

UTRECHT UNIVERSITY

DOCTORAL THESIS

---

# Human Activity Recognition Using Accelerometer Data

---

*Author:*

R.Q. VLASVELD

*Supervisor:*

Dr. James SMITH

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

*in the*

Research Group Name  
Department or School Name

March 2013

# Declaration of Authorship

I, R.Q. VLASVELD, declare that this thesis titled, 'Human Activity Recognition Using Accelerometer Data' 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:

---

*“Thanks to my solid academic training, today I can write hundreds of words on virtually any topic without possessing a shred of information, which is how I got a good job in journalism.”*

Dave Barry

UNIVERSITY NAME (IN BLOCK CAPITALS)

# *Abstract*

Faculty Name

Department or School Name

Master of Science

**Human Activity Recognition Using Accelerometer Data**

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...

# *Acknowledgements*

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

# Contents

<b>Declaration of Authorship</b>	<b>i</b>
<b>Abstract</b>	<b>iii</b>
<b>Acknowledgements</b>	<b>iv</b>
<b>List of Figures</b>	<b>vii</b>
<b>List of Tables</b>	<b>viii</b>
<b>Abbreviations</b>	<b>ix</b>
<b>Physical Constants</b>	<b>x</b>
<b>Symbols</b>	<b>xi</b>
<b>1 Introduction</b>	<b>2</b>
1.1 Literature review . . . . .	2
1.1.1 Segmentation and classification . . . . .	2
1.1.2 Temporal Segmentation . . . . .	5
1.1.3 List of papers . . . . .	8
1.1.3.1 Pure segmentation . . . . .	12
1.2 Problem statement . . . . .	13
<b>2 Techniques</b>	<b>14</b>
2.1 Signal pro-processing and sensor fusion for timed patterns . . . . .	14
2.1.1 Signal fusion . . . . .	14
2.1.2 Signal smoothing . . . . .	16
2.1.3 Feature extraction . . . . .	16
2.2 Temporal Segmentation . . . . .	16
2.2.1 Aims of segmentation . . . . .	16
2.2.2 Formal definition . . . . .	17
2.2.3 Application in research fields . . . . .	18
2.2.4 PCA Based Methods . . . . .	18
2.2.5 Statistical methods . . . . .	20

---

2.2.6	Hidden Markov Models Based Methods . . . . .	20
2.2.7	Bayesian Methods . . . . .	20
2.2.8	Principal Component Analysis . . . . .	20
2.2.9	Hidden Markov Models . . . . .	22
2.2.10	Segmentation as clustering . . . . .	24
2.3	Clustering . . . . .	25
2.3.1	Formal definition . . . . .	27
2.3.2	Application in research fields . . . . .	27
2.4	Temporal pattern recognition . . . . .	27
2.4.1	Dynamic Time Warping . . . . .	27
2.4.2	k-means clustering . . . . .	27
2.4.3	Self-organizing Map . . . . .	28
2.4.4	Support Vector Machine . . . . .	28
2.4.5	Naïve Bayes . . . . .	28
2.5	Unsupervised clustering of temporal patterns . . . . .	28
<b>3</b>	<b>Experiments</b>	<b>29</b>
3.1	. . . . .	29
<b>4</b>	<b>Results</b>	<b>30</b>
4.1	. . . . .	30
<b>5</b>	<b>Discussion</b>	<b>31</b>
5.1	. . . . .	31
<b>6</b>	<b>Conclusions</b>	<b>32</b>
6.1	. . . . .	32
<b>A</b>	<b>Appendix Title Here</b>	<b>33</b>
	<b>Bibliography</b>	<b>34</b>

# List of Figures

2.1	K-means . . . . .	15
2.2	K-means . . . . .	27



# List of Tables

1.1	Overview of researches with the focus on segmentation and classification techniques. . . . .	5
-----	--	---

# Abbreviations

**LAH** List Abbreviations **H**ere

**PCA** Principle Component **A**nalysis

# Physical Constants

$$\text{Speed of Light } c = 2.997\,924\,58 \times 10^8 \text{ ms}^{-\text{s}} \text{ (exact)}$$

# Symbols

$a$	distance	m
$P$	power	W ( $\text{Js}^{-1}$ )
$\omega$	angular frequency	$\text{rads}^{-1}$

*For/Dedicated to/To my...*

\*\*\* [GRAPHICS diagram: sensors → signal pre-processing → feature extraction → analysis [\* segmentation, \* labeling] → conclusions] \*\*\* Give context problems, means to reach (e.g. monitoring of patients, getting conclusions over health and activity).

# Chapter 1

## Introduction

### 1.1 Literature review

This section will provide a literature review per coherent subject of this thesis.

#### 1.1.1 Segmentation and classification

In recent years much research is done in the field of segmentation and classification of signals into human activities. Due to the nature of many applied algorithms and the field of research (computer vision, data mining, robotics), often there is no clear distinction between the segmentation and classification phase. As a result, these two subjects will be discussed together although, where possible, the different techniques for the two will be mentioned. All the reviewed papers have a global setup in common; a single- or multiple- sensory device is worn by subjects. From the raw signals features are extracted which are used to train some models. During the testing phase the same extracted features are compared with the trained models to classify the activity. Some methods rely on clear segmentations of the signals, whilst other ignore this possible feature and classify on sliding windows. The content of this section is summarized in table [1.1](#).

During the last ten years much research is performed, as a result of smaller and cheaper measuring devices. \*\*\* cite needed \*\*\*. With the introduction of smartphones with sensory hardware the research has shifted to these devices, as they create a more natural and ubiquitous setting for the test subject. \*\*\* cite needed \*\*\*.

A distinction can be made in the sensor setup, like wearing multiple sensors, a single sensor and the location of the sensors. Research have been performed with multiple

sensors, referred to as Body Sensor Networks, as done in [1–4]. Locations used include the upper arm, waist, thigh, wrist and ankle. Single sensors have been worn at different locations, including the upper arm [5], the front leg pocket [6–9], the waist [10–12], the chest [13, 14], the shoulder [15] and the dominant wrist [16, 17]. Other approaches have not used body-worn sensors at all, but instead relied on other measures, such as motion capture data [18, 19] and visual [20] and acoustic data in office surroundings [21]

In [3] the effectiveness of different body locations for sensors is discussed. It is shown that using two locations (the hip and wrist or thigh and wrist) instead of five does not significantly reduce the activity recognition. The (dominant) wrist is the preferred location when only a single sensor is used.

Some experiments have only been performed in an artificial or laboratory setting [1, 3, 4, 9, 10, 16, 20, 22], thus probably resulting in less robust trained models. In [6, 8] the experiments are performed in a uncontrolled environment. This results in a more natural setting with more noise due to free movement. It also introduces transition-states in the data, because there is no clear distinction between activities. The method of [5] explicitly filters the signals for transitions. Others have provided the test subjects with a measuring device for a longer period of time, some up to a few days [5, 17].

A more objective comparison can be made based on the segmentation and classification technique used. In [3] and [10] extensive comparisons are made between much used algorithms which are widely implemented, such as in the WEKA toolkit [23]. These include the base-classifiers Decision Trees (also used in [6, 7]), Decision Tables, K-nearest neighbors (also used [7, 8]), Support Vector Machines (also used in [9, 24, 25]) and Naive Bayes (also used in [17]). The research of [3] also implemented a number of meta-level classifiers, of which Plurality Voting resulted in the highest accuracy. A very well known and applied algorithm from the field of data mining is  $k$ -means clustering. A drawback of this method is that the number of expected clusters needs to be predetermined. To overcome this shortcoming an iterative implementation can be used, which selects the best clustering. A variant is known as Fuzzy  $c$ -means clustering, in which a data point can belong to one or more clusters, each with a degree of membership. The method of  $k$ -means clustering (or a modified version) is used in [5, 12, 19].

Very much used in machine learning problems are neural networks. \*\*\* cite needed \*\*\* Although these kind of networks typically perform well on spatial-typed problems \*\*\* cite needed \*\*\*, adjustments have been made to apply them for temporal classification. The methods of [6, 16] create a multi-level neural network to generate complex discriminators as activity classifiers. A variant of neural networks known as Kohonen Self Organizing Maps is used in [5] to create a vector codebook, which is further processed by a  $k$ -means clustering.



Due similarity with the problem of continuous speech recognition, Hidden Markov Models have been applied to activity recognition tasks. A layered approach proved to be robust for small changes in [20]. In [26] five activities have been modeled as HMMs. The method of [27] introduces HMMs with Markov Jump Linear Systems to segment a signal, with a fixed number of models. In [11] a layer of HMMs of on top of a layer of static classifiers to estimate the current activity. This has a commonality with [15] in which HMMs are used to capture regularities and smooth activity transitions. The application of [2] constructs separate HMMs for each activity and then joins them in a single model to estimate the likelihood of an activity.

A common property of the above mentioned classification algorithms is that there is no explicit segmentation phase. Data points can be segmented a posteriori, by joining adjacent data points which have the same label. An approach in which the data points are processed per segment include [22], in which motifs over sensor data from actions are discovered. In [1] the signal is segmented (without classification) by applying adaptive threshold values on the energy of the signal. The estimates of the direct density-ratio for probability distributions are compared to create change-points in [28]. By making minimal model assumptions a non-parametric segmentation is obtained, without classification. A further application of probability distributions is discussed in [13], which uses an non-parametric Bayesian method to create an Infinite Gaussian Mixture Model to represent activities. Very interesting is the ability to recognize past and new activities without training. The method of [18] considers the intrinsic dimensionality of poses by using Probabilistic Principal Component Analysis and creates cut-points. This method also is independent of trained models and thus also is a non-parametric approach. The method of [29] creates a piecewise linear representation of the signal as thus approximates it. This does not result in a segmentation in the sense of some homogeneity over the data points, but it can support one. In contrast, the method of [14] splits and merges segments the signal in a constant number  $k$  segments using a cost function to minimize the heterogeneity of a segment.

Each research outcome is defined in some kind of correctness to segment and classify the signal in the correct and preferred manner. In many methods an error was measured as the difference between the algorithms outcome and a domain-experts ground truth. In case of segmenting the time series, the method of [1] measures the delayed number of seconds before a new segment is created. For classification the measure takes on many forms and names, such as recognition rates [3] or hit precision [13]. Confusion matrices often are used to give insight in which activities are hard to discriminate from each other, as for instance in [6]. The term false alarm rate is often used to indicate incorrectly new introduced elements, such as activities in [13] and [28].

Refs.	#	Location	Activities	Env.	Algorithms	Results
[22]	Single	Wrist	Signal motifs	Artificial data	Motif discovery, HHMs, DTW	overall 87%
[1]	Multiple	wrist, legs, waist	12 daily routine	Controlled	Adaptive threshold, segmenting	85%
[26]	Single	Unknown	walking, upstairs, falling, running, standing	Controlled	FFT, HHMs	90% - 100%
[28]	-	-	-	Artificial data	Non-parametric Density-Ratio estimation	???
many more...						

TABLE 1.1: Overview of researches with the focus on segmentation and classification techniques.

The overall results of many methods are fairly high and promising. Many methods, for a fixed number of activities or an iterative approach, obtain precision up to the range from 95% to 100% [6–9, 12, 22, 26]. Problems for which the number of activities is unknown perform up to 95% [13, 18]. The most problems arise from discriminating between similar activities, such as taking the stairs up or down and walking when the movement of legs or waist is considered [6, 7].

Some researchers have made their used data set explicitly public available and other data sets are just available, providing valuable labeled activity patterns. \*\*\* This should be in experiments setup? \*\*\* \*\*\* List here the datasets used \*\*\*

### 1.1.2 Temporal Segmentation

Many different approaches have been applied to segmenting time series data. Due to the nature of this research, only segmenting techniques with an on-line approach (or easy translation to one \*\*\*\*?) will be considered. Interesting are algorithms which have been applied to multi-dimensional data, since we are considering segmenting accelerometer data provided by three-axial sensors. Nonetheless, online segmenting algorithms on single dimensional data will also be considered.

\*\*\* Piecewise linear approximation \*\*\*

Segmentation methods can roughly be categorized in three methods in the way the data is processed, as discussed by Avci *et al.* [30]:

- **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. These methods process the data points recursively, which results in a complexity of  $O(n^2k)$ .

- **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. Given the average segment length  $L$ , the complexity of this method is  $O(nL)$ .
- **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. Because of the simple approach, the data is only processed very locally which can yield in poor results [29]. The complexity is equal to the bottom-up approach,  $O(nL)$ , where  $L$  is the average segment length.
- **Sliding Window And Bottom-up (SWAB)**, as introduced by Keogh *et al.* [29], joins 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, and pass this to a buffer with space for a few segments. The buffer then processes the data Bottom-up and returns the first (left-most) segment as final segment. Because this second stage retains some (semi-)global view of the data, the results are comparative with normal Bottom-up. It is stated by Keogh *et al.* that the complexity of SWAB is a small constant factor worse than that of regular Bottom-up.

It is clear that for the application of this research sliding-window and preferably SWAB-based algorithms should be considered.

The SWAB method proposed by Keogh *et al.* [29] is dependent on an user setting, providing the maximum error when performing both stages. Each segment is approximated by using piecewise linear representation (PLR), an often used method. The user 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.* [1] and an adaptive CUSUM-based test by Alippi *et al.* [31], 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.* extracts features from the raw sensor signal to base the segmentation on other properties than the pure values. The method of Bernecker *et al.* [32] 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.* [33] 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). Whereas the before mentioned algorithms use a linear representation, this method 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 [34], 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 a stochastic process \*\*\* is this a correct term? \*\*\*. Probabilistic models can be constructed on that signal, yielding in probabilistic and Bayesian based segmentation methods. The CUSUM-method relies on the log-likelihood ratio 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.* [28] proposes to estimate the ratio of probability densities (known as the *importance*) directly, without explicit estimation of the densities. They claim this results in a robust approach for real-world scenarios. Although this is a model-based method, no strong assumptions (parameter settings) are made on the models.

Another application of PCA is to characterize the data by determining the dimensionality of a sequence of data points. The proposed method of Berbič *et al.* [18] 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.

The method of Adams and MacKay [35] 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 [36] 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.

\*\*\* TODO \*\*\* aan het einde van de lopende tekst richting-zinnen maken, die het einde aanduiden of een erg beknopte samenvatting geven / conclusie van literatuuronderzoek.

### 1.1.3 List of papers

The following is a list of used papers. This list is for my own personal convenience. It is possible this will be in tabular form in the final thesis document.

\*\*\* NOTE: this list will not be in the final thesis (in this form at least) \*\*\*

- “Discovering characteristic actions from on-body sensor data” [22], cited: 63. Motifs, HHM, DTW, 96%, overall 87% accuracy. Discovers motif (actions) from data. Worn on wrist, single. Environment artificial.
- “Recognition of human activities using layered hidden Markov models” [20], cited: 3. HMM, layered (primitive and abstract actions). Uses vision for workplace activities. Direct classification. No sensors. Artificial environment for region labeling.
- “Accelerometer-Based Gait Analysis, A survey”, [24], cited: 4. Compares methods to distinguish walking of normal, fast, slow. Focus on gait, but compares classification methods such as SVM, PCA, KSOM.
- “An Automatic Segmentation Technique in Body Sensor Networks based on Signal Energy”, [1], cited: 13. Automatic segmenting, adaptive threshold. Pure segmenting, no classification. Nice and clean method? Weakness: single action segmented as multiple. Used multiple sensors.
- “Towards HMM based Human Motion Recognition using MEMS Inertial Sensors”, [26], cited: 13. Uses HHM, Fourier transform for features. 5 fixed activities, trained HMM. Correct rates from 90% to 100%. Classification, not just segmentation. Single sensor.
- “Change-Point Detection in Time-Series Data by Direct Density-Ratio Estimation”, [28], cited: 41. Non-parametric approach, online method, no strong model assumptions. Uses multiple datasets, origin referenced. Much follow-up work done. Artificial datasets. (no sensors)

- “Unsupervised, Dynamic Identification of Physiological and Activity Context in Wearable Computing”, [5], cited: 106(!). Combination of KSOM, iterative K-means clustering, transient states removal (using Markov model). On-line algorithm. Single sensor on upper arm. Tested for a few days and controlled environment.
- “Activity Recognition using Cell Phone Accelerometers”, [6], cited: 122 (!), 2011. Uses WEKA: decision trees, logistic regression and multilayer neural networks. Results up to 98% (jogging). Much difficulty with stairs up and down. Single accelerometer, front leg pocket. Offline method. Many recent references. Controlled environment.
- “Activity recognition from accelerometer data”, [10], cited: 434 (!), 2005. Compares 18 (offline) classifiers, base level (WEKA) and meta-level. Single device, pelvic region. Controlled setting.
- “Activity recognition from user-annotated acceleration data”, [3], cited: 1008 (!!), 2004. Five bi-axial sensors. Decision tables, instance-based learning, C4.5 (trees, highest accuracy 89.30%) and naive Bayes (weka), twenty activities. Controlled environment, common room.
- “Using Hierarchical Clustering Methods to Classify Motor Activities of COPD Patients from Wearable Sensor Data”, [4], cited: none, 2005. Uses Linear Discriminant Analysis for cluster classification (using Cluster Quality Index to determine k) and simple rule-based separation for high level ambulatory. Hierarchical Dendrogram to merge clusters when similarity to high. Sensors on forearms and legs. Controlled clinical environment.
- “Non-Parametric Bayesian Human Motion Recognition Using a Single MEMS Tri-Axial Accelerometer”, [13], cited: none, 2012. Recognizes the number of human activities, single sensor on the chest. No training data. Infinite Gaussian Mixture Model, collapsed Gibbs sampler. Compares with parametric Fuzzy C-mean (data point belongs to multiple clusters), unsupervised K-means, non-parametric mean-shift. Outperforms all significantly, high hit rate and low false alarms.
- “An Online Algorithm for Segmenting Time Series”, [29], cited: 487, 2001. Reviews algorithms to get a piecewise linear representation, proposes a novel sliding-window and bottom up approach, on-line. Mere segmentation. No classification; only approximation of signal. (Dataset available).
- “Segmenting Motion Capture Data into Distinct Behaviors”, [18], cited: 242, 2004. On-line methods: cut-points on increased intrinsic dimensionality, distribution of

poses is observed to change. Batch: cut when consecutive frames belong to different Gaussian mixture models. Fixed fourteen motions, compared with manual segmentation. Uses PCA and Probabilistic PCA (models the non-pc subspaces as noise). PPCA is the best. (paper anne)

- “Aligned Cluster Analysis for Temporal Segmentation of Human Motion”, [19], cited: 54, 2008. Uses extension on kernel k-means clustering and Dynamic Time Alignment Kernel (kernel of DTW) for temporal invariance, robust temporal matching metric. Coordinate descent algorithm solves ACA. Fixed number of clusters, trapped in local minima. (paper anne). Uses motion capture data.
- “Bayesian Nonparametric Methods for Learning Markov Switching Processes”, [27], cited: 10, 2010. Uses HMM for state-space model for segmentation, Markov Jump Linear systems. Unbounded number of Markov modes (parameters). Number of models is fixed, though? (paper anne). No sensors.
- “Layered Representations for Human Activity Recognition”, [21], cited: 215, 2002. Layers of Hidden Markov Models HMM, multiple levels of granularity (Based on intuition) and context. Less retraining, only lower models. Classified also. No sensors, visual and akoustic data in office surroundings.
- “Time series segmentation for context recognition in mobile devices”, [14], cited: 151, 2001. Splits and merges segments, while keeping k constant using cost function on segments for internal heterogeneity. Multiple instances determine number of k. Microphone and accelerometer data, in front of chest. Artificial data
- “A practical approach to recognizing physical activities”, [11], cited: 259, 2006. Uses method of [15], activity classification algorithm. Selects most useful features and then recognizes walking, sitting, etc. First layer static classifier on features or data (energy, mean, variance, correlation, etc), then a layer of HMM to estimate activity. All offline. Shows trained model is robust to location of wearing on body.
- “A hybrid discriminative/generative approach for modeling human activities”, [15], cited: 252, 2005. Used in method above. Combines boosting to select and reduce useful features and learn static classifiers with HMM to capture regularities and smooth activities (quick switching is unlikely). Single sensor, multiple measures (accelerometer, audio, etc), worn on the shoulder.
- “Using acceleration measurements for activity recognition: An effective learning algorithm for constructing neural classifiers”, [16], cited: 67, 2008. Separates dynamic from static activities. Uses multilayer feedforward neural networks to generate complex discriminating surfaces as activity classifiers. Feature subset

selection approach is developed. Neural pre-classifier with constant threshold criterion. Uses Common Principal Component Analysis. Recognizes eight activities with 95% overall accuracy. Laboratory environment. Worn on dominant wrist.

- “Single-accelerometer-based daily physical activity classification”, [17], cited: 46, 2009. Uses Naive Bayes classifier, sensor worn on wrist. Five activities. PCA to reduce dimensions and create independence for NB. Comparable results as Decision Trees, but with flexibility to add activities. Performed outside laboratory, 10 hours long.
- “Distributed Continuous Action Recognition Using a Hidden Markov Model in Body Sensor Networks”, [2] cited: 17, 2009. Single HMM, sensor network cluster movements, HMM constructs continuous actions using postures and actions, merges left-to-right HMM for each action to a single. Transcript generation per sensor uses Gaussian Mixed Model, multiple models with  $m$  mixtures are trained with Expectation-Maximization (EM), best is chosen.
- “Offline and online activity recognition on mobile devices using accelerometer data”, [7], cited: none, 2012. Compares offline and online methods, cellphone in pocket. Reduces feature extraction (mean and standard deviation (for resources of cellphone)). Training is done by database and labeling, classification with K-nearest neighbors, decision trees (C4.5) and decision rules. Online and offline KNN-3 is best, 99% and 97%. Six daily activities. Trousers pocket.
- “Recognizing human activities user-independently on smartphones based on accelerometer data”, [8], cited: none, 2012. Classifiers K-nearest neighbors and QDA (quadratic discriminant analysis). Online recognition implemented on phone. five daily activities. Model training by decision tree (active/inactive). Results for QDA, online: 95%. Activities performed outside. Trousers pocket.
- “Activity recognition from acceleration data based on discrete cosine transform and SVM”, [9], cited: 25, 2009. Uses Discrete Cosine Transform (DCT) to get features, PCA to reduce them an multiple SVMs to classify windows. Four daily activities, 97% precision. Trousers pocket. Laboratory setting.
- “Physical Activity Recognition Using a Single Tri-Axis Accelerometer”, [12], cited: 16, 2009. Five daily activities, FFT, fuzzy c means classification, offline. 99% accuracy. No smartphone, other device. Outside, although restricted (standing, sitting, lying, walking, running)



### 1.1.3.1 Pure segmentation

- “Bayesian Online Changepoint Detection”, [35], cited: 67, 2007. Uses a priori and posteriori distribution prediction of the current run-length. Online. Uses well-drilling data, Dow Jones return value and coal-mine accidents rate, which are all used by others. Does not compare results. Decouples algorithm from model (?).
- CUSUM Method: piece wise segments of gaussian mean with noise. More refs needed.
- “An Automatic Segmentation Technique in Body Sensor Networks based on Signal Energy”, [1], cited: 13. Automatic segmenting, adaptive threshold, standard deviation, niet bayesian. Pure segmenting, no classification. Nice and clean method? Weakness: single action segmented as multiple. Used multiple sensors.
- “Segmenting Motion Capture Data into Distinct Behaviors”, [18], cited: 242, 2004. On-line methods: cut-points on increased intrinsic dimensionality, distribution of poses is observed to change. Batch: cut when consecutive frames belong to different Gaussian mixture models. Fixed fourteen motions, compared with manual segmentation. Uses PCA and Probabilistic PCA (models the non-pc subspaces as noise). PPCA is the best. (paper anne). From computer graphics / motion capture data.
- “Activity recognition using inertial sensing for healthcare, wellbeing and sports applications: A survey”, [30], cited: 22, 2010. Discusses segmentation as form of sliding window, top-down, bottom-up or swab. Interesting overview of five steps: preprocessing, segmentation, feature extraction, dimensionality reduction and classifiers. Classifiers only with fixed number of activities (no non-parametric).
- “Event Detection From Time Series Data”, [36], cited: 291, 1999. Iterative (online, incremental) method to fit model to segment, likelihood to determine if partitioned further (number of change points), when batch. Maximum likelihood methods. Also model selection, loss and risk functions. With incremental (online) method, every new data point had likelihood criteria of being a new segment or current one. If no segment for long time, computations become increasingly expensive. Solution: sliding window for last  $w$  points only.
- “An adaptive approach for online segmentation of multi-dimensional mobile data”, [33], cited: 1, 2012. Adaptive model for selecting suitable dimensions to build model and select model (linear, quadratic, cubic). Multiple metric for results evaluation; traditional regression error, information retrieval (precision, recall, F-measure) and segmentation delay. Uses three-axis accelerometer data. Extends the SWAB method. Natural setting, no restriction on body location.

- “Activity Recognition on 3D Accelerometer Data (Technical Report)”, [32], cited: 0, 2012. First makes distinction between periodic and non-periodic activities, directly on accelerometer data. Autocorrelation as a measure for self-similarity (periodicity). After segmentation, features are extracted for clustering. Then dimensionality reduction and (re)classification. Interesting post-processing; small segments with neighbors of the same activity are merged. Removes small mis-activities.
- “Change-Point Detection in Time-Series Data by Direct Density-Ratio Estimation”, [28], cited: 41. Non-parametric approach, online method, no strong model assumptions. Uses multiple datasets, origin referenced. Much follow-up work done. Artificial datasets. (no sensors) Uses likelihood/density ratios estimator, not the densities themselves. Probability distributions that generate the time series.
- “Adaptive Filtering and Change Detection” (book), [37], cited: 682 (!!), 2000. Lots of algorithms, hard to read.
- “Online segmentation of time series based on polynomial least-squares approximations”, [34], cited: 17, 2010. Creates higher-order polynomial model of data, sliding and growing window. Works well for traditional time series (stocks), but not for non-stationary such as daily life.
- “Novel online methods for time series segmentation”, [38], cited: 17, 2008. (Step-wise) Feasible Space Window. Compares with SWAB. Referenced in [33], but that one choses SWAB because it should be lower time complexity and fast. (S)FSW creates less segments, while SWAB has lower error. FSW is fastest and reliable (due to this paper self).
- “An adaptive cusum-based test for signal change detection”, [31] cited: 16, 2006. CUSUM-test with adaptive parameter settings. Can detect non-stationary processes by drift of signal. Applicable for extracted features, e.g. mean.

## 1.2 Problem statement

## Chapter 2

# Techniques

### 2.1 Signal pro-processing and sensor fusion for timed patterns

Describe how data is collected and processed, the form of the data (continuous to discrete, amplitude, energy, frequencies, etc). Include feature extraction, like PCA, FFT, energy, mean, but discuss it further in subsection. Only the concepts will be explained and the scheme of the computations (or simply formulas), no precise implementations. Mention only shortly with meaning, no further explanation.

---

#### 2.1.1 Signal fusion

Signal fusion is a broad field of research and application and many different interpretations exist. The function and interpretation used in this thesis follows the definition from [39], stating:

**Sensor fusion** is the combining of sensory data or data derived from sensory data such that the resulting information is in some sense better than would be possible when these sources were used individually.

This can essentially be considered as a form of synergy, in which the combination of the sensors provide more information than the sum of the raw output. For the application of this thesis we will focus on direct fusion of raw data, which can be applied to both homogeneous and heterogeneous sensors. The fusion stage will combine the data into

a single representation, which can be interpreted by the controlling system. This differs from mere multisensor integration, in which the controlling system is responsible for processing and interpreting the data from different sources. Figure 2.1 shows the difference.

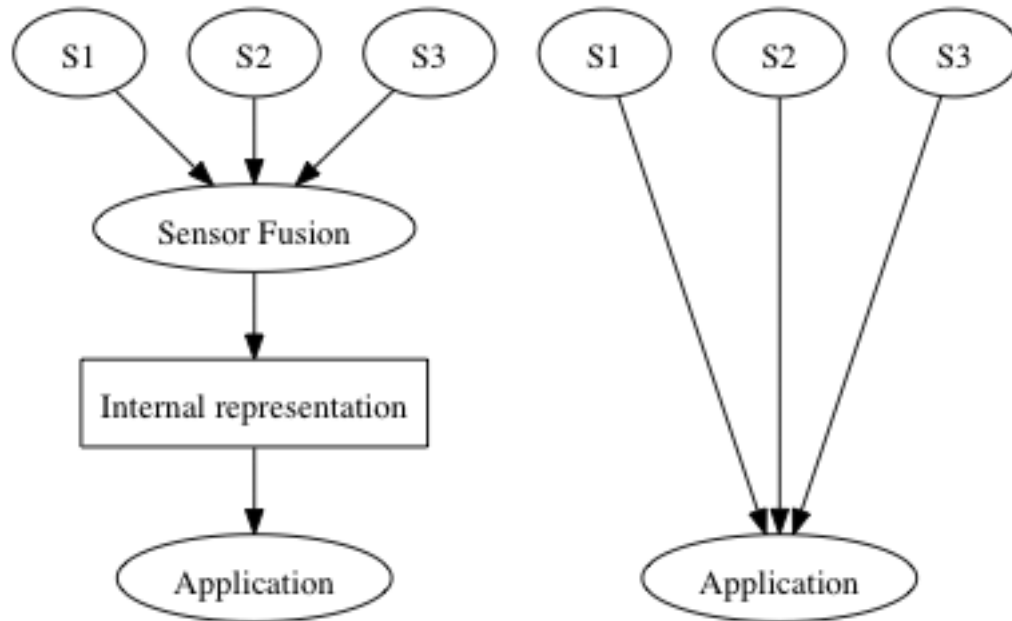


FIGURE 2.1: On the left *Sensor Fusion* and on the right *Multisensor integration*

Using the single representation, the computational and system complexity can be reduced. The representation forms a standardized and normalized interface from the sensors to the system, which create independence and modularity of sensors. Other benefits as a result from sensor fusion include [39]:

- **Robustness en reliability:** When the system incorporates multiple sensors as redundancy, the system becomes more stable in case of partial failure.
- **Extended spatial and temporal coverage:** With multiple sensors a broader spectrum can be analyzed, both in the case of homogeneous and heterogeneous combinations.
- **Increases precision:** Sensors typically have a range and domain of measurements. By combining a set of heterogeneous sensors, the overall range of observation is increased.
- **Reduced ambiguity and uncertainty:** In case of high uncertainty about a observation, the observation of other sensors can confirm it.

Sensor fusion can be applied on multiple levels. A rough diversion over three levels is:

- **Low-level:** Raw sensor data is combined which is considered to be more informative than the original sources.
- **Intermediate-level:** Also known as *feature fusion* combines the features which can be extracted from the data. See also section 2.1.3. This data would be applicable to segmentation \*\*\*[???].
- **High-level:** Fusion of this level will support decision making, optionally followed by actions taken by the system to the environment or conclusion being drawn from the environment.

### 2.1.2 Signal smoothing

What, why and how.

### 2.1.3 Feature extraction

What, why and how.

## 2.2 Temporal Segmentation

This section will give an introduction and in-depth analysis of temporal segmentation.

### 2.2.1 Aims of segmentation

When processing and analyzing time series of data, e.g. motion measurements, stock market fluctuations or natural language, first a low-level division between the discriminative parts of the stream must be made. One can view this as splitting the series into the *atoms*, which are the building blocks of the total stream. These building blocks will be the aggregation of non-overlapping, internally homogeneous segments [14]. This means that the data points inside a segment should have some resemblance relation to each other and their difference lies between some boundary. The process of segmenting can be viewed as a subproblem to context analysis of time series. Temporal segmentation is closely related to temporal clustering, although it is a stricter, and simpler, process. Whereby clustering only restricts the data points on their distance relation (as used in a Voronoi diagram), within a segment the data points must also be contiguous.

The task of segmentation can be performed in a manual matter, by cutting and labeling parts of the stream into coherent parts. This would require human (expert) knowledge

and does not yield a clear cut because of ambiguity. With increasing storage abilities and easier motion capture systems, there is a desire for automated systems which perform the segmentation task unsupervised. Some algorithms used have the (often desired) side-effect of also clustering the segments, such that classes of segments can be discovered in the time series. These algorithms would not only be able to make a distinction between walking, sitting and walking, but would also recognize the reappearance of the walking activity.

### 2.2.2 Formal definition

Formally, temporal segmentation is dividing a time series  $s$ , which consists of  $N$  samples  $\mathbf{x}(1), \mathbf{x}(2), \dots, \mathbf{x}(N)$  from  $\mathbf{R}^d$ . Individual *segments* are referenced by  $s(a, b)$ , consisting of the consecutive samples  $\mathbf{x}(a), \mathbf{x}(a+1), \dots, \mathbf{x}(b)$ ,  $a \leq b$ . Let  $s_1 = s(a, b)$  and  $s_2 = s(b+1, c)$  be two segments, then their concatenation is  $s_1 s_2 = s(a, c)$ . A segmentation  $S$  of  $s$  consists of a sequence of  $k$  non-empty segments  $s_1 s_2 \dots s_k = s$ . This notation is adopted from [14].

As stated, informally each segment should be internally homogeneous. This can formally be measured with an cost function  $F$ , indication the heterogeneity of a segment. The overall aim is to minimize the cost  $F$ . The cost of a segment is a function from the data points and the number of data points  $n = b - a + 1$  and is expressed as

$$\text{cost}_F(s(a, b)) = F(\mathbf{x}; n | \mathbf{x} \in s(a, b)) \quad (2.1)$$

The cost of a  $k$ -segmentation  $S$  is the summation of the costs of the  $k$  segments:

$$\text{Cost}_F(s_1 s_2 \dots s_k) = \sum_{i=1}^k \text{cost}_F(s_i) \quad (2.2)$$

With the objective of minimizing the cost function, the optimal  $k$ -segmentation  $S_F^{\text{opt}}(s; k)$  is the segmentation with minimal  $\text{Cost}_F(s_1 s_2 \dots s_k)$  over all possible  $k$ -segmentations.

The cost function, to calculate the heterogeneity of a (set of) segment(s), can be any function. A simple and natural function would be the sum of variances of the segments. The overall cost function would then be

$$\text{Cost}_V = \frac{1}{N} \sum_{i=1}^k \sum_{j=c_{i-1}+1}^{c_i} \|\mathbf{x}(j) - \mu_i\|^2 \quad (2.3)$$

where  $\mu_i$  is the mean vector of data points in segment  $s_i$ .

### 2.2.3 Application in research fields

[CHARACTERISTICS OF HUMAN MOTION] temporal variability, invariance over time, metrics over actions.

[COMPUTER VISION]

[GRAPHICS/VIDEO]

[DATA-MINING]

[MODEL BASED]

To analyze time series it is often preferred to divide the stream in segments of correlated data. After dividing, each segment represent a period in time in which the same activity is performed. Or, stated otherwise, it results in transitions moments between activities.

### 2.2.4 PCA Based Methods

Many fields of research have been active in the unsupervised segmentation of data. Many authors rely on a form of Principal Component Analysis (PCA), as used a.o. in [18]. Often PCA is used to reduce the dimensionality of the data being processed [REFERENCE] by only using the top  $r$  dimensions to describe the data set. It is observed that data series of simple motions have a lower dimensionality then complexer motions. When a simple (repetitive) motion is about to end and fluently transforms in a new motion, there will be a window of time in which a high dimensionality will be present, due to the new motion. After this period of transition, the dimensionality will decrease, since only the new simple motion is present in the window of time. The first algorithm of [18] is based on this principle.

Given a set of data points, a lower dimensional hyperplane can be constructed to which the data points can be projected. This projection on a lower dimension introduces a error to the original position. When the error is fixed, less dimensions are needed for simple motions in which movements of body parts are highly correlated. For segments in which the data points are lesser correlated, e.g. because of transition state, a higher degree of dimensions of the hyperplane is needed to represent the data with equal error degree. When the dimensionality is reduced from  $d$  to  $r \leq d$ , the ratio of error  $E_r$  can be calculated as

$$E_r = \frac{\sum_{j=1}^r \sigma_j^2}{\sum_{j=1}^d \sigma_j^2} \quad (2.4)$$

where  $\sigma_j$  are the singular values as a result from Singular Value Decomposition (SVD), which is closely related to PCA [40].

In [18] the stream of frames is analyzed on cut-frames to find transitions between action. First, for a number of  $k$  frames (e.g. the equivalent of 2.5 seconds) the required dimensionality  $r$  to keep the error  $E_r$  below some threshold  $\tau$  is calculated. This will yield in an error  $e_i$  for the first  $i$  frames. When more frames are added to the window, the error will increase with a low constant when it is still in the same activity, due to noise in the activity. When a new activity starts, e.g. at frame  $j$ , the error at frame  $j$  will start increasing faster. This can be expressed by the derivative of the error rate  $d_i = e_i - e_{i-l}$ , where  $l$  is a constant to remove noise. From this derivative the mean and standard deviation can be calculated, for each point. When a derivative  $d_j$  rises more than a factor  $k_\sigma = 3$  standard deviations from the mean, a transition point is encountered. The previous frames are then cut from the sequence (as a segment) and the algorithm starts over.

\*\*\* [figure to illustrate derivative and standard deviation]

A second approach in [18] uses the probabilistic variant of PCA (PPCA) to model the data set as a Gaussian distribution instead of ignoring the frames which do not fit in the subspace. Over windows of frames the mean and variance are calculated. In a forward manner the Mahalanobis distance of a new window of frames is calculated, which represents the likelihood of the new window belonging to the same segment as the original window. When the distance decreases, the likelihood increases which happens when the motions in the becomes more homogeneous. When a peak in the distance is reached, the new window of frames indicates a heterogeneous collection of motions in the window and thus a low likelihood of membership and a indication of a transition. In order to distinct activities and sub-activities (which require a subset of motions is a distinct activity) the algorithm is also processed backward over the data series.

\*\*\* [figure to illustrate mahalanobis distance measure and peaks]

The third algorithm in [18] is based on the observation that data points (frames) tend to form clusters in the space. These clusters are represented by  $k$  Gaussian distributions for which each the Expectation-Maximization (EM) algorithm estimates the mean  $m_j$ , covariance matrix  $\sum_j$  and prior  $\pi_j$ . With all the Gaussian distributions estimated, the data points are assigned to the cluster with the highest membership likelihood. When two consecutive frames  $x_i$  and  $x_{i+1}$  belong to different clusters, a transition of activities is recognized. Note that this algorithm succeeds in segmenting the data and also labels the similar simple activities.



A drawback in this system, and many others which implement a variant of the  $k$ -means algorithm, is that the number of clusters  $k$  need to be predetermined. To cope with this, often the algorithm is performed multiple times for different values of  $k$ . Using some criteria, e.g. the Bayesian Information Criterion [41] or the Davies-Bouldin Index which guides  $k$ -means clustering as used in [5].

### 2.2.5 Statistical methods

An method to process data, or create a basis for further processing, is by analyzing the statistical properties of a data set. When there is a continuous stream of data, it is possible to extract events by comparing the statistical properties of sliding time windows. By definition, events have a different statistical profile then their background. When two consecutive actions are regarded as the background, then a transition is an event between them.

In [1] a automatic system is used to segment a continuous stream of data. By calculating features as activity level from the standard deviation of a sliding window and applying Adaptive Thresholds, the data can be segmented in active and rest parts. There is no model representation.

### 2.2.6 Hidden Markov Models Based Methods

### 2.2.7 Bayesian Methods

### 2.2.8 Principal Component Analysis

When working with simple 2- or 3-dimensional data it is often, for humans, easy to discover patterns in the set. The data can be plotted and the lines or planes among which the data points lie gives an indication of the pattern. With Principal Component Analysis (PCA) this process is also possible for automated systems. With this method the overall form of a set of data points can be represented, the most discriminative dimensions can be found, the dimensionality of a set can be reduced (which yield in data compression) or the result can be used to characterize and differentiate sets of data points.

The following steps are performed, which are explained below and illustrated in figure \*\*\* [add figure with plot of sets]:

1. **Gather data** and represent them in a chosen number of features,

2. **Subtract mean** to center the data around the axis. The new data set will have mean zero,
3. **Generate covariance matrix** by calculating all pairwise feature variances,
4. **Calculate Eigenvectors and Eigenvalues** of the covariance matrix. The Eigenvectors are then normalized to make further calculations easier,
5. **Generate feature vector** which indicated which features are characteristic or meant to keep in the data set, by comparing the Eigenvalues.

The workings of the PCA [42] relies on the concepts of standard deviation, variance, covariance, eigenvectors and eigenvalues of matrices. These concepts will be discussed very briefly <sup>1</sup>. The standard deviation and variance of a set are measures for the spread around the mean for a single dimension, or feature. The covariance between two features (dimensions of the data points) indicates how they are related; when positive the two features will increase together; when negative one will decrease when the other increases and when zero they are unrelated. The covariance matrix gives all the covariances for all pairs of features. This matrix is symmetrical about the diagonal and the values on the diagonal are the variances for each feature. For this matrix the eigenvectors can be computed. Eigenvectors have the characteristic that when they multiply a matrix, the resulting vector is a multiple of the Eigenvector. The amount by which it is the multiple is the Eigenvalue. This is illustrated in formulae 2.5 and 2.6. The second formula shows an Eigenvector  $\begin{pmatrix} 3 \\ 2 \end{pmatrix}$  and its associated Eigenvalue, 4. The vector in the first formula is not an Eigenvector.

$$\begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \times \begin{pmatrix} 1 \\ 3 \end{pmatrix} = \begin{pmatrix} 2 \cdot 1 + 3 \cdot 3 \\ 2 \cdot 1 + 1 \cdot 3 \end{pmatrix} = \begin{pmatrix} 11 \\ 5 \end{pmatrix} \quad (2.5)$$

$$\begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \times \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} 2 \cdot 3 + 3 \cdot 2 \\ 2 \cdot 3 + 1 \cdot 2 \end{pmatrix} = \begin{pmatrix} 12 \\ 8 \end{pmatrix} = 4 \times \begin{pmatrix} 3 \\ 2 \end{pmatrix} \quad (2.6)$$

It are these Eigenvectors and Eigenvalues that makes PCA possible. Given a matrix of size  $n \times n$  (only square matrices have Eigenvectors),  $n$  Eigenvectors can be found. Each Eigenvector is perpendicular, or orthogonal, to the others. The Eigenvectors combined describes the lines over which the data is plotted. The Eigenvector with the largest Eigenvalue is the *principal* component of the data set. It is said that it is the most significant feature. Thus the Eigenvectors, and thereby the features or dimensions of the data set, can be sorted on significance by the Eigenvalues.

---

<sup>1</sup>For a more in-depth discussion we would like to refer the reader to [43]

With this ordering there are a few applications possible. The first is to just make a (comprehensible) representation of the data. The Eigenvectors describe the cloud of data points and thus are a compressed representation of the data. When the original data points are to be compressed, it is possible to remove the least significant features and reconstruct the data points from the resulting Eigenvectors. This will yield in a lossy compression. \*\*\* [figure to illustrate]. An extension of this application is to determine the number of features needed to keep the compression within a certain error criterion, as used in [18]. Sets of data can then be distinguished by the number of features needed to describe the points.

\*\*\* [graph, a bit of formulae with sets and dimensions]

### 2.2.9 Hidden Markov Models

Many dynamic systems can be viewed as some sequence of consecutive states which are observable by the world. These kind of systems are known as Markov Chains. When the state in which the system is can not be observed with complete certainty, a system known as a Hidden Markov Model (HMM) can be constructed. It has been shown that HMMs are effective on in the task of speech recognition [44]. Because of some resemblance between natural speech construction and human activities, modifications have been able to classify human activities [2] and extensions are implemented which create layers of HMMs to handle different granularity over time and reduce retraining [20, 21]. A HMM is a state-space model representing a single model for each label which needs to be classified.

Formally, a Hidden Markov Model is characterized by the following elements [44]:

1. A number of  $N$  states in the model. These states are a representation of the underlying mechanism of the system and can not be directly observed from outside (hence the *hidden* model). Individual states are referenced by  $S = \{S_1, S_2, \dots, S_N\}$  and a single state in time  $t$  is denoted by  $q_t$ .
2. An alphabet of  $M$  discrete observable symbols for each state. An observed symbol is the only indication of the systems current situation. Each individual symbol is referenced by  $V = \{v_1, v_2, \dots, v_m\}$ .
3. A state transition probability distribution  $A = \{a_{ij}\}$  indication the probability of being in state  $j$  at time  $t + 1$  when the state on time  $t$  was  $i$ :

$$a_{ij} = P[q_{t+1} = S_j | q_t = S_i], \quad 1 \leq i, j \leq N \quad (2.7)$$

Note that by setting the probability  $a_{ij} = 0$  for some  $i$  and  $j$  indicates that  $j$  is not reachable from state  $i$ .

4. A symbol observation probability distribution in state  $j$ ,  $B = \{b_j(k)\}$ , where

$$\begin{aligned} b_j(k) &= P[v_k \text{ at } t | q_t = S_j], \quad 1 \leq j \leq N \\ &\quad 1 \leq k \leq M \end{aligned} \quad (2.8)$$

Note that every symbol can be observed in every state.

5. The initial state distribution in which the system starts  $\pi = \{\pi_i\}$  where

$$\pi_i = P[q_1 = S_i], \quad 1 \leq i \leq N \quad (2.9)$$

This system is often abbreviated to

$$\lambda = \{\pi_i, a_{ij}, b_j(k)\} \quad (2.10)$$

The HMM will act as a generator giving a sequence of  $T$  observations, each being a symbol from  $V$ :

$$O = O_1, O_2 \dots O_T \quad (2.11)$$

Given the characteristics of a specific HMM, there are three major application which serve the classification of an activity. The first application is, given a sequence of observations  $O = O_1, O_2, \dots, O_T$ , to label this sequence with the activity with the highest associated probability for that label, e.g. compute  $P(O|\lambda)$  for each model  $\lambda$ . The second task a HMM can perform is to *explain* a observation sequence  $O$ , by choosing an optimal corresponding state sequence  $Q = q_1, q_2, \dots, q_T$ . The last application is more of a problem statement, being the task to train the HMM to adjust the parameters  $\lambda = \{\pi_i, a_{ij}, b_j(k)\}$  to maximize  $P(O|\lambda)$ . In the overall task of activity recognition, one would first use the last application to train a set of HMMs with the provided and labeled training data. The second application can be used to further optimize the HMM (set the number of states, adjust vector codebook etc.) and improve the capability of modeling the activities. Finally, when the system is provided with unlabeled data the first application is used to identify the observed sequence with trained examples, by comparing the probability result of multiple HMMs.

Generally a HMM can be fully connected, or ergodic model. Other configurations are also possible, of which the left-right model, or Bakis model, is popular for temporal applications [44]. This because the structure of the state probability distribution resembles the construction of possible states over time. This left-right model can be generalized

by setting

$$a_{ij} = 0, \quad j > i + \Delta \quad (2.12)$$

for some  $\Delta$ , indicating the possible *distance* of the jumps over states.

\*\*\* [INSERT GRAPHICS SHOWING STRUCTURE OF HMM]

\*\*\* [discrete and continuous HHMs]

\*\*\* Applications of HHM: [26] uses 10 coefficients extracted with FFT from 6 axis (3 times accelerometer and 3 times gyroscope). Example of continuous HMM (?). Recognition rates fo 50%-100%.

In [21] it is stated that HHMs are robust mechanisms with respect to changes in temporal segmentation of observations, although they suffer from a few problems when applied to reason about longer and complexer sequences over time. They state that, especially with limited training data, the HHM tend to be overfitted and tend to lack structure and have an excess of parameters. To overcome this limitation, they introduce a Layered Hidden Markov Model (LHHM) which they show is more robust to temporal changes in observations. They make use of the hypothesis from the field of psychology, that human behavior is hierarchically structured [45]. In [20] a layered construction of HHMs is used which feeds the low-level data from a sliding window to the lower layer. This layer will recognize simple motions, such as moving the arm in a pre-defined segment. The results from these lower HHMs are then fed to a higher-level layer of HHMs which are able to recognize abstract motions, such as adjusting a monitor and picking up a pen. The results showed that the layered system was better at recognition than a similar single-layered setup. It was also robust in the sense that when the environment was changed, only the relatively simple low-level HHMs needed to be retrained.

\*\*\* [DIAGRAM WITH EXAMPLE OF LOW- AND HIGH-LEVEL HHMs, SUCH AS MOVING ARM, PICKING UP PEN]

\*\*\* Viterbi algorithm / path; find the most likely sequence of cause-states for the observed measures.

### 2.2.10 Segmentation as clustering

In the previous section the discussed methods all relied on the stream of data points and tried to find cuts the discriminate between successive different type of activities. Another approach is to consider the tasks of segmentation as a type of clustering [19]. In clustering the objective is to assign labels, or classes, to all the data points indication a

similar type of activity. A clustering  $\mathcal{L}$  is thereby more informative than a segmentation but is also harder to produce.

A clustering  $\mathcal{L}$  is generated from a sequence of elements  $\mathbf{X}$  which is decomposed in  $m$  disjoint segments, each belonging to one of the  $k$  classes. A segment  $\mathbf{Y}_i \hat{=} \mathbf{X}_{[s_i, s_{i+1})}$  is composed of frames from position  $s_i$  to  $s_{i+1}$ . A vector  $g_{ci} = 1$  indicates class membership if  $\mathbf{Y}_i$  belongs to class  $c$ , otherwise  $g_{ci} = 0$ .

When regarding segmentation of human motion as a task of clustering the difficulty is to model the temporal variability of actions and defining a robust metric between temporal actions. To overcome this, [19] introduces Aligned Cluster Analysis (ACA), by minimizing

$$J_{ACA}(\mathbf{G}, \mathbf{s}) = \sum_{c=1}^K \sum_{i=1}^m g_{ci} \text{dist}_c(\mathbf{X}_{[s_i, s_{i+1})}) \quad (2.13)$$

The characteristic of ACA is that it enables segments to span over different number of data points, whereas the standard kernel  $k$ -means algorithm results in equally sized segments. [!!! TRUE?!!] The second difference is that the kernel used in  $\text{dist}_c$  to measure the distance from a segment to the class which it is assigned to uses the Dynamic Time Alignment Kernel [REFERENCE?] to measure between time series.

## 2.3 Clustering

This section will give an introduction and overview to clustering of data. When processing time series of data the line of distinction between segmentation and clustering is very fine. This section will introduce clustering for the purpose as used in this project.

Data clustering can be used with two different applications, exploratory or confirmatory [46]. In the former, the purpose is to discover clusters in a point cloud of data points. Here the cluster is defined as a group of homogeneous data points measured by some coherence measure, often a distance function with the data points being defined as vectors. The produced clusters, which together form a representation of the data examined, can be used in the latter application to assign new provided data points to a cluster and thereby classify them.

Depended on the precise setup, clustering can be used as a successive step after or an implementation of temporal segmentation. When temporal segmentation is used to create successive coherent data points, clustering can be used to recognize the same activities in different points of time. The data points provided to the clustering will be some representation of each segment and each segment will be labeled with a activity.

Another mechanism could be to extract features from the raw data points and use these to create clusters directly \*\*\* [temporal modifications]. Both mechanisms requires the exploratory and confirmatory phases.

The exploratory phase roughly consists of the following five steps [46]

1. data point representation,
2. definition of coherence,
3. grouping data points,
4. cluster abstraction,
5. output qualification

The first step is to find a *representation* of the data points and clusters, considering the number, type and scale of the features and the number of desired classes to cluster in. When dealing with high-dimensional data points a *feature selection* can be performed to use only the most discriminating features. When the raw features are not useful enough, *extraction* can be used to create new synthetic features. \*\*\* [qualitative versus quantitative/conceptual]

When the data points are represented in a meaningful way, the measure of *coherence* between pairs of points must be defined. Often this is implemented as the Euclidean or Mahalanobis distance when the data points are represented as vectors or other similarity measures for conceptual patterns.

In the most critical step, the *grouping*, the data points with the highest coherence are linked together to form a cluster. The result can be a hard partition over the data or a fuzzy membership degree to each cluster for each data point. When applying rules for merging and splitting cluster, a hierarchical partition is constructed. Partitional clustering algorithms assign data points to clusters by optimizing some criterion, e.g. the mean squared distance from each data point to the clusters' centroid.

Especially in case of large data sets, the resulting clusters will be defined by many data points. To form a compacter and simpler representation, *abstraction* can be used. This will simply analysis by humans or automated systems. A common abstraction is to use the clusters' centroid [47] or parameters of Gaussian patterns.

The final step which can be applied to the resulting partition is some *qualification*. The quality of multiple partitions can be compared and the validity of the partition can be determined. A partition is valid if reasonably it could not been constructed by change

or as some random process. External references of models can be used to compare, or the internal data points can be examined. \*\*\* [better wordings]

### 2.3.1 Formal definition

### 2.3.2 Application is research fields

## 2.4 Temporal pattern recognition

### 2.4.1 Dynamic Time Warping

Used to measure similarity between time sequences. Exact matching is high-cost, so approximations such as Minimum Bounding Rectangles are used.

### 2.4.2 k-means clustering

[EXPLAIN] divide data set  $n$  into  $k$  clusters. Follows the outline in figure 2.2

Among many unsupervised clustering techniques,  $k$ -means is successfully applied to large data sets. It is simple to implement and linear in time complexity so computationally attractive [46]. A drawback of the method is that the results of the algorithm greatly depends on the initial configuration (the data points which will act as centroids) and the number of cluster  $k$  must be determined beforehand.

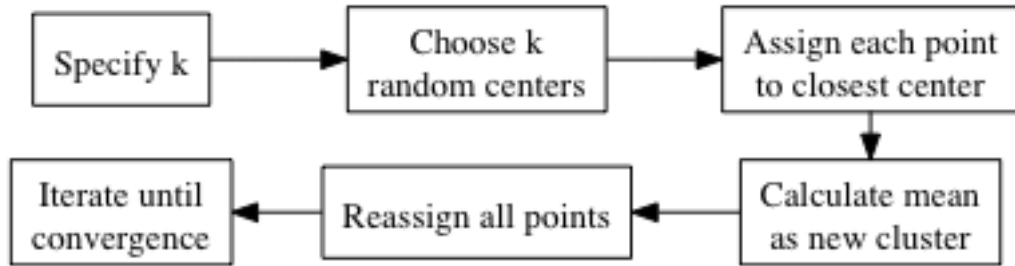


FIGURE 2.2: Outline of the  $k$ -means algorithm.

Generally, the  $k$ -means methods will minimize the squared error for a clustering  $\mathcal{L}$  criterion which is defined as the distance from the data points to the centroid for each cluster in  $\mathcal{K}$ . This is expressed as optimizing to a local optimum the energy function

$$e^2(\mathcal{K}, \mathcal{L}) = \sum_{j=1}^K \sum_{i=1}^{n_j} \|\mathbf{x}_i^{(j)} - \mathbf{c}_j\|^2 \quad (2.14)$$



There are several limitation on the  $k$ -means method. One of these is that only spherical shapes of cluster can be generated. One of the extensions is kernel  $k$ -means [48], which implicitly projects the data points to a higher dimension and thereby is able to form irregular shaped cluster.

### **2.4.3 Self-organizing Map**

### **2.4.4 Support Vector Machine**

### **2.4.5 Naïve Bayes**

## **2.5 Unsupervised clustering of temporal patterns**

## Chapter 3

# Experiments

### 3.1

## Chapter 4

# Results

### 4.1

## Chapter 5

# Discussion

### 5.1

## Chapter 6

# Conclusions

### 6.1

## Appendix A

# Appendix Title Here

Write your Appendix content here.

# Bibliography

- [1] 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.
- [2] 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.
- [3] Ling Bao and Stephen Intille. Activity recognition from user-annotated acceleration data. *Pervasive Computing*, pages 1–17, 2004.
- [4] Delsey Sherrill, Marilyn Moy, John Reilly, and Paolo Bonato. Using hierarchical clustering methods to classify motor activities of copd patients from wearable sensor data. *Journal of NeuroEngineering and Rehabilitation*, 2(1):16, 2005. ISSN 1743-0003. doi: 10.1186/1743-0003-2-16. URL <http://www.jneuroengrehab.com/content/2/1/16>.
- [5] A. Krause, D.P. Siewiorek, A. Smailagic, and J. Farringdon. Unsupervised, dynamic identification of physiological and activity context in wearable computing. page 88, 2003.
- [6] 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.
- [7] 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.
- [8] Pekka Siirtola and Juha Rönning. Recognizing human activities user-independently on smartphones based on accelerometer data. *IJIMAI*, 1(5):38–45, 2012.

- [9] Zhenyu He and Lianwen Jin. Activity recognition from acceleration data based on discrete cosine transform and svm. In *Systems, Man and Cybernetics, 2009. SMC 2009. IEEE International Conference on*, pages 5041–5044. IEEE, 2009.
- [10] Nishkam Ravi, Nikhil Dandekar, Preetham Mysore, and Michael L Littman. Activity recognition from accelerometer data. In *Proceedings of the national conference on artificial intelligence*, volume 20, page 1541. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2005.
- [11] Jonathan Lester, Tanzeem Choudhury, and Gaetano Borriello. A practical approach to recognizing physical activities. *Pervasive Computing*, pages 1–16, 2006.
- [12] Mi-hee Lee, Jungchae Kim, Kwangsoo Kim, Inho Lee, Sun Ha Jee, and Sun Kook Yoo. Physical activity recognition using a single tri-axis accelerometer. *Lecture Notes in Engineering and Computer Science*, 2178.
- [13] M Ejaz Ahmed and Ju Bin Song. Non-parametric bayesian human motion recognition using a single mems tri-axial accelerometer. *Sensors*, 12(10):13185–13211, 2012.
- [14] 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.
- [15] Jonathan Lester, Tanzeem Choudhury, Nicky Kern, Gaetano Borriello, and Blake Hannaford. A hybrid discriminative/generative approach for modeling human activities. In *Proceedings of the 19th international joint conference on Artificial intelligence*, pages 766–772, 2005.
- [16] 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.
- [17] Xi Long, Bin Yin, and Ronald M Aarts. Single-accelerometer-based daily physical activity classification. In *Engineering in Medicine and Biology Society, 2009. EMBC 2009. Annual International Conference of the IEEE*, pages 6107–6110. IEEE, 2009.
- [18] J. Barbič, A. Safonova, 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.



- [19] 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.
- [20] S. Perdakis, D. Tzovaras, and M.G. Strintzis. Recognition of human activities using layered hidden markov models. In *Cognitive Information Processing Workshop*, pages 114–119. Citeseer, 2008.
- [21] N. Oliver, E. Horvitz, and A. Garg. Layered representations for human activity recognition. In *Multimodal Interfaces, 2002. Proceedings. Fourth IEEE International Conference on*, pages 3–8. IEEE, 2002.
- [22] David Minnen, Thad Starner, M Essa, and Charles Isbell. Discovering characteristic actions from on-body sensor data. In *Wearable Computers, 2006 10th IEEE International Symposium on*, pages 11–18. IEEE, 2006.
- [23] Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, and Ian H Witten. The weka data mining software: an update. *ACM SIGKDD Explorations Newsletter*, 11(1):10–18, 2009.
- [24] M.O. Derawi. Accelerometer-based gait analysis, a survey. *Norsk informasjonssikkerhetsskonferanse (NISK)*, 2010.
- [25] 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.
- [26] G. Shi, Y. Zou, Y. Jin, X. Cui, and W.J. Li. Towards hmm based human motion recognition using mems inertial sensors. In *Robotics and Biomimetics, 2008. ROBIO 2008. IEEE International Conference on*, pages 1762–1766. IEEE, 2009.
- [27] Emily B Fox, Erik B Sudderth, Michael I Jordan, and Alan S Willsky. Bayesian nonparametric methods for learning markov switching processes. *Signal Processing Magazine, IEEE*, 27(6):43–54, 2010.
- [28] 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.
- [29] 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.

- [30] 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.
- [31] 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.
- [32] Thomas Bernecker, Franz Graf, Hans-Peter Kriegel, Christian Moennig, Dieter Dill, and Christoph Tuermer. Activity recognition on 3d accelerometer data (technical report). 2012.
- [33] 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.
- [34] 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.
- [35] Ryan Prescott Adams and David JC MacKay. Bayesian online changepoint detection. *arXiv preprint arXiv:0710.3742*, 2007.
- [36] 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.
- [37] Fredrik Gustafsson and Fredrik Gustafsson. *Adaptive filtering and change detection*, volume 1. Wiley Londres, 2000.
- [38] Xiaoyan Liu, Zhenjiang Lin, and Huaqing Wang. Novel online methods for time series segmentation. *Knowledge and Data Engineering, IEEE Transactions on*, 20(12):1616–1626, 2008.
- [39] W. Elmenreich. An introduction to sensor fusion. 2001.
- [40] J. Shlens. A tutorial on principal component analysis. *Systems Neurobiology Laboratory, University of California at San Diego*, 2005.
- [41] D. Pelleg, A. Moore, et al. X-means: Extending k-means with efficient estimation of the number of clusters. In *Proceedings of the Seventeenth International Conference on Machine Learning*, volume 1, pages 727–734. San Francisco, 2000.

- 
- [42] L.I. Smith. A tutorial on principal components analysis. *Cornell University, USA*, 51:52, 2002.
  - [43] I. Jolliffe. *Principal component analysis*. Wiley Online Library, 2005.
  - [44] L.R. Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.
  - [45] J.M. Zacks and B. Tversky. Event structure in perception and conception. *Psychological bulletin*, 127(1):3, 2001.
  - [46] A.K. Jain, M.N. Murty, and P.J. Flynn. Data clustering: a review. *ACM computing surveys (CSUR)*, 31(3):264–323, 1999.
  - [47] E. Diday and JC Simon. Clustering analysis. *Digital Pattern Recognition*, 10:47–94, 1976.
  - [48] B. Schölkopf, A. Smola, and K.R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural computation*, 10(5):1299–1319, 1998.