams by similarity search. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMCCAP)*, 3(3):16, 2007.
- [47] Kun-Lun Li, Hou-Kuan Huang, Sheng-Feng Tian, and Wei Xu. Improving one-class svm for anomaly detection. In *Machine Learning and Cybernetics, 2003 International Conference on*, volume 5, pages 3077–3081. IEEE, 2003.
- [48] Song Liu, Makoto Yamada, Nigel Collier, and Masashi Sugiyama. Change-point detection in time-series data by relative density-ratio estimation. *Neural Networks*, 2013.
- [49] Xiaoyan Liu, Zhenjiang Lin, and Huaiqing Wang. Novel online methods for time series segmentation. *Knowledge and Data Engineering, IEEE Transactions on*, 20(12):1616–1626, 2008.
- [50] Junshui Ma and Simon Perkins. Time-series novelty detection using one-class support vector machines. In *Neural Networks, 2003. Proceedings of the International Joint Conference on*, volume 3, pages 1741–1745. IEEE, 2003.
- [51] Markos Markou and Sameer Singh. Novelty detection: a review part 1: statistical approaches. *Signal processing*, 83(12):2481–2497, 2003.
- [52] Markos Markou and Sameer Singh. Novelty detection: a review part 2: neural network based approaches. *Signal processing*, 83(12):2499–2521, 2003.
- [53] Giorgos Mountrakis, Jungho Im, and Caesar Ogole. Support vector machines in remote sensing: A review. *ISPRS Journal of Photogrammetry and Remote Sensing*, 66(3):247–259, 2011.
- [54] MM Moya, MW Koch, and LD Hostetler. One-class classifier networks for target recognition applications. Technical report, Sandia National Labs., Albuquerque, NM (United States), 1993.
- [55] Zineb Noumir, Paul Honeine, and Cedric Richard. On simple one-class classification methods. In *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium on*, pages 2022–2026. IEEE, 2012.
- [56] ES Page. Continuous inspection schemes. *Biometrika*, 41(1/2):100–115, 1954.
- [57] Roberto Perdisci, Guofei Gu, and Wenke Lee. Using an ensemble of one-class svm classifiers to harden payload-based anomaly detection systems. In *Data Mining, 2006. ICDM’06. Sixth International Conference on*, pages 488–498. IEEE, 2006.

- [58] Bernhard Schölkopf and Alexander J Smola. *Learning with kernels*. The MIT Press, 2002.
- [59] Bernhard Schölkopf, Robert C Williamson, Alex J Smola, John Shawe-Taylor, and John C Platt. Support vector method for novelty detection. In *NIPS*, volume 12, pages 582–588, 1999.
- [60] Alex J Smola, Bernhard Schölkopf, and Klaus-Robert Müller. The connection between regularization operators and support vector kernels. *Neural Networks*, 11(4):637–649, 1998.
- [61] Masashi Sugiyama, Taiji Suzuki, Shinichi Nakajima, Hisashi Kashima, Paul von Bünau, and Motoaki Kawanabe. Direct importance estimation for covariate shift adaptation. *Annals of the Institute of Statistical Mathematics*, 60(4):699–746, 2008.
- [62] Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. *Density ratio estimation in machine learning*. Cambridge University Press, 2012.
- [63] Jun-ichi Takeuchi and Kenji Yamanishi. A unifying framework for detecting outliers and change points from time series. *Knowledge and Data Engineering, IEEE Transactions on*, 18(4):482–492, 2006.
- [64] David MJ Tax. One-class classification. 2001.
- [65] David MJ Tax and Robert PW Duin. Support vector domain description. *Pattern recognition letters*, 20(11):1191–1199, 1999.
- [66] David MJ Tax and Robert PW Duin. Uniform object generation for optimizing one-class classifiers. *The Journal of Machine Learning Research*, 2:155–173, 2002.
- [67] David MJ Tax and Robert PW Duin. Support vector data description. *Machine learning*, 54(1):45–66, 2004.
- [68] David MJ Tax and Pavel Laskov. Online svm learning: from classification to data description and back. In *Neural Networks for Signal Processing, 2003. NNSP'03. 2003 IEEE 13th Workshop on*, pages 499–508. IEEE, 2003.
- [69] D.M.J. Tax. Ddtools, the data description toolbox for matlab, Oct 2013. version 2.0.2.
- [70] Ferdinand Van Der Heijden, Robert Duin, Dick De Ridder, and David MJ Tax. *Classification, parameter estimation and state estimation: an engineering approach using MATLAB*. Wiley. com, 2005.
- [71] Vladimir Vapnik. Pattern recognition using generalized portrait method. *Automation and remote control*, 24:774–780, 1963.

- [72] Vladimir Vapnik. Statistical learning theory. 1998, 1998.
- [73] R.Q. Vlasveld. Continuously recorded human activities data set, 2014. URL <http://rvlasveld.github.io/blog/2014/02/05/continuously-recorded-human-activity-inertial-sensor-data/>.
- [74] J.A. Ward, P. Lukowicz, G. Troster, and T.E. Starner. Activity recognition of assembly tasks using body-worn microphones and accelerometers. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 28(10):1553–1567, 2006.
- [75] Makoto Yamada, Taiji Suzuki, Takafumi Kanamori, Hirotaka Hachiya, and Masashi Sugiyama. Relative density-ratio estimation for robust distribution comparison. *Neural computation*, 25(5):1324–1370, 2013.
- [76] A.Y. Yang, S. Iyengar, S. Sastry, R. Bajcsy, P. Kuryloski, and R.afari. Distributed segmentation and classification of human actions using a wearable motion sensor network. In *Computer Vision and Pattern Recognition Workshops, 2008. CVPRW’08. IEEE Computer Society Conference on*, pages 1–8. IEEE, 2008.
- [77] Jhun-Ying Yang, Jeen-Shing Wang, and Yen-Ping Chen. Using acceleration measurements for activity recognition: An effective learning algorithm for constructing neural classifiers. *Pattern recognition letters*, 29(16):2213–2220, 2008.
- [78] Jie Yin, Qiang Yang, and Jeffrey Junfeng Pan. Sensor-based abnormal human-activity detection. *Knowledge and Data Engineering, IEEE Transactions on*, 20(8):1082–1090, 2008.
- [79] F. Zhou, F. Torre, and J.K. Hodgins. Aligned cluster analysis for temporal segmentation of human motion. In *Automatic Face & Gesture Recognition, 2008. FG’08. 8th IEEE International Conference on*, pages 1–7. IEEE, 2008.