

Gait sequence modelling and estimation

using Hidden Markov Models



Presented by:
Kouame Hermann Kouassi

Prepared for:
Fred Nicolls
Dept. of Electrical and Electronics Engineering
University of Cape Town

Submitted to the Department of Electrical Engineering at the University of Cape Town
in partial fulfilment of the academic requirements for a Bachelor of Science degree in
Electrical and Computer Engineering

November 12, 2017

Declaration

1. I know that plagiarism is wrong. Plagiarism is to use another's work and pretend that it is one's own.
2. I have used the IEEE convention for citation and referencing. Each contribution to, and quotation in, this report from the work(s) of other people has been attributed, and has been cited and referenced.
3. This report is my own work.
4. I have not allowed, and will not allow, anyone to copy my work with the intention of passing it off as their own work or part thereof.

Signature:.....

Kouame H. Kouassi

Date:.....

Acknowledgments

Abstract

Hidden Markov Models (HMMs) are doubly embedded stochastic processes with a rich underlying statistical structure. In the present work, gait sequence modelling and estimation was performed with HMMs using inertial measurement unit (IMU) data. More specifically, from IMU measurements, a dog's footfalls were correctly identified with up to 95% precision. The continuous-valued IMU measurements were modelled with Gaussian mixtures. The originality of this project lies in the following three reasons.

Firstly, the similar work found in literature only uses 6-dimensional accelerometer and gyroscope measurements [4] [2]. The gait estimation was performed using 18-dimensional data. Thus, feature extraction and feature subset selection techniques were used for dimensionality reduction. The experiments demonstrated that up to 78% performance increase is obtainable by using dimensionality reduction when, the training data is not large.

Secondly, the available IMU dataset was very small in size. Two different methods were employed to overcome this limitation.

In the first method, the reverse gait measurements were appended to the original data to increase the training data. In addition, the IMU measurements of the dog running, walking, and trotting were aggregated. Subsequently, it was possible to build more than 90% accurate HMM model with 18-dimensional IMU measurements with low variance error.

Thirdly, the HMM gait estimation model was designed by building separate models for the front and the back legs of the dog. It, therefore, follows that, the gait estimation algorithm developed in this projet is applicable to both quadrupeds and bipeds.

Finally, the designed HMM models were successfully used to perform motion type recognition. Overall, this undergraduate final year project may be deemed successful.

Contents

1	Introduction	1
1.1	Background to the study	1
1.2	Objectives of this study	2
1.3	Scope and limitations	3
1.4	Plan of development	3
2	Literature review and basic theory	4
2.1	Hidden Markov Models	4
2.1.1	HMM parameters specification	5
2.1.2	The three basics problems for HMM design	7
2.2	Dimensionality reduction	8
2.2.1	Motivations for dimensionality reduction	8
2.2.2	Feature selection: filters and wrappers	9
2.2.3	Feature extraction	12
3	Hidden Markov Model and dimensionality reduction design	14

3.1	Design overview	14
3.2	Description of available dataset	14
3.3	Quadruped Gait sequence modelling with HMM	15
3.3.1	HMM model elements: states and observations properties	15
3.3.2	Splitting the 16-states HMM into two 4-states HHMs	16
3.3.3	State transitions	16
3.4	Pre-processing and increasing the dataset	17
3.5	Dimensionality reduction	18
3.5.1	Feature ranking with separability index (SI)	18
3.5.2	Determining optimal number of features using SI and MCE	18
3.5.3	Feature ranking	20
3.5.4	Forward feature selection	20
3.5.5	Full-search feature selection	21
3.5.6	PCA and LDA	21
3.6	Solving the three basics HMM problems	22
3.6.1	Solution to the training problem	23
3.6.2	Solution to the evaluation problem	27
3.6.3	Solution to the decoding problem	27
3.7	Evaluation of the footfall identification accuracy	27
3.8	Practical implementation	28

3.8.1	Model's parameters estimation implementation	29
3.8.2	Decoding and log-likelihood computation implementation	30
4	Results	33
4.1	Experiment 1: Observation sequence dimensionality's effect on performance	33
4.1.1	Aim of the experiment	33
4.1.2	Experiment apparatus	33
4.1.3	Experiment procedure	34
4.1.4	Experiment results	34
4.1.5	Analysis and discussion of results	35
4.1.6	Conclusions and recommendations of the experiment	38
4.2	Experiment 2: The impact of dimensionality reduction on performance .	39
4.2.1	Aim of the experiment	39
4.2.2	Experiment apparatus	39
4.2.3	Experiment procedure	40
4.2.4	Experiment results	40
4.2.5	Analysis and discussion of results	40
4.2.6	Conclusions and recommendations of the experiment	44
4.3	Experiment 3: The necessity of combining the front and back IMU sensor measurements	44
4.3.1	Aim of the experiment	44

4.3.2	Experiment apparatus	44
4.3.3	Experiment procedure	45
4.3.4	Experiment results	45
4.3.5	Analysis and discussion of results	45
4.3.6	Conclusions and recommendations of the experiment	48
4.4	Experiment 4: Motion type recognition	49
4.4.1	Aim of the experiment	49
4.4.2	Experiment apparatus	50
4.4.3	Experiment procedure	50
4.4.4	Experiment results	51
4.4.5	Analysis of results	51
4.4.6	Conclusions of the experiment	51
5	Discussion	52
5.1	Data dimensionality and bias analysis	52
5.2	Data aggregation and mirroring	53
5.3	State duration time	53
6	Conclusions	54
7	Recommendations	56
A	Additional Files and Schematics	62

A.1	Implementation and experiment code	62
A.2	Additional results	62
B	Addenda	79
B.1	Ethics Forms	79

List of Figures

2.1	Example of a 3-states ergodic HMM	4
2.2	The procedure of filters in dimensionality reduction	9
2.3	The procedure of wrappers in dimensionality reduction	10
2.4	The procedure of hybrid filter-wrapper in dimensionality reduction	11
3.1	Misclassification error vs feature number using separability index and KNN	19
3.2	Optimal PCA component number	22
3.3	Finding best number of mixture component using AIC	26
4.1	The effect of CHMM's observation dimensionality the state sequence decoding accuracy	35
4.2	The effect of CHMM's observation dimensionality the log-likelihood	36
4.3	Scatter plot of 5-dimensional observations grouped according the ground-truth state sequence	36
4.4	Scatter plot of 5-dimensional observations grouped according to the estimated state sequence	37
4.5	Scatter plot of 18-dimensional observations grouped according to the ground-truth state sequence	37

4.6	Scatter plot of 18-dimensional observations grouped according to the estimated state sequence	38
4.7	The effect of training datasize on the prediction accuracy	41
4.8	The effect of training datasize on the log likelihood	42
4.9	Bias-Variance tradeoff analysis	43
4.10	Front footfalls prediction accuracy of both IMUs vs with only the front IMU	46
4.11	Front footfalls prediction log-likelihood with both IMUs vs with only the front IMU	47
4.12	Back footfalls prediction accuracy with both IMUs vs with only the front IMU	48
4.13	Back footfalls prediction log-likelihood with both IMUs vs with only the front IMU	49
A.1	Scatter plot of 1-dimensional observations grouped according the ground-truth state sequence	63
A.2	Scatter plot of 1-dimensional observations grouped according to the estimated state sequence	63
A.3	Scatter plot of 2-dimensional observations grouped according to the ground-truth state sequence	64
A.4	Scatter plot of 2-dimensional observations grouped according to the estimated state sequence	64
A.5	Scatter plot of 3-dimensional observations grouped according the ground-truth state sequence	65
A.6	Scatter plot of 3-dimensional observations grouped according to the estimated state sequence	65

A.7	Scatter plot of 4-dimensional observations grouped according to the ground-truth state sequence	66
A.8	Scatter plot of 4-dimensional observations grouped according to the estimated state sequence	66
A.9	Scatter plot of 6-dimensional observations grouped according the ground-truth state sequence	67
A.10	Scatter plot of 6-dimensional observations grouped according to the estimated state sequence	67
A.11	Scatter plot of 7-dimensional observations grouped according to the ground-truth state sequence	68
A.12	Scatter plot of 7-dimensional observations grouped according to the estimated state sequence	68
A.13	Scatter plot of 8-dimensional observations grouped according the ground-truth state sequence	69
A.14	Scatter plot of 8-dimensional observations grouped according to the estimated state sequence	69
A.15	Scatter plot of 9-dimensional observations grouped according to the ground-truth state sequence	70
A.16	Scatter plot of 9-dimensional observations grouped according to the estimated state sequence	70
A.17	Scatter plot of 10-dimensional observations grouped according to the ground-truth state sequence	71
A.18	Scatter plot of 10-dimensional observations grouped according to the estimated state sequence	71
A.19	Scatter plot of 11-dimensional observations grouped according the ground-truth state sequence	72

A.20 Scatter plot of 11-dimensional observations grouped according to the estimated state sequence	72
A.21 Scatter plot of 12-dimensional observations grouped according to the ground-truth state sequence	73
A.22 Scatter plot of 12-dimensional observations grouped according to the estimated state sequence	73
A.23 Scatter plot of 13-dimensional observations grouped according the ground-truth state sequence	74
A.24 Scatter plot of 13-dimensional observations grouped according to the estimated state sequence	74
A.25 Scatter plot of 14-dimensional observations grouped according to the ground-truth state sequence	75
A.26 Scatter plot of 14-dimensional observations grouped according to the estimated state sequence	75
A.27 Scatter plot of 15-dimensional observations grouped according the ground-truth state sequence	76
A.28 Scatter plot of 15-dimensional observations grouped according to the estimated state sequence	76
A.29 Scatter plot of 16-dimensional observations grouped according to the ground-truth state sequence	77
A.30 Scatter plot of 16-dimensional observations grouped according to the estimated state sequence	77
A.31 Scatter plot of 17-dimensional observations grouped according to the ground-truth state sequence	78
A.32 Scatter plot of 17-dimensional observations grouped according to the estimated state sequence	78

List of Tables

3.1	Main variables of the dataset	15
3.2	Main design components together with toolboxes and functions used for their implementation	29
4.1	Footfall prediction confusion matrix (% accuracy)	51
4.2	Footfall sequence log-likelihood	51

Chapter 1

Introduction

Human motion analysis is an important research topic motivated by medical, artistic or scientific purposes [3]. Among others reasons, it is investigated for human activity monitoring, fall detection, gesture recognition, balance control evaluation, and abnormal behaviour detection [2]. To carry out these studies, micro-electro-mechanical inertial measurement units (IMUs) are increasingly used to capture motion data. Primarily, because they are lightweight and cheap [4] [3] [2].

Gait analysis requires the identification of latent gait states [2]. A gait state being a particular stance in a given gait activity. For this reason, it is extensively used in biologically inspired robot design. More specifically, understanding animals gait mechanism is very invaluable in bio-robotics.

Although, human gait analysis has been successfully performed using IMU, the threshold-based algorithms, and the fuzzy logic methods used, lack robustnes [2]. Thus, a more robust probabilistic technique is investigated in the present work. In this final year undergraduate project, we use Hidden Markov Model (HMM) to perform gait sequence modelling and estimation with IMU measurements of a moving dog.

1.1 Background to the study

Bio-inspired robotics identifies useful mechanisms from nature to inform real-world engineering systems. In this vein, the University of Cape Town (UCT) set out to investigate how a cheetah uses its tail for stability during high acceleration, quick turns and sudden braking for sophisticated robots design.

To achieve this purpose, inertial measurement unit (IMU) data for a dog running, walking

1.2. OBJECTIVES OF THIS STUDY

and trotting was acquired and labelled with the corresponding footfalls.

This dataset can, therefore, be used to perform gait sequence modelling and analysis using Hidden Markov Models (HMM), given its relatively small size.

1.2 Objectives of this study

The objective of the present work as stated in its outline is to "design, implement, and test (Hidden Markov Models) for estimating the gait sequence from IMU data, so that specific models can be formulated and their parameters estimated and interrogated." Thus, four main sub-objectives are formulated.

- Formulate HMM models to estimate the gait sequence of a dog from the labelled IMU measurements.
- Implement, test and evaluate the formulated models in order to estimate their parameters.
- Investigate how dimensionality reduction affects the performance of the designed models.
- Using the implemented models, perform useful gait pattern analysis such as motion recognition.

To achieve the above objectives the hypotheses have been formulated:

- HMMs can successfully model gait sequence dynamics using IMU data, in the absence of huge training data.
- Data aggregation and mirroring techniques can be used to overcome training data size limitations.
- Dimensionality reduction can increase the performance of an HMM, in the absence of a large training set.
- HMMs can be used to successfully perform gait activity recognition.

1.3 Scope and limitations

The project was proposed with a very broad scope. In order to ensure completion of the project in the allocated time, its scope was therefore reduced to the most essential aspects. Thus the project focuses on the design, the implementation, and the performance evaluation of HMM models to identify the dog's footfalls from IMU measurements.

Moreover, particular attention was paid to the design of dimensionality reduction techniques and how they affect the performance of the HMM models when the training dataset is not large.

Furthermore, using the implemented models, motion type recognition was performed using the implemented HMM models. However, the capturing and calibration of the IMU data is not part of the scope of this project.

The main constraint of this project was the relatively small size of the available IMU data.

1.4 Plan of development

Here you tell the reader how your report has been organised and what is included in each chapter.

Chapter 2

Literature review and basic theory

2.1 Hidden Markov Models

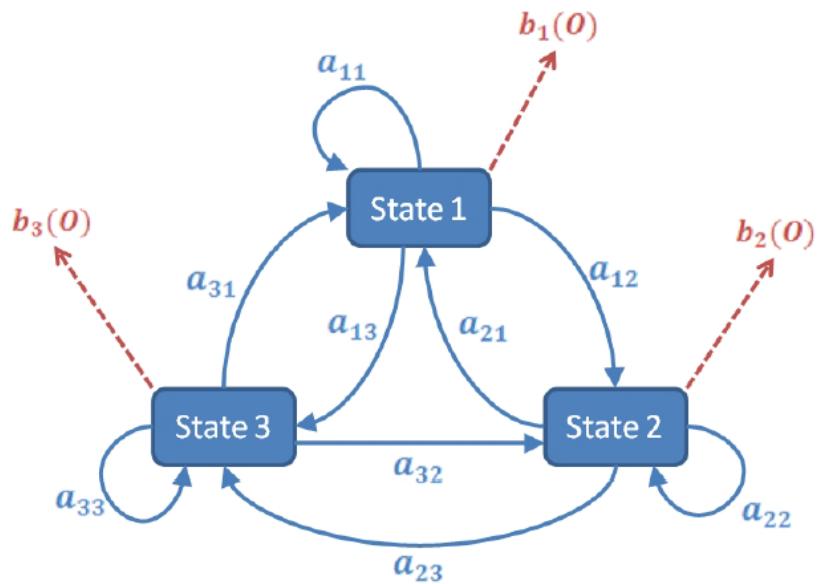


Figure 2.1: Example of a 3-states ergodic HMM
[8]

Hidden Markov Models (HMMs) are doubly embedded stochastic processes with a rich underlying statistical structure. Introduced at the end of the 1960s by Baum and colleagues, they have become one of the preferred techniques in speech recognition after the implementation of Baker and Jelinek in 1970s. HMMs have been successfully applied to various other engineering problems in pattern recognition for classification and fraud detection purposes, amongst others [1] [9] [3] [11].

The type of HMM depends on the possible connections between the states. Thus, an

HMM, in which a state can transition to any other state is said ergodic. Other types such as the Left-Right model or Bakis do not allow all possible transitions between the states. Figure 2.1 is an illustration of an ergodic HMM of 3 states (for simplicity). The annotations are explained in the next section.

2.1.1 HMM parameters specification

An HMM is fully specified by the following parameters

1. N, the number of distinct states of the model. Together they form the set of individual states $S = \{S_1, S_2, \dots, S_N\}$.
2. T, the number of observations. A sample observation sequence is denoted as $O = \{O_1, O_2, \dots, O_T\}$.
3. $Q = q_t$, the set of states with q_t denoting the current state at time instance, t such that $q_t \in S$ and $t = 1, 2, \dots, T$.
4. K, the number of distinct observation symbols per state.
5. $V = \{v_1, v_2, \dots, v_K\}$, the feature set of K dimensions.
6. $A = \{a_{ij}\}$, the state transition probabilities. a_{ij} denotes probability of transitioning from state S_i to state S_j .
7. $\Phi = \{\phi_j(k)\}$, the probability distribution of observation symbols in state j.
8. The initial state distribution, $\pi = \pi_i$

For continuous HMM (CHMM), i.e, HMM with continuous-valued observations, Φ consists in a probability distribution function. Many applications have successfully modelled such distributions with mixtures of Gaussian distributions [2] [4]. As such, ϕ is approximated by a weighted sum of M multivariate Gaussian distributions η . For a given, observation sequence, ϕ and η are therefore given by equations 2.1 and 2.2,

$$\phi(O_t) = \sum_{m=1}^M \beta_{jm} \eta(\mu_{jm}, \Sigma_{jm}, O_t), \quad (2.1)$$

$$\eta(\mu, \Sigma, O) = \frac{1}{\sqrt{(2\pi)^K |\Sigma|}} \exp\left(-\frac{1}{2}(O - \mu)' \Sigma^{-1} (O - \mu)\right) \quad (2.2)$$

$$1 \leq j \leq N; 1 \leq m \leq M; \beta_{jm} \geq 0; \sum_{m=1}^M \beta_{jm} = 1$$

where β_{jm} is the mixture composition coefficient; μ_{jm} , Σ_{jm} , respectively the mean vector and covariance matrix of state j ; M is the number of mixture components and K is the dimensionality of O . In practice, the log-likelihood of $\phi(O_t)$ is rather computed to avoid overflow when implemented on a machine.

As a summary, the compact specification of a continuous valued observation HMM is defined by 2.3 and that of a discrete HMM in 2.4.

$$CHMM = \lambda_C = (A, \beta_{jm}, \mu_{jm}, \Sigma_{jm}, \pi) \quad (2.3)$$

$$DHMM = \lambda_D = (A, b_j(k), \pi) \quad (2.4)$$

Basic assumptions of HMMs theory

HMM theory is built on three basic assumptions listed below.

1. *The Markov assumption:* HMM assumes that the probability of being in the current state at any instance of time t , is uniquely dependent on the previous state, at time, $t - 1$. More specifically, $a_{ij} = P[q_t = S_j | q_{t-1} = S_i]$. This assumption makes it unsuitable for long-range correlation capturing applications.
2. *The stationary assumption:* Furthermore, HMM state transition probabilities are assumed to be time-independent. Thus, the transition probabilities of two distinct time, t_1 and t_2 are identical, $P[q_{t_1} = S_j | q_{t_1-1} = S_i] = P[q_{t_2} = S_j | q_{t_2-1} = S_i]$. HMMs can therefore effectively model mechanisms with stationary observations.
3. *The output/observation independence assumption:* The current observation also known as emission symbol is statistically independent of the previous observations. It is "emitted" only by the current state, $P[O|q_1, q_2, \dots, q_T, \lambda] = \prod_{t=1}^T P[O_t|q_t, \lambda]$.

The three assumptions make an HMM a relatively simple graphical model to implement. This simplicity naturally comes with some limitations in modelling more complex problems which, however, may be modelled with higher-order HMMs. Furthermore, the three assumptions are very similar to those of a Markov chain. This is because the stochastic process of an HMM pertaining to the hidden states can be reduced to a Markov chain. In fact, an HMM is an extension of a Markov Chain. The essential difference between

the two is that, with the former, there is no a one-to-one mapping between the states and the observation symbols [7].

2.1.2 The three basics problems for HMM design

In [1], Lawrence argued that an HMM design needs to answer three basic problems. They are the *training problem*, the *evaluation problem*, and the *decoding problem*. Each problem and its solution are discussed in greater details next.

The evaluation problem

The evaluation problem is about answering this question: *Given the observation sequence $O = O_1 O_2 O_T$, and a model λ , how do we efficiently compute $P(O|\lambda)$, the probability of the observation sequence?* [1] [4] [2]. The naive answer to this question is simply computing the $P(O|\lambda)$ according to equation 2.5:

$$P(O|\lambda) = \sum_{q_1}^{q_T} \pi_{q_1} b_{q_1}(O_1) a_{q_1 q_2} b_{q_2}(O_2) \dots a_{q_{T-1} q_T} b_{q_T}(O_T) \quad (2.5)$$

This approach has two issues, it is not only, computationally expensive because of the exponential complexity with respect to T, but also, intractable for very long sequence. In practice, $P(O|\lambda)$ is computed by an algorithm called *forward-backward* procedure, which is a more efficient method.

The decoding problem

The decoding problem can be reduced to this interrogation: *Given the observation sequence $O = O_1 O_2 O_T$, and the model λ , how do we choose a corresponding state sequence $Q = q_1 q_2 \dots Q_T$ which is optimal in some meaningful sense i.e, best "explains" the observations?* Simply put, this problem is about deciphering the most likely hidden states that emitted the visible observation sequences. This is done dynamically using the Viterbi algorithm.

The training problem

Given the model, λ , the training problem raises the following question: *how do we adjust the model parameters λ to maximise the $P(O|\lambda)$, the probability of the probability of the observation sequence?* This problem is usually solved by iterative learning algorithms called expectation-maximisation. Examples of this algorithms are Baum-Welch method or any gradient based method. [1] [9] [7]

When using Baum-Welch algorithm, the parameters are initialised by guesses then re-estimated iteratively to find the parameters with maximum likelihood. This method is vulnerable to local maxima issues. To avoid such cases, it is advised to run it multiple times with different initial values in order to keep the estimation with the highest likelihood value.

2.2 Dimensionality reduction

Dimensionality or dimension reduction is used pattern recognition, machine learning and statistics to find the most compact representation of the dataset by removing redundant and irrelevant information [10]. It is achieved by extracting principal features, i.e, feature extracting or by selecting the most relevant subset of the initial feature vector, i.e, feature selection, using a supervised or an unsupervised approach.

2.2.1 Motivations for dimensionality reduction

When building a model, the need for dimensionality is supported by several reasons. Some of the important ones are presented in three points. Firstly, by reducing the feature space's dimension, we can build a model with higher quality[12]. In most classification problems, the feature domains contain variables with very little to no information for the purpose at hand. Thus, removing these features reduces the complexity of the problem which can in return, increases the model's accuracy.

Secondly, working with hundreds to thousands of features can be difficult to conceptualise and visualise. By using dimensionality reduction, we can better understand the model and present it to others by comprehensive visualisation[12].

The third reason is about efficiency in terms of computational time and storage. In

general, pattern recognition and machine learning algorithms computationally intensive. Besides, storage capacity is limited in some engineering applications such as embedded systems. So, solving the problem only with the relevant features can alleviate these two problems. Consequently, the computation speed of the model can be increased with dimensionality reduction [13].

Various dimensionality reduction have been developed in literature [14] [10] [13] [12], the next section will present a handful of the ones used in the present work.

2.2.2 Feature selection: filters and wrappers

Filters and wrappers are two major categories of feature selection methods [10] [12] [13] [14]. The structure of both methods are illustrated by figure 2.2 and figure 2.3, respectively. They both require a mechanism for generating a subset of the original

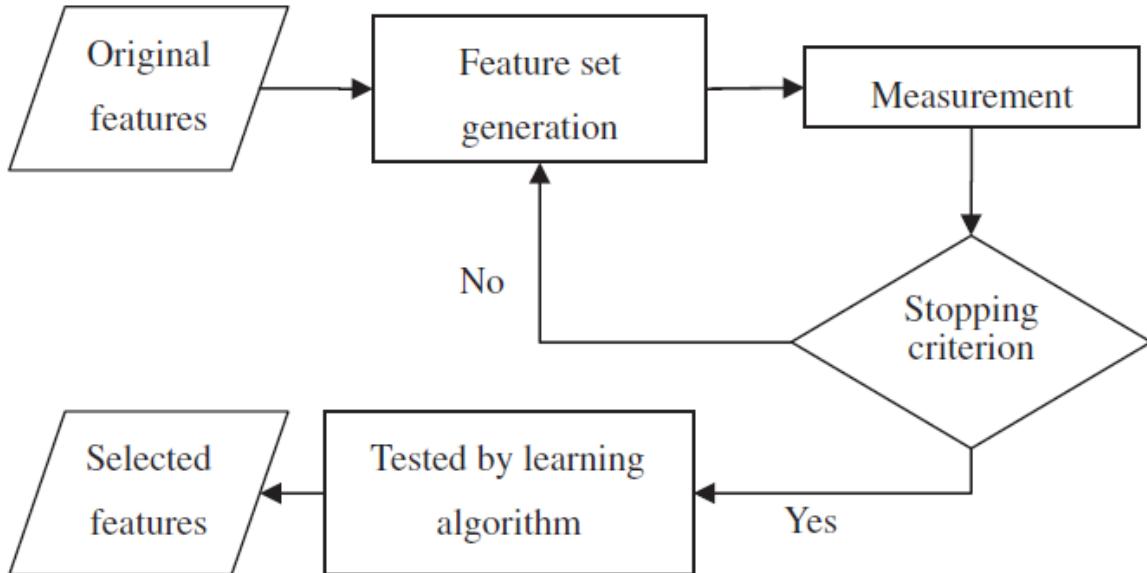


Figure 2.2: The procedure of filters in dimensionality reduction
[13]

feature set and a stopping criterion, which can be a simple distance measure or a more sophisticated separability measure between the features. However, wrappers identify the best feature subset using a learning algorithm usually by fully searching the feature space, whereas, filters use a simple measurement metric based on mutual information, correlation and other distance criteria [13] [14]. As a result of the two different approaches, filters are fast and do not guarantee optimal classification accuracy. On the other hand, wrappers

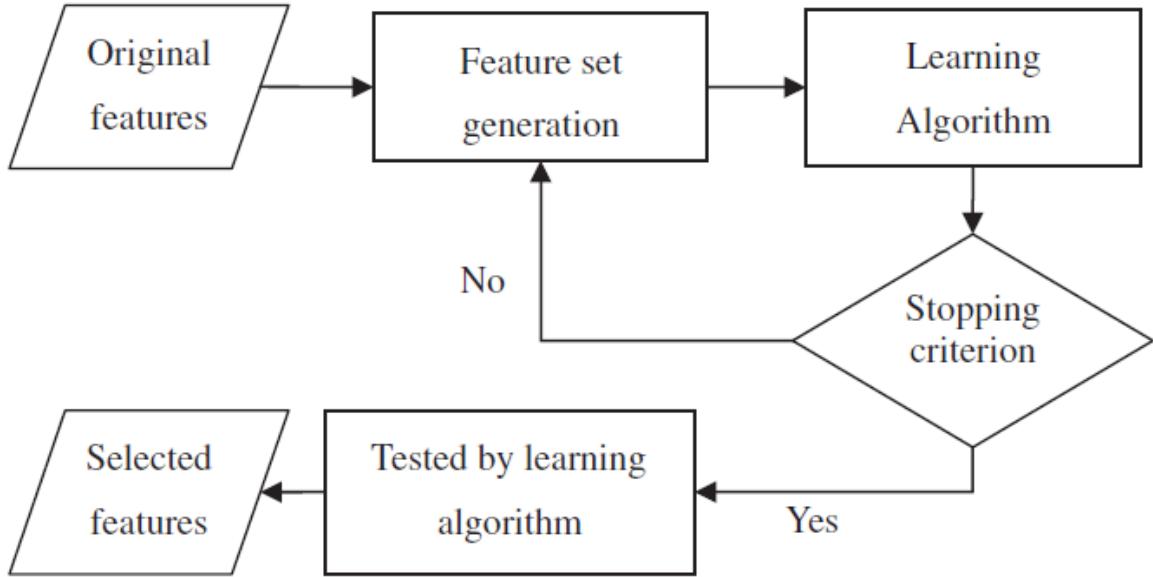


Figure 2.3: The procedure of wrappers in dimensionality reduction
[13]

give accurate prediction results but are very slow. Both approaches are often combined to build effective and efficient feature selection methods. One such approach is shown in figure 2.4. In this approach, the filters are used as a preliminary stage to discard irrelevant features. A feature is deemed irrelevant if it cannot discriminate between the different classes or if it contains redundant information [13]. The output of the filter stage is a smaller set of relevant features which is fed into the wrapper to find the optimal final subset of features. Thus, the filter stage effectively reduces the search time of the wrapper. In classification problems, using the separability index matrix to determine the classification content or degree of the feature generally results in very good results [14]. The next section gives a particular attention to this approach.

Feature ranking using separability index matrix

In this section, we look into a systematic approach to determine the 'classifiability of a feature' as present in [14]. In their paper, Jeong-Su, Sang Wan Lee and Zeungnam Bien, proposed a new criterion called separability index matrix to identify features that can discriminate between the different classes of a classification problem. The important concepts of this method are here defined together with their significance.

1. Separability Degree Matrix (SDM)

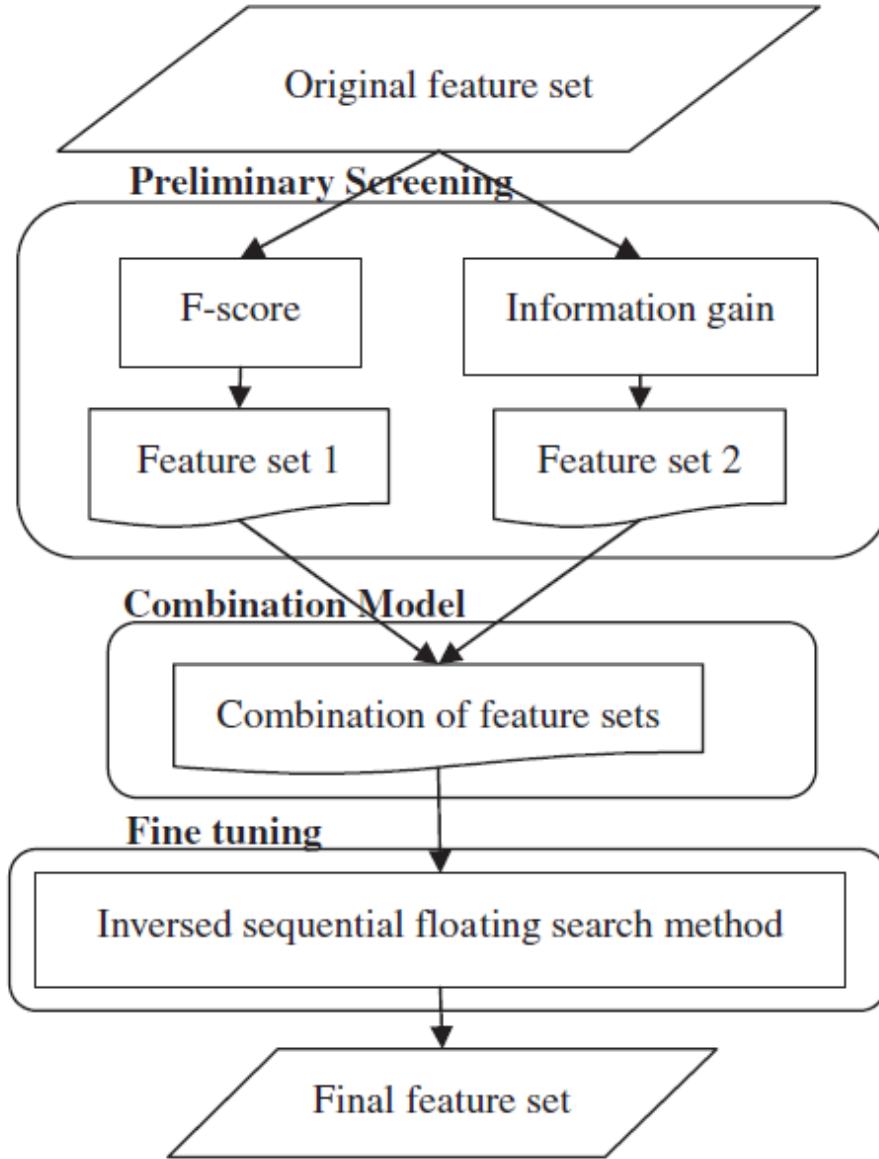


Figure 2.4: The procedure of hybrid filter-wrapper in dimensionality reduction [13]

SDM_k is a $C \times C$ symmetric matrix of the separability measures between the different classes of classification problem of C distinct classes given a particular feature x_k . It is defined in 1

$$SDM_k = [J(w_i, w_j; x_k)]$$

where $J(w_i, w_j; x_k)$ denotes the separability or the distance value between class w_i and class w_j when the criterion function $J(\cdot)$ such as mahalanobis [15] or euclidian, is applied to the feature x_k [14]. SDM_k is symmetric matrix with zero diagonal values, because for a given feature x_k , $J(w_i, w_j; x_k) = J(w_j, w_i; x_k)$ and $J(w_i, w_i; x_k) = 0$. This observation can be exploited to half the computation

required to calculate $SDM_v = SDM_1, SDM_2, \dots, SDM_N$, the set of all SDM for the feature set of size N.

2. Separability Index Matrix (SIM)

$SIM_k = c_{ij}$ is a CxC matrix of binary values 0, 1. if $c_{ij} = 0$, then the classes w_i and w_j are not separable by the feature x_k .

SIM_k is obtained by applying an threshold function to SDM_k . For instance,

$$\begin{aligned} SIM_k(i, j) &= 0 \quad \text{if } SDM_k(i, j) < SDM_{avg}(i, j) \\ SIM_k(i, j) &= 1 \quad \text{if } SDM_k(i, j) \geq SDM_{avg} \end{aligned}$$

where SDM_{avg} is the element-wise average of SDM_v . SIM_k is significant because it can be used to systematically determine the irrelevant features, i.e, features whose $SIM_k = 0$, and/or redundant ones. [14]

3. Classifiability: $G(x_k)$

Although, conventional distance criterion reveal the separability of the feature distribution, [14] we are often more interested in, how effectively can a particular feature distinguish one class from another. This information is denoted $G(x_k)$, the ' classifiability of a the feature x_k '. $G(x_k)$ is computed by 2.6.

$$G(x_k) = \sum_{i=1}^C \sum_{j=1}^C (SIM_k * WM_k) \quad (2.6)$$

$$\text{with } WM_k = SIM_k / \sum_{i=1}^N SIM_i \quad (2.7)$$

where * and / denote respectively element wise matrix multiplication and division.

This method can be used to effectively and efficiently rank the features in a classification for subset selection or a preliminary step in a hybrid filter-wrapper feature selection solution.

2.2.3 Feature extraction

In this work, two distinct feature extraction tools were explored namely, the linear discriminant analysis and the principal component analysis. Each technique is further explained in the next two subsections.

Linear discriminant analysis

Linear discriminant analysis maps a K-dimensional dataset, $O \in \mathbb{R}^K$ into a subspace that maximizes the class between-scatter with regards to the class within-scatter [16]. In other words, LDA finds the linear transformation - of the dataset into a lower dimensional space - that maximises the between class discrimination [17] It is a well-received dimensionality reduction tool in the computer vision and pattern recognition community.[17].

LDA is a supervised dimensionality reduction method, it does require the dataset to be labelled. On the other hand, PCA, the next feature extraction method is unsupervised method [16].

Principal component analysis

Principal component analysis (PCA) is an orthogonal linear transformation of a dataset of T K-dimensional observations, $O \in \mathbb{R}^{T \times K}$, into a so-called PC-space defined by the principal components (PC) [18]. Naturally, the dimension of the new space is smaller or equal to K, the original set's dimensionality. The components are ordered according to the variance of the features of the observation set. The principal components represent the eigenvectors of the covariance matrices of the features and the eigenvalues are measures of the features' variances. Thus, the first component is the feature with the highest variance [18]. In effect, PCA not only transforms the original observation set into a different space but, it also ranks the features according to their variability. It can therefore be used as an unsupervised dimensionality reduction technique by simply discarding the features with low variance after finding the principal components.

Chapter 3

Hidden Markov Model and dimensionality reduction design

3.1 Design overview

3.2 Description of available dataset

The available dataset was acquired from a moving dog using inertial measurement units. Two inertial measurements units (IMU) were strapped to the front and back of a dog. Each unit has an accelerometer, a gyroscope, and a magnetometer. The dataset contains calibrated measurements of a dog running, walking, and trotting then walking; together with the footfalls. The footfall is represented by a binary value that indicates the state of the dog's leg: if it is on or above the a particular instant in its gait sequence. More specifically, the value 0 means leg up and the value 1 means leg down. The four variables representing the footfalls effectively constitute the ground truth, informing us about the state in which the dog is, at a given time in its movement.

The dataset can be retrieved from nine different Matlab files. Each file contains twenty four Matlab variables. The variables of interest are listed in table 3.1.

	Observations			State	
Body part	Accelerometer	Gyroscope	Magnetometer	Left leg	Right leg
Front	accFrontX	FrontPitch	magFront_cal	LF	BF
	accFrontY	FrontRoll	magFront_cal2		
	accFrontZ	FrontYaw	magFront_cal3		
Back	accBackX	BackPitch	magBack_cal	LB	RB
	accBackY	BackRoll	magBack_cal2		
	accBackZ	BackYaw	magBack_cal3		

Table 3.1: Main variables of the dataset

3.3 Quadruped Gait sequence modelling with HMM

One of the objectives of this project is to effectively model the gait sequence dynamic of the dog from IMU measurements using HMM. Quadrupedes gait can be modelled as a succession of latent states observed through the 'visible' IMU measurements. The states representing the footfalls and the observations, the outputs of the accelerometer, gyroscope and the magnetometer. Similar to human gait mechanism, it is sound to assume that the current state of a quadruped is conditionally dependent on its previous state. This inference combined with the statistical robustness of HMM makes it the best model candidate when the available dataset is not too large.

3.3.1 HMM model elements: states and observations properties

The problem at hand requires 16 distinct states that make up the state vector S , shown in equation 3.1

$$S = S_i = \{(LF, RF, LB, RB)\} = \{0000, 0001, 0010, \dots, 1111\}. \quad (3.1)$$

$$|S| = N = 2^4 = 16$$

$$i = 1, 2, \dots, 16$$

3.3. QUADRUPED GAIT SEQUENCE MODELLING WITH HMM

The 16 distinct states are derived from the combination of the four binary footfalls. In practice, the dataset may not reveal all the 16 states.

The stream of IMU measurement forms the observation sequence. An observation instance is a row vector of K dimensions. The initial K value before any dimensionality reduction is 18, from the 18 IMU measurements. Thus, an observation sequence O is a TxK matrix of continuous-valued voltages as presented in 3.2. T is the total number of the successive measurements.

$$O_t = \{O_t^k\} = O_t^1, O_t^2, \dots, O_t^8, t \quad (3.2)$$

$$k = 1, 2, \dots, 18. \quad (3.3)$$

$$t = 1, 2, \dots, T. \quad (3.4)$$

3.3.2 Splitting the 16-states HMM into two 4-states HHMs

In order to simplify the problem, it was decided to split the four legs in two sub-parts: two front legs and two back legs. This decision exploits the fact that "*fore and hind quarters of dogs behaved like two independent bipeds*" as demonstrated in [19]. As a result, the initial 16-states HMM becomes two distinct 4-states HMMs. These two models may be combined to reconstruct the holistic 16-states model. With the two separate HHMs, we are faced a simpler task of discriminating between 4 distinct states instead of 16. A further benefit of this decision is that this work may well be applied to human gait estimation using IMU measurements. From here onward, we will focus on the 4-states HMM model.

3.3.3 State transitions

It was assumed that in its movement, a dog may transition from one state S_i to any other state S_j where $i, j = 1, 2, 3, 4..$. For instance, if a dog has its left leg above ground and its right leg on the ground, at the time instance t, it may move to any of the 3 other possible positions or remain in the same state, in the next time instance, $t + 1$. Thus, we are dealing with an ergodic type HMM, where all the states transitions are allowed (Please refer back to figure 2.1 for illustration on an ergodic HMM).

3.4 Pre-processing and increasing the dataset

Very little data pre-processing was required given that a relatively clean data was already available in Matlab files. So, this stage of the design consisted of three simple steps.

1. **Extracting and labeling the feature vector:** In this step, the $T \times K$ observation matrix was formed using the 18 different variables listed in 3.1, for a given gait sequence. Using, the footfalls as ground-truth, i.e, (LF, RF) and (LB, RB), each observation instance was correctly labelled. The labeling is useful to stratify the different into the four states whenever required in the training process.
2. **Aggregating the gait sequences:** The individual gait sequence samples - for a walking, running and trotting dog - are very small in size, it was therefore decided to aggregate them in a single longer sequence. As such, the observation sequence becomes that of a dog walking, trot, then running or a similar combination of the three actions. The impact of such a decision on the model's parameter is that, instead of fitting a single action, it fits all the three types of actions. In theory, this loss of specificity can negatively affects the model's precision. Nevertheless, the overall design can still be applied to a specific dog's action dataset if there is enough training data.
3. **Increasing the data by mirroring:** : The aggregated observation sequence amounts to 2695 different samples. In order, to further increase the dataset, the mirrored sequence was appended to the initial sequence. This decision works on the assumption that, given a gait sequence, the reverse sequence is also a valid sequence from the subject matter. Practically speaking, each mirrored observation sample is a duplicate, no new sample is added. Moreover, the state distribution remains unchanged. However, as demonstrated in the results section 4.2, this increase in data size caused an increase in the model's performance.

As a summary, the output of the pre-processing stage is an 18×5390 matrix where each row represents an observation sample and each column a particular feature. This dataset was fed directly into the HMM model or in the dimensionality reduction module to remove irrelevant or redundant features. The design of the dimensionality reduction stage is now explained.

3.5 Dimensionality reduction

Five different dimensionality reduction methods were considered in this work: three feature selection methods and two are of feature extraction type. The feature selection methods are: a feature ranking based on separability of Index, a distance-based forward feature selection, and a full-search feature selection with pre-defined output feature vector dimension. Principle component analysis (PCA) and linear discriminant analysis (LDA) were considered regarding the feature extraction. The two different types were considered to determine the most suitable one in gait sequence analysis with IMU measurements. For the feature selection methods, a preliminary experiment was performed to determine the optimal number of features. Before presenting this experiment, it is necessary to explain the feature ranking method based on separability of index.

3.5.1 Feature ranking with separability index (SI)

The ability of a feature X_k to classify the different states was defined as the 'between-class classifiability', $G(X_k)$. Refer back 3 for greater details. This quantity can be used a metric for ordering the different features. The greater the $G(X_k)$, the better feature X_k is at classifying the different states. Applying this concept to our problem at hand, the different $G(X_k) = G(X_1), G(X_2), \dots, G(X_{18})$ were computed using the Mahalanobis distance [15] [30]. . Then, the 18 features (the IMU measurements) were ranked accordingly, in a descending order.

The ranking is advantageous because it can be used as feature selection method 3.5.3. This was the case when determining the minimum feature size discussed in the next subsection.

3.5.2 Determining optimal number of features using SI and MCE

The objective of this experiment was to determine the minimum number of features that best represent the dataset in a classification problem.

The experiment was performed with a K-Nearest Neighbour classifier together with feature ranking using separability index, by following the steps below:

1. **Step 1: *Features ranking*:** Firstly, the features were ranked in a descending

order according to the 'classifiability of each feature' as defined in 3.

2. **Step 2: *Classifier design***: Then, a 1-Nearest Neighbour (1-NN) classifier was designed and trained using half of the available dataset with the first K different features ($K \leq 18$, the initial feature vector size). The first K features are assumed to be the best possible candidate features as explained in 3.5.1.
3. **Step 3: *Classifier's performance***: Furthermore, the classifier was tested with the remaining dataset and the misclassification error (MCE) was recorded.
4. **Step 4: *Iteration over the feature subset size***: Finally, Step 2 and 3 were repeated while varying the subset size from 1 to 18.

Figure 3.1 shows the percentage misclassification error (MCE) as a function of the feature subset size.

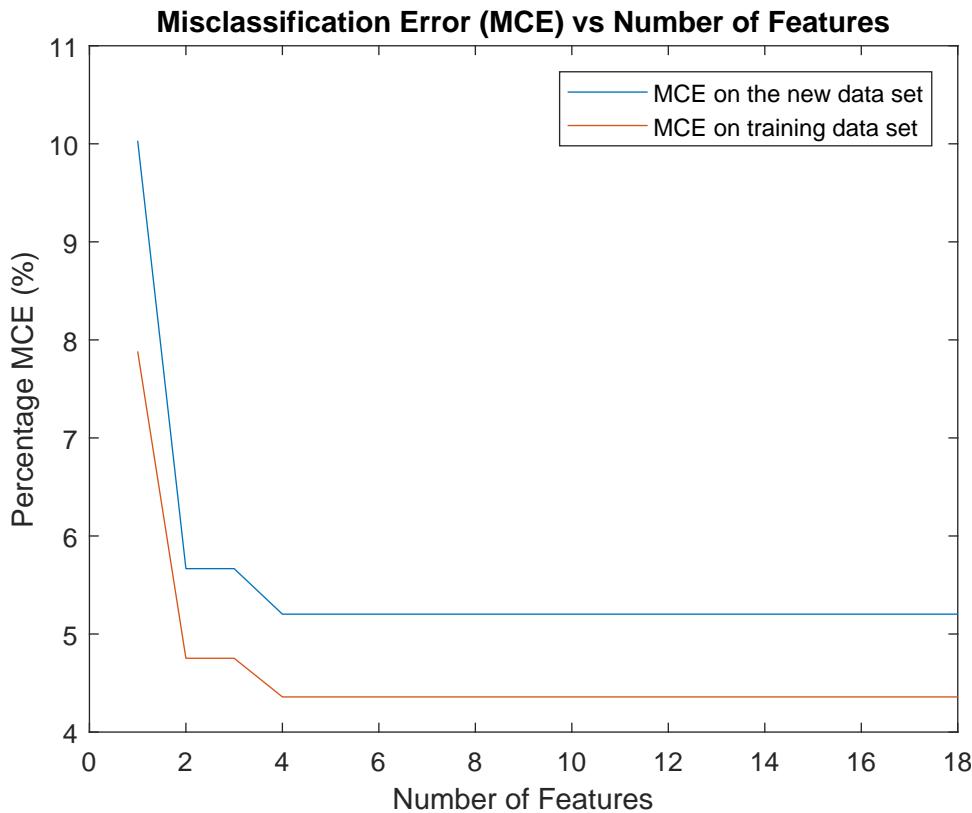


Figure 3.1: Misclassification error vs feature number using separability index and KNN

The MCE decreases from 1 feature to 3 and stays constant from 4 features to 18. The minimum number of features that best distinguishes the four states is therefore four. The following feature subset selection methods will consequently focus on determining the best combination of the four features.

3.5.3 Feature ranking

This ranking filter method is based on ranking the features using separability of index as presented in 3.5.1. It is simple and efficient dimensionality reduction method. It can be implemented in two very simple steps. Firstly, the features are ranked in a descending order using a ranking criterion such as the 'feature classifiability'. Secondly, the first K features are selected where K is the desired number of output features. It is built on the assumption that the combination of features with higher 'classifiability content' is the best candidate. The drawback of this assumption is the possibility of redundant features. For this reason, this method can be used as a preliminary step to a more sophisticated feature selection approach. In this present work, it was used as a separate filter method and its performance was compared to other methods in 4.2. The next section presents a more sophisticated but computationally expensive method: the forward selection.

3.5.4 Forward feature selection

Forward feature selection is a sequential feature subset selection, "in which features are sequentially added to an empty candidate set until the addition of further features does not decrease" [20]. a stopping condition is met. The stopping condition can be a pre-defined number of features as a design constraint or until the addition of further features leaves the evaluation criterion unchanged.

In this project, the evaluation criterion used is 1-Nearest Neighbour (1-NN) [21]. Furthermore, since the purpose of the feature selection is to reduce the 18, initial features down to a much smaller number, the output feature number was limited to four, the optimal number of features, according to the outcome of the previous experiment performed in 3.5.2. Because this method is based on a sequential search, it is computationally very expensive. In order to make the algorithm faster, a preliminary step was introduced to select the relevant features. Thus, the features, with zero classification ability, were discarded. These features have their 'classifiability content', $G(X_k) = 0$. This step can significantly decrease the computational time of the overall algorithm.

Although, the forward feature selection is more sophisticated than the feature ranking, it does not guarantee the best possible subset of features. This is because it does not do a full search of all possible combinations. The next filter method takes advantage of the relatively small initial feature vector's size, i.e, 18 to perform a full-search feature selection of pre-defined number of output features.

3.5.5 Full-search feature selection

Given the pre-defined output feature vector's dimensionality (4) and its relatively small initial value (18), a full-search feature selection could be explored. Thus, by eliminating the irrelevant features as explained in the above method 3.5.4, the number of combinations to be searched can be drastically reduced. In fact, nine IMU measurements had no 'classification content', reducing the number of combination to just $\binom{9}{4} = 126$. Because a full-search will ensure the best possible combination of features, this feature selection method was implemented. Here, the evaluation criterion was the prediction accuracy of the CHMM model. Effectively, making this method a pseudo-hybrid filter-wrapper with pre-defined output feature vector dimension.

3.5.6 PCA and LDA

The design of the PCA and LDA were very simple. The former was designed as follows. Firstly, the covariance matrices of the features were computed. Then, the eigenvectors of these matrices were computed along with their corresponding eigenvalues. Given a feature, its corresponding eigenvalue represents the variance in the feature. Thus, the eigenvectors were ranked in a descending order according to their eigenvalues. The first K eigenvectors or components represented the dataset in the new PC space, where K dimension of the new observation sequence.

Optimal PCA component number

The optimal K value was determine with the HMM model itself. In fact, different HMMs were trained with varying PC size. The number of components which results in the highest precision is taken as the optimal K value. One such experiment was performed with 539 dataset while varying the number of PC components from 1 to 10. The results is shown in 3.2

Thus, the optimal component number is 4. It is worth noting that after 4, the model's precision significantly decreases. The overall performance is very low, this fact will be further explained in the results chapter. The relative performance is sufficient to determine the optimal PC number.

LDA is a supervised feature extraction method. Thus, the observation sequence was labelled with the corresponding state sequence. This served as ground-truth to maximises

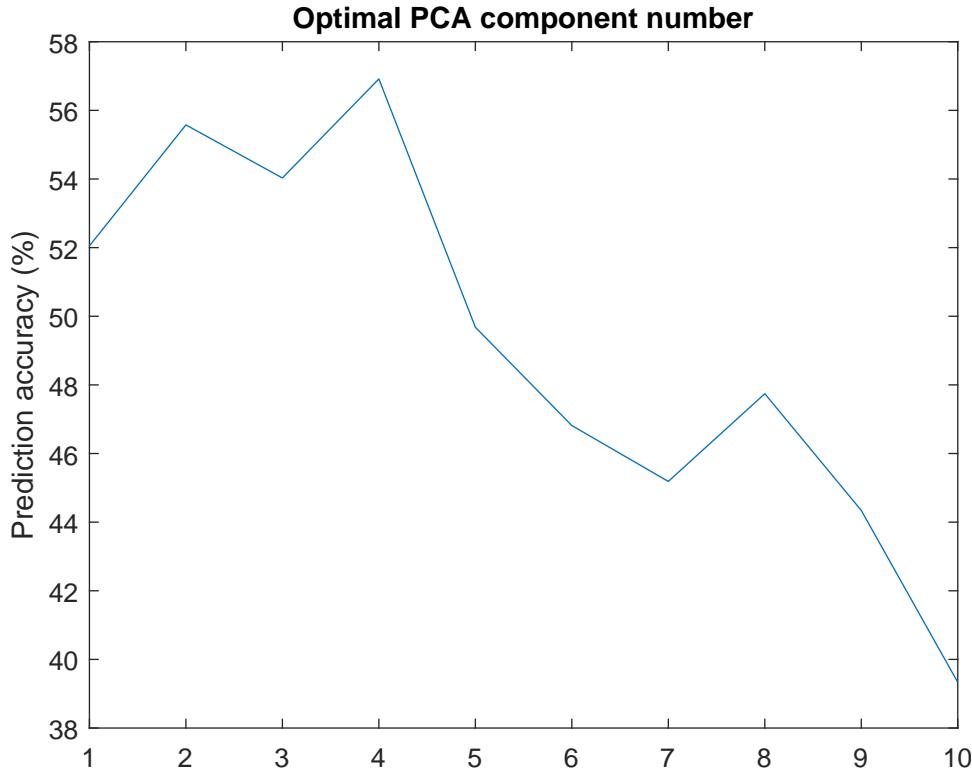


Figure 3.2: Optimal PCA component number

the between-scatter over the within-scatter [37] [38].

Because, the number of classes in question is 4, the dimensionality of the reduced output data was set at 3, the highest possible dimension with LDA of 4 classes [32]. Number 3 was chosen empirically instead of 1 or 2.

3.6 Solving the three basics HMM problems

As a reminder, a continuous HMM model is completely specified by its initial state distribution: π transition matrix: A ; the mean covariance matrices: μ, Σ which can be combined into Φ . When the observations are modelled with Gaussian mixture distributions, an addition parameter is required to represent the initial mixture distribution: β . Having clearly defined the model's parameters, we will now focus on solving each of the three basic problems of an HMM design.

3.6.1 Solution to the training problem

As a reminder, the training consists in estimating the parameters the parameters of the HMM model. The important steps of this algorithm are explained in the points below.

Estimation of the transition matrix: A

For each of the front and back 4-state HMM, the state transition matrix A , is a 4-by-4 matrix. Two different approaches were considered in the estimation of A . The two methods make use of the expectation maximisation algorithm but differ in the input arguments considered.

1. ***Approach 1: EM algorithm using Gaussian mixture model:*** This method is the standard approach found in literature using expectation algorithm such as Welch-Baum algorithm any gradient-based method [9]. So, prior to maximising the transition probabilities -in the M-step-, each observation is first labelled as having been emitted by a particular state S_i - in the E-step - using the Gaussian mixture model used to model the observations. This process may introduce some degree of error, depending on how well the Gaussian mixture models the observation sequence. The resulting transition matrix may not be very representative of the real underlying state transition mechanism, especially when the dataset is not large enough. For this reason, the second approach, which eliminates the labeling step by exploiting available ground-truth, was developed in this work.
2. ***Approach 1: EM algorithm using ground-truth:*** This technique takes advantage of the ground truth provided by the labelled dataset to reduce the HMM model to its underlying Markov process.

This was achieved by setting the observation sequence to the state sequence, i.e, $O = O_1O_2\dots O_T = S_1S_2\dots S_T = S$. In effect, each continuous observation is replaced by its label. Effectively, the continuous-valued observations are transformed into discrete values, for the purpose of accurately determining the transition matrix. This transformation is sound because the transition matrix is only concerned with the probability of transitioning from one state to the other. It does not care about the actual value of the observation.

Up to this point, we have a "discrete observation sequence O with known states, S ". The next and final step is simply about computing the maximum likelihood of the transition using EM algorithm, as one would do for a discrete HMM with

known state sequence [22].

To make the transition matrix more representation of the state transition mechanism, prior knowledge, also known as pseudo-count was added to the estimated transition probability as follows: $A = \{a_{ij}\} = \{a_{ij}^{estimated} + prior_j\}$ [7] [22]. In this work, the pseudocount, $PseudoA_j$ for a given state S_j was set to the number of occurrences of S_j , in the training data plus a constant value C: $PseudoA = |S_j| + C$. The additional constant C is to take care of theoretically possible states and transitions that are not revealed by the dataset. It can be chosen empirically until the transition matrix does not contain zero. Although, this method is very simple and more effective, it has two limitations. It not only assumes that the number of states is known but also requires the training data to be labelled. Since these two constraints are met in this design, it was chosen as a better option.

Estimation of the mean matrix, covariance matrices, and mixture distribution: μ , Σ and β

”Gaussian mixture models (GMMs) are powerful in modelling any desired continuous distribution and are used for example in speech processing successfully” [4], and IMU based HMM for gait human classification [2] [4]. So, in this project, the IMU observations were modelled with GMM density functions. These functions represented by a KxMxN mean matrix: μ , a KxKxMxN co-variance matrices: Σ , and a MxN, mixture prior distribution β , where: K is *the number of features*; M, *mixture number*; and N, *number of distinct states*.

More specifically, the observations were grouped into N classes based on the four distinct states. Subsequently, the three GMM parameters were estimated for each state, by using the expectation maximisation algorithm. The important steps of this algorithm are explained in the points below.

- **Parameters initialisation by k-means++:** The initial GMM parameters were estimated using a variant of k-means called k-means++ which uses a heuristic to find clusters in a dataset [23], [24].
- **EM-algorithm step:** The iterative EM algorithm is thereafter applied to optimise the initial parameters. The algorithm can be optimised by turning some parameters for both speed and accuracy. The main parameters considered in this projects are explained in the next points. In general, model accuracy was given preference over speed.

- ***Replication***: Given that the EM-algorithm may suffer from local maxima problem, it was repeated multiple times, i.e, 10 times, with different initial values at each iteration. The model with the highest log-likelihood value is therefore chosen [24].
- ***Maximum number of iterations***: Because the EM-algorithm does not always converge, a maximum number of iterations needs to be provided. In the present work, it was experimentally set at 1000 iterations by favouring accuracy over speed.
- ***Termination tolerance***: The EM-algorithm terminates upon the log-likelihood equating to a certain value with a tolerance called termination tolerance. This value was set at 10^{-10}
- ***Regularization value***: This is a positive value added to the diagonal of each covariance matrix to ensure that the estimates are positive-definite [23] [24].
- ***Unbiased variance estimation***: The variance estimator used was normalised with N-1 to make the co-variance matrices unbiased.

In literature, the number of mixture of IMU data is set empirically to 2 or 3 mixture components [4] [2], in this project, it was estimated using akaike information criterion (AIC), which is a measure of information loss [25] [26]. So, distinct Gaussian mixture models were built using EM algorithm with varying mixture component number: from 1 to 18, a randomly stop at the maximum feature number. Since AIC is a measure of information loss, the model with the minimum AIC best represents the dataset. The number of mixture M, is therefore set at the mixture number of this model. This algorithm is outlined in 3.1 and a typical result is shown in 3.3.

```

1 function [ numComponents , AIC ] = optimal_mixture_component_by_AIC ( data )
2 X = data . observ ;
3 max_mixt_num = 18 ;
4 AIC = zeros ( 1 , max_mixt_num ) ;
5 BIC = zeros ( 1 , max_mixt_num ) ;
6 GMModels = cell ( 1 , max_mixt_num ) ;
7 options = statset ( 'MaxIter' , 2000 ) ;
8 for k = 1 : max_mixt_num
9     GMModels { k } = estimate_model_parameter ( X , k ) ;
10    AIC ( k ) = GMModels { k } . AIC ;
11    BIC ( k ) = GMModels { k } . BIC ;
12 end
13
14 [ minAIC , numComponents ] = min ( AIC ) ;
15 numComponents
16 minAIC

```

17 end

Listing 3.1: Pseudo-code for finding the optimal number of mixture using AIC

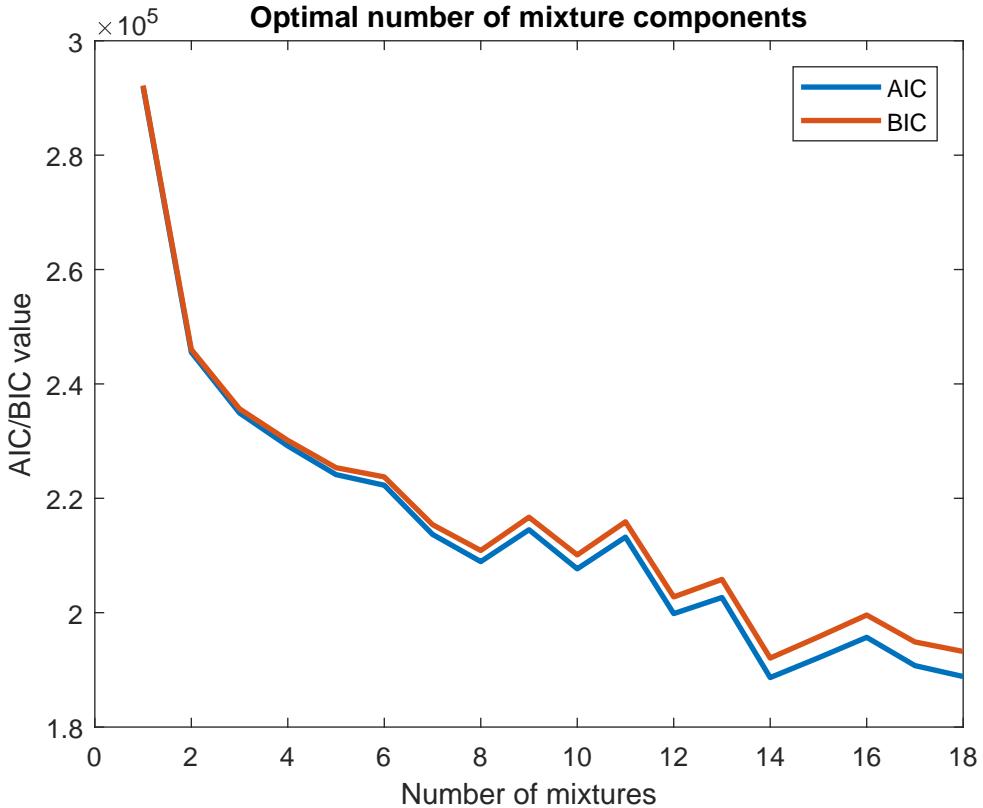


Figure 3.3: Finding best number of mixture component using AIC

As shown in 3.3 the same experiment can be performed using Bayesian Information Criterion (BIC), an equivalent but stricter information loss measure [24] [27]. The more the number of components, the less information is lost. In this example, 14 is the optimal number of components for the range 0 to 18. As expected, the trend in BIC confirms this fact.

Initial state distribution: π estimation

In this project, the initial state distribution was estimated using the probability of occurrence of each state in the training data sample. Thus, for a given set of length T, $\pi = \{\pi_i\} = \frac{|S_i|+C}{T}$, where C is a constant to cater for any possible state that is not present in the training dataset. Similar to the transition matrix estimation, this could have been determined by an iterative process if the state sequence was unknown. However, given that there is an available state sequence, it was deemed unnecessary.

3.6.2 Solution to the evaluation problem

The evaluation problem is about efficiently computing the probability of a given an observation sequence $O = O_1O_2\dots O_T$. In other words, the likelihood of the sequence having been generated by the HMM model. This problem was solved in two parts described below. Note that the log-likelihood or log-posterior probabilities were computed to avoid overflow or underflow problems.

1. ***Part 1: Computing the log-probabilities:*** In this step, the estimated Gaussian model is used to calculate the log- probabilities of the dog being in a given state according to 2.1 and 2.2.
2. ***Part 2: Determining the log-likelihood with log-forward-backward:*** The log-probabilities together with the state transition matrix, A and the initial state distribution π_i are fed into the log-forward-backward algorithm to compute the log-likelihood[28].

3.6.3 Solution to the decoding problem

The decoding problem consists in predicting the hidden states, i.e, the positions of the dog's limbs that "generated" the sequence of IMU measurements. After successfully decoding the initial state, $Q_0 = q_0$ the problem is reduced to predicting the next state $Q_{t+1} = q_{t+1}$ when in state $Q_t = q_t$. This simplification follows from the *Markov assumption* 2.1.1. Fortunately, the state transition matrix A gives insight into determining the most likely next state Q_{t+1} from the current Q_t with the initial state distribution π_i as bias. This fact can be represented by the following equation $Q_{t+1} = AQ_t + \pi$. Practically speaking, the state sequence was decoded with the Log-Viterbi algorithm using A, π_i and the log-probabilities obtained from the Gaussian model 2.1.1. This algorithm dynamically computes the posterior log-probabilities of the observations for all four states. The state with the highest probability is most likely to have generated the observation in sequence.

3.7 Evaluation of the footfall identification accuracy

This brief section deals with the evaluation of the HMM model's decoding accuracy. The decoding is about successfully identifying the correct footfalls, given an observation

instance. So, the identification accuracy was evaluated as a percentage of the number of correctly classified footfalls over the total number of footfalls to be identified:

$$\text{accuracy} = 100 \times \frac{\text{number of correctly identified footfalls}}{\text{length of observation sequence}} \quad (3.5)$$

Differently put, the classification error is the percentage of misclassified footfalls. However, state (footfall) transition misclassifications were tolerated. These errors are simply due to the HMM model predicting a transition occurrence one step early or one step late. Listing 3.2 shows the pseudo-code of the accuracy evaluation.

```

1 function accuracy = evaluate(ground_truth, estimation)
2   accurate_estimation = sum(ground_truth == estimation) +
3     count_wrong_transitions(ground_truth, estimation);
4   accuracy = 100 * accurate_estimation / size(ground_truth, 1);
5 end
6
7 function count = count_wrong_transitions(ground_truth, estimation)
8   count = 0;
9   for i = 1:size(ground_truth, 1) - 2
10    if (ground_truth(i) ~= ground_truth(i + 1) || estimation(i) ~= estimation(i + 1)) ... % transition in ground_truth or estimation
11      && ((ground_truth(i + 1) ~= estimation(i + 1)) ... % estimation lagged
12        or led the transition
13        && (ground_truth(i + 2) == estimation(i + 2)) ... % by just one step
14        , i.e., it got it correct in the next step
15        && (ground_truth(i + 1) == ground_truth(i + 2) || estimation(i + 1)
16        == estimation(i + 2)) %make sure it was just a transition in ground
17        truth or in estimation
18      count = count + 1;
19    end
20  end
21 end
22
```

Listing 3.2: Pseudo-code for computing the model's accuracy

3.8 Practical implementation

The prototype design was implemented with three main MATLAB toolboxes: *MATLAB Statistics and Machine Learning Toolbox* [36], the pattern recognition tool named *prtoools* [35] and an open-source HMM implementation by Qiugiang Kong called *matlab-hmm* [28].

Table 3.2 summarises each of the main design components and the corresponding function used to implement it. Please refer to their references for more details.

Design component	Toolbox	Function
Transition matrix	Statistics and Machine Learning	hmmpredict[22]
Gaussian mixture	Statistics and Machine Learning	fitgmdist[24]
State posterior probability	matlab-hmm	Gmm_logp_xn_given_zn[28]
Log-likelihood	matlab-hmm	LogForwardBackward [28]
State sequence decoding	matlab-hmm	LogViterbiDecode citematlab-hmm
Forward feature selection	prtools	featself[29]
Separability degree matrix	prtools	distmahacitemaha
PCA	prtools	pcam[31]
LDA	prtools	fisherm [32]
K-NN classifier	prtools	knn[21]
Data partitioning	MATLAB and prtools	cvpartition [33] and gendat[34]

Table 3.2: Main design components together with toolboxes and functions used for their implementation

The next two sub-sections show the codes used to train and test the HMM model. The rest of the implementation code can be found in A.1

3.8.1 Model's parameters estimation implementation

The listing 3.3 is the implementation code used to build the model. It is the implementation of the solution to the training problem described here 3.6.1.

```

1 function HHM_model = GmmHMMBuild(observ_seq , state_seq)
2     if nargin < 1
3         [observ_seq , state_seq] = get_aggregated_data('front');
4     end
5     state_num = 4; %Number of states
6     feat_num = size(observ_seq , 2); %feature size
7     mix_num = optimal_mixture_component(observ_seq , state_seq);
8     C = 1; % To account for state not in training dataset
9     observ_groups = group_observation_per_state(observ_seq , state_seq);
10    % state initial probability distribution
11    state_distr = [size(O1, 1) size(O2, 1) size(O3, 1) size(O4, 1)] + C;
12    model.pi = state_distr/(C + size(observ_seq , 1));
13    % state transition probabilities estimation
14    PSEUDOTR = ones(state_num , state_num)*C + state_distr ;
15    model.A = hmmpredict(state_seq , state_seq , 'PSEUDOTRANSITIONS' ,PSEUDOTR) ;
16    % GMM optimisation parameters
17    replicate_num = 10;
18    options = statset('Display','off' , 'MaxIter' ,1000 , 'TolFun' ,1e-10);

```

```

19 regularization = 1e-10;
20 % GMM parameters estimation
21 covariances = zeros(feat_num, feat_num, mix_num, state_num);
22 means = zeros(feat_num, mix_num, state_num);
23 mix_prob = zeros(mix_num, state_num);
24 for state = 1:state_num
25     O = observ_groups(state);
26     if isempty(O) % if state not available generate a random distribution
27         O = (-10 + (10+10))*rand(mix_num + C, feat_num);
28     end
29     gmm = fitgmdist(O, mix_num, 'Replicates', replicate_num, 'Options',
30 options, 'Regularize', regularization);
31     means(:, :, state) = gmm.mu.';
32     covariances(:, :, :, state) = gmm.Sigma;
33     mix_prob(:, state) = gmm.ComponentProportion;
34 end
35 phi.mu = means;
36 phi.Sigma = covariance;
37 phi.B = mix_prob;
38 model.phi = phi;
39 end

```

Listing 3.3: Implementation code for estimating HMM model's parameters

3.8.2 Decoding and log-likelihood computation implementation

Listing 3.4 is the implementation code of both solutions to the decoding 3.6.3 and the evaluation problem 3.6.2.

```

1 function [accuracy, loglik, post_prob, path] = GmmHMMPredict(model, testdata)
2 if nargin < 1
3     model = GmmHMMBuild();
4     [testdata.observ_seq, testdata.state_seq, testdata.feat_names] =
5         get_aggregated_data();
6 end
7 % Solving the evaluation problem in the next two lines
8 % 1-Computing the log-probabilities of the dog being in a given state
9 logp_xn_given_zn = Gmm_logp_xn_given_zn(testdata.observ_seq, model.phi);
10 % 2-Computing the log-likelihood with log-forward-backward
11 [~,~, loglik] = LogForwardBackward(logp_xn_given_zn, model.pi, model.A);
12 % Solving the decoding problem with log-viterbi
13 [path, post_prob] = LogViterbiDecode(logp_xn_given_zn, model.pi, model.A);
14 % calculate the decoding accuracy
15 accuracy = evaluate(testdata.state, path);

```

15 end

Listing 3.4: Implementation code for decoding hidden states and computing sequence log-likelihood

3.8. PRACTICAL IMPLEMENTATION

Design

Chapter 4

Results

4.1 Experiment 1: Observation sequence dimensionality's effect on performance

4.1.1 Aim of the experiment

The aim of this experiment is to investigate how the number of features impacts the accuracy of a Hidden Markov Model with continuous emission symbols (CHMM), in the absence of enough training data. Thus, the hypothesis under investigation is:

In the absence of enough training data, a CHMM with observations of high dimensionality performs poorly.

4.1.2 Experiment apparatus

To perform this expereiment, the following materials are required:

- λ , a continuous Hidden Markov Model specified by $\lambda = (A, \beta_{jm}, \mu_{jm}, \Sigma_{jm}, \pi)$.
- At least two sample data sets for training the model and testing it.
- A criterion to rank and select subsets of features.
- A measure to evaluate the performance of the CHMM model.

4.1. EXPERIMENT 1: OBSERVATION SEQUENCE DIMENSIONALITY'S EFFECT ON PERFORMANCE

- Finally, a way to visualise the results of the experiments

4.1.3 Experiment procedure

The experiment was performed with just the back limps by following the steps listed below:

1. Step 1 - Preliminary data pre-processing: This step consisted in the data pre-processing as described in 3.4.
2. Step 2 - Partitioning data into training and test sets: Here, the dataset was randomly sampled into training and test sets. The training set was made relatively small, it was a sequence of 539 observations.
3. Step 3 - Feature ranking: The features were ranked using separability of index 3.5.3.
4. Step 4 - Constructing and training the models: 18 different CHMM models, $\lambda = \lambda_i = \lambda_1, \lambda_2, \dots, \lambda_{18}$ were built and trained using the first i features of the same training dataset.
5. Step 5 - Testing and evaluation the models: After training the models, they were tested with the same test dataset with the correct feature subsets and their performance were evaluated in terms of decoding accuracy and log-likelihood value.
6. The different accuracies and the sequence log-likelihood values were finally plotted as a function of the observation dimensionality. Moreover, the observations were grouped based on the corresponding hidden state sequence and scattered in a 2-dimensional principal component space. This is to compare the decoded states against the ground-truth.

4.1.4 Experiment results

The results of the experiments are presented in figure 4.1, 4.2, 4.3, 4.4, 4.5, 4.6.

Figure 4.1 and 4.2 respectively show how the hidden state decoding accuracy and the test sequence log-likelihood estimated by the CHMM model varies as the the number of features considered increases.

Figure 4.3 and 4.4 are visualisations of the 5-dimensional observation sequence grouped respectively, according to the actual state sequence and the decoded state sequence.

4.1. EXPERIMENT 1: OBSERVATION SEQUENCE DIMENSIONALITY'S EFFECT ON PERFORMANCE

Figure 4.5 and 4.6, are for an 18-dimensional observation. 5 and 18 dimensions were presented because they resulted in the two extreme prediction accuracies. The results for other dimensions may be found in the appendix from figure A.1 to figure A.32,

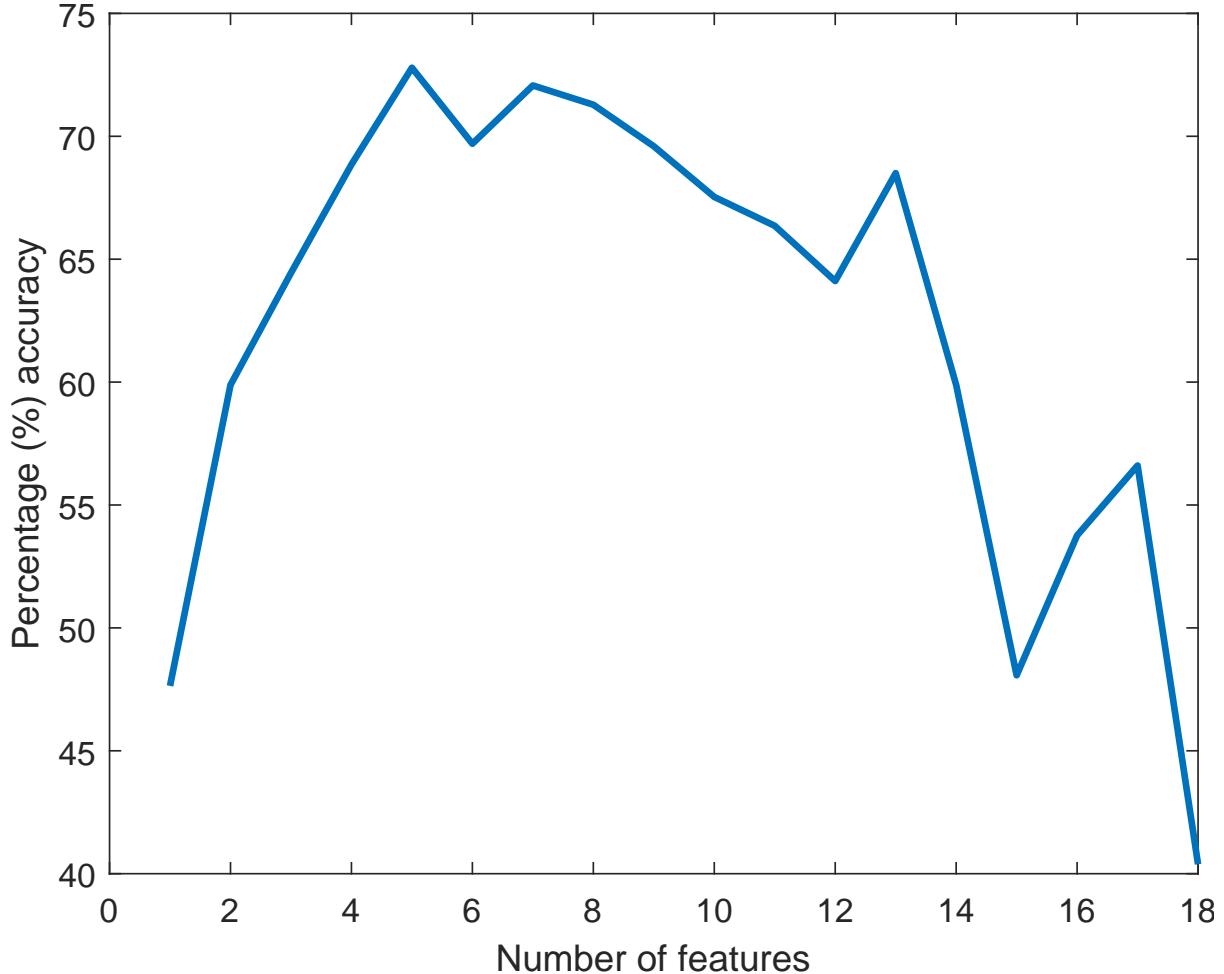


Figure 4.1: The effect of CHMM's observation dimensionality on the state sequence decoding accuracy

4.1.5 Analysis and discussion of results

Starting with 1 feature, the CHMM model's state decoding accuracy in 4.1 increases. It reaches its maximum performance at about 73% accuracy with 5 features. After 5 features, the accuracy generally declines to about just 40% accuracy with all 18 features. The decline is however not consistent, for example, there is a local peak at 13 and 17 features. A performance increase of up to 78% was obtained by reducing the feature set's dimension by from 18 to 5. In addition, the sequence log-likelihood value is close to zero from 1 to 8 features. It then slowly declining - in overall - from 9 to 15 features and suddenly gets to its lowest value at 18 dimensions. The log-likelihood being a measure

4.1. EXPERIMENT 1: OBSERVATION SEQUENCE DIMENSIONALITY'S EFFECT ON PERFORMANCE

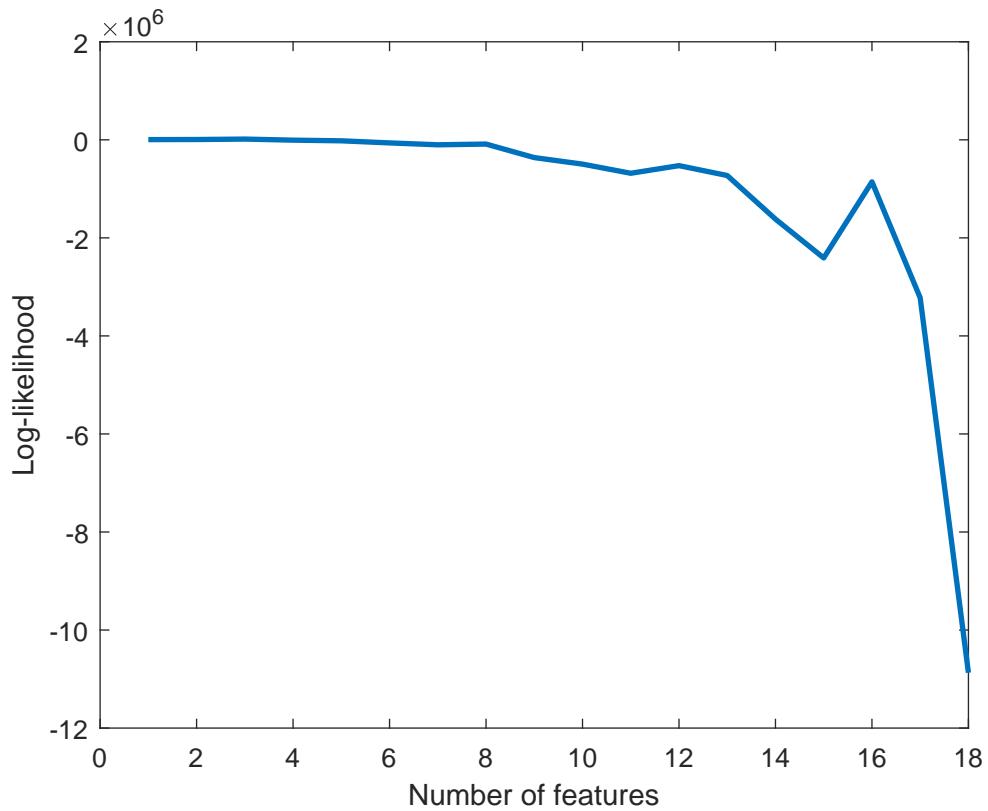


Figure 4.2: The effect of CHMM's observation dimensionality the log-likelihood

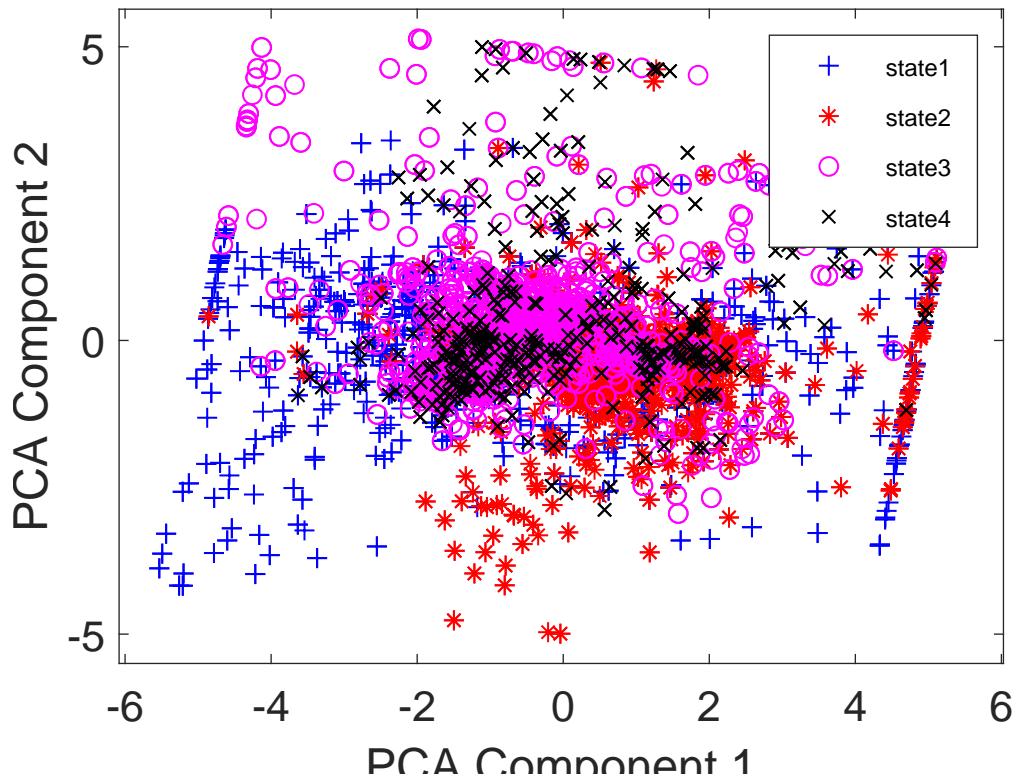


Figure 4.3: Scatter plot of 5-dimensional observations grouped according the ground-truth state sequence

4.1. EXPERIMENT 1: OBSERVATION SEQUENCE DIMENSIONALITY'S EFFECT ON PERFORMANCE

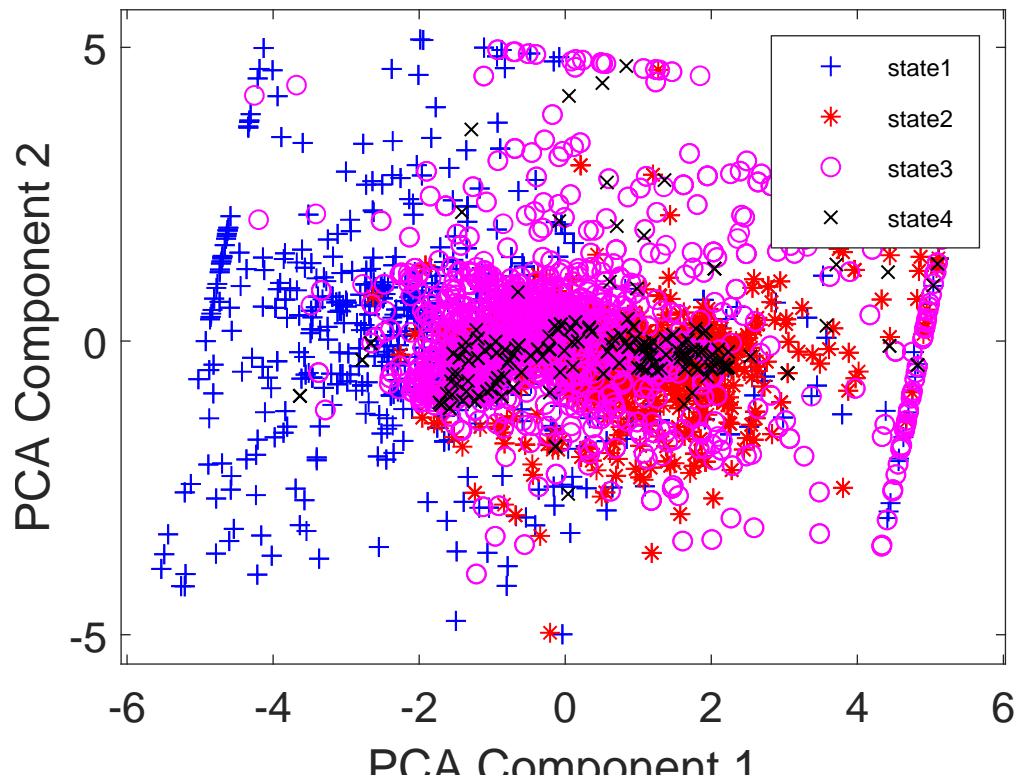


Figure 4.4: Scatter plot of 5-dimensional observations grouped according to the estimated state sequence

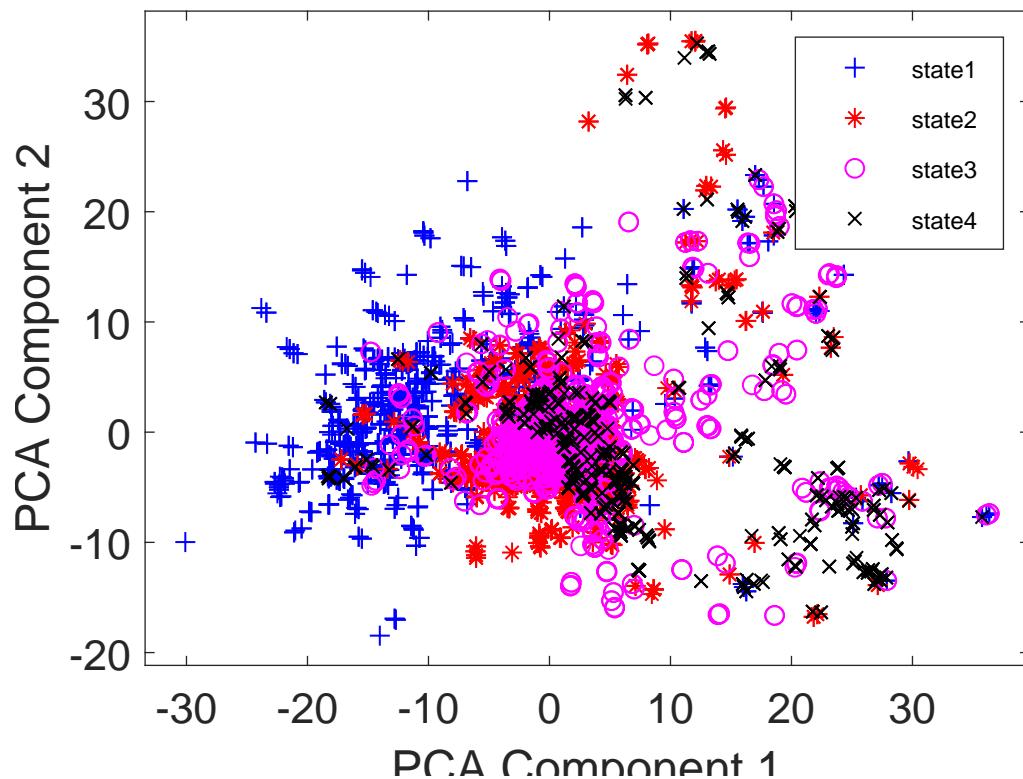


Figure 4.5: Scatter plot of 18-dimensional observations grouped according to the ground-truth state sequence

4.1. EXPERIMENT 1: OBSERVATION SEQUENCE DIMENSIONALITY'S EFFECT ON PERFORMANCE

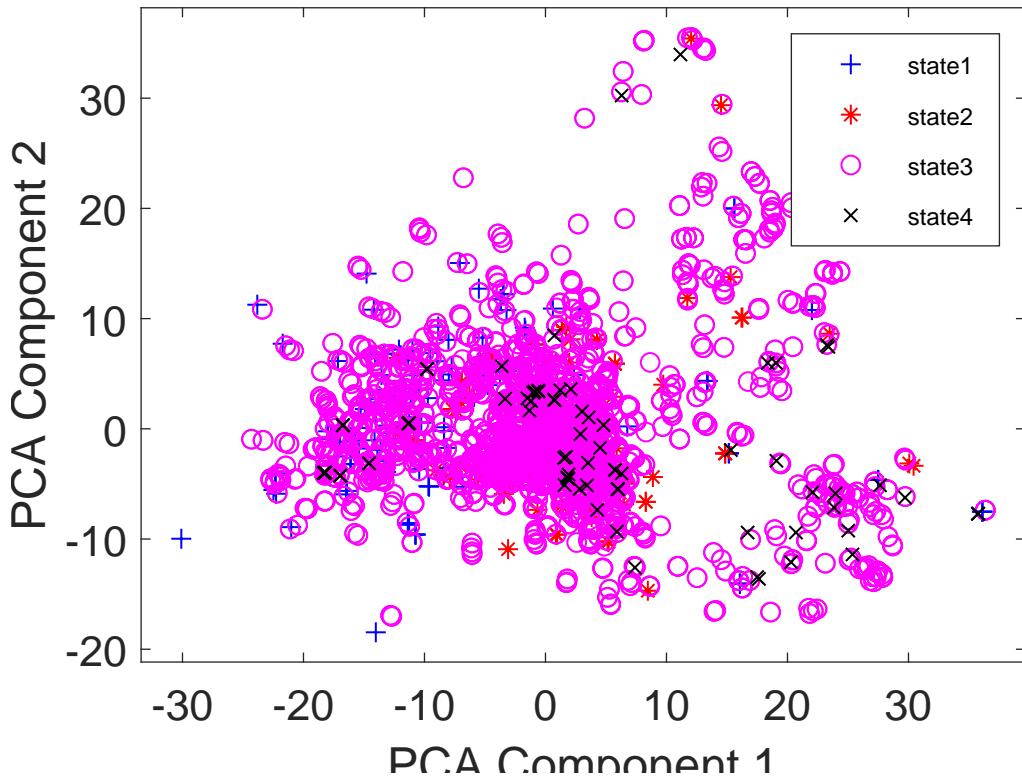


Figure 4.6: Scatter plot of 18-dimensional observations grouped according to the estimated state sequence

of the probability that the test sequence was generated by the model, it is clear that the CHMM model better recognises test sequence with smaller feature sizes.

Thus, the smaller the feature size, the better the CHMM performance with a relatively small training dataset. The model is both better at decoding the hidden states as well as recognising IMU measurements generated by the dog's gait mechanism. This fact is further shown by 4.3, 4.4, 4.5 and 4.6. With 5 features, the model correctly attributes most IMU measurements to the correct footfalls except the state 4 observations which were mostly classified as state 3 emissions. In opposite, with 18 features, the HMM model predicted wrongly that the dog remained in state 3 during all the 539 sequences. In fact, with 18 dimensions, the model is not twice as good as predicting with random guesses, i.e. $0.4 < 2 \times \frac{1}{4}$.

4.1.6 Conclusions and recommendations of the experiment

Based on the above results and discussions, it can be concluded that in the absence of enough training data, a CHMM performs very poorly. This validates the hypothesis formulated. It is therefore necessary to further explore the effect of dimensionality reduction with more robust techniques. This investigation is performed in the next

4.2. EXPERIMENT 2: THE IMPACT OF DIMENSIONALITY REDUCTION ON PERFORMANCE

experiment with various dimensionality reduction technique and varying data sizes.

4.2 Experiment 2: The impact of dimensionality reduction on performance

4.2.1 Aim of the experiment

The aim of this experiment is to investigate the effect of dimensionality reduction on the performance of a continuous Hidden Markov Model (CHMM). The hypothesis under investigation is therefore the following: *Dimension reduction can cause an increase in a CHMM's performance when there is not enough training data.* Thus, the performances of the CHMM with and without dimensionality reduction are compared to test the hypothesis.

4.2.2 Experiment apparatus

The assets needed to perform the experiment are listed below.

- λ , a continuous Hidden Markov Model specified by $\lambda = (A, \beta_{jm}, \mu_{jm}, \Sigma_{jm}, \pi)$.
- Four distinct dimensionality reduction methods. Two feature extraction methods and two feature selection methods were considered. The feature extraction methods were Principle Component Analysis (PCA), Linear Discriminant Analysis (LDA). The two feature selection methods were feature ranking with separability index denoted SI Ranking, and a combination of forward feature selection and separability index denoted SI-forward.
- A performance measure to compare the different models with and without dimensionality reduction. Again, the metrics used were the hidden state decoding accuracy and the ability to recognise a sequence generated by the model namely the log-likelihood.

4.2. EXPERIMENT 2: THE IMPACT OF DIMENSIONALITY REDUCTION ON PERFORMANCE

4.2.3 Experiment procedure

The experiment was performed as follows. Firstly, the dataset was partitioned into two different set for training and testing using random sampling.

Using the same training dataset, five specific models were built and trained. They are denoted as: $\lambda_{NoReduction}$, λ_{PCA} , λ_{LDA} , λ_{SI} , $\lambda_{SI-forward}$, respectively, for the model with all 18 features, the models built with PCA, LDA, SI Ranking and SI-forward. Regarding the last four models, the dimension of the training set was reduced with the corresponding technique before feeding it into the actual CHMM model. (Please refer to 3.5 for more details on the dimension reduction methods design)

Next, the different models were tested with the same test dataset - naturally applying the dimensionality reduction technique where required- and the prediction accuracy and the log-likelihood values were recorded.

This above procedure was repeated while varying the proportion of training data used from 10% to 90% of the total dataset.

The prediction accuracies and log-likelihoods were finally plotted as a function of the training data size for each model. These findings are presented in the figures 4.7 and 4.8.

In order to verify, the variance in the five distinct models, a separate an additional experiment was performed for a bias-variance error analysis. Thus, five models were trained and tested 100 times with different datasets. For each model, the average accuracy over the 100 tests was taken as its bias error and variance in the 100, as the variance error.

4.2.4 Experiment results

Firstly, figure 4.7 illustrates how the performances of the five CHMMs compare against each other as the training data size increases. Secondly, the log-likelihoods presented in 4.8 shows how effectively each model can recognise an observation sequence generated by the underlying mechanism. These results are discussed in the next sub-section. Finally, 4.9 is the bias-variance errors of the different models.

4.2.5 Analysis and discussion of results

In the sub-sections below, the effect of each dimensionality reduction technique is discussed.

4.2. EXPERIMENT 2: THE IMPACT OF DIMENSIONALITY REDUCTION ON PERFORMANCE

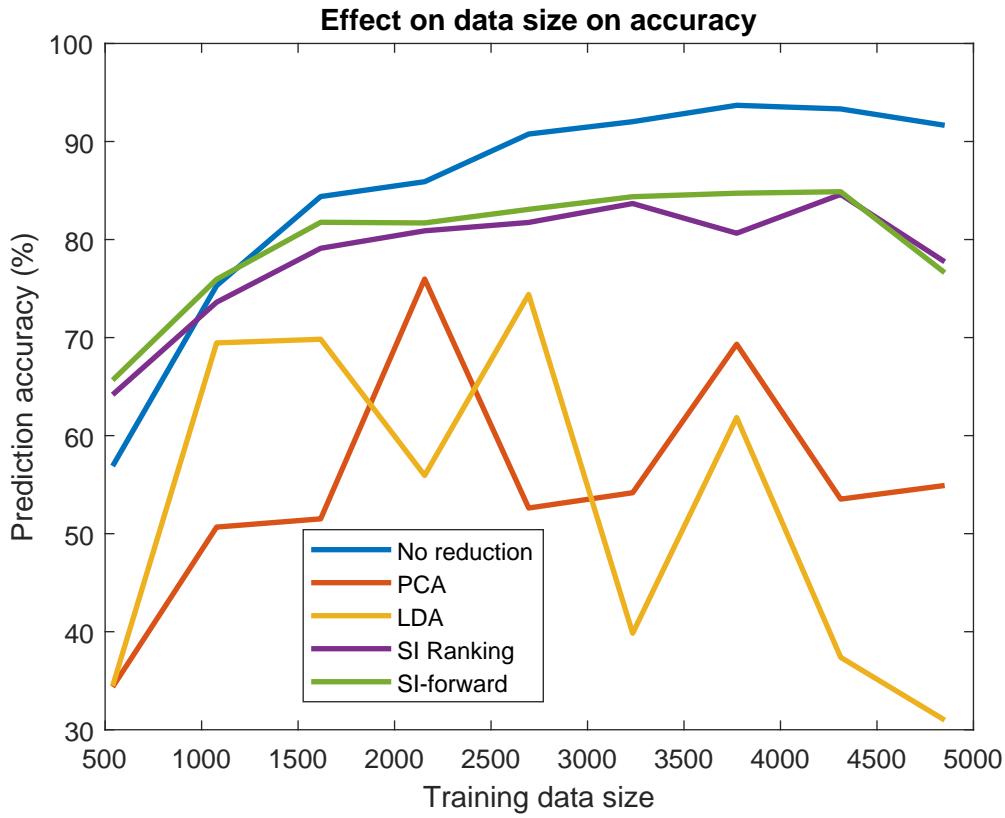


Figure 4.7: The effect of training datasize on the prediction accuracy

No reduction vs feature Ranking: SI Ranking

The two graphs under considering here are the purple and the blue graphs in 4.1 and 4.2. For a relatively small training data size, i.e, under 1000 samples, the model with feature subset selection using feature ranking outperforms the model with dimensionality reduction. This fact is true for both the prediction accuracy and log-likelihood values. At 539 sample data, there is at least 8% improvement after reducing the feature size in classification. Based on the linear nature of the graphs before 1000, it could be extrapolated that this improvement will be more significant for observation sequences with less than 539 samples. After, 1000 samples, the accuracies of both models increases but, the model built with the 18-dimensional observation sequence increasing performs better.

No reduction vs forward feature selection: SI-forward

Here, we are discussing the green and the blue graphs in 4.1 and 4.2. From these two graphs, we can note that the forward feature selection has the similar effect to the feature ranking. This is because they are both feature subset selection methods. It is therfore a

4.2. EXPERIMENT 2: THE IMPACT OF DIMENSIONALITY REDUCTION ON PERFORMANCE

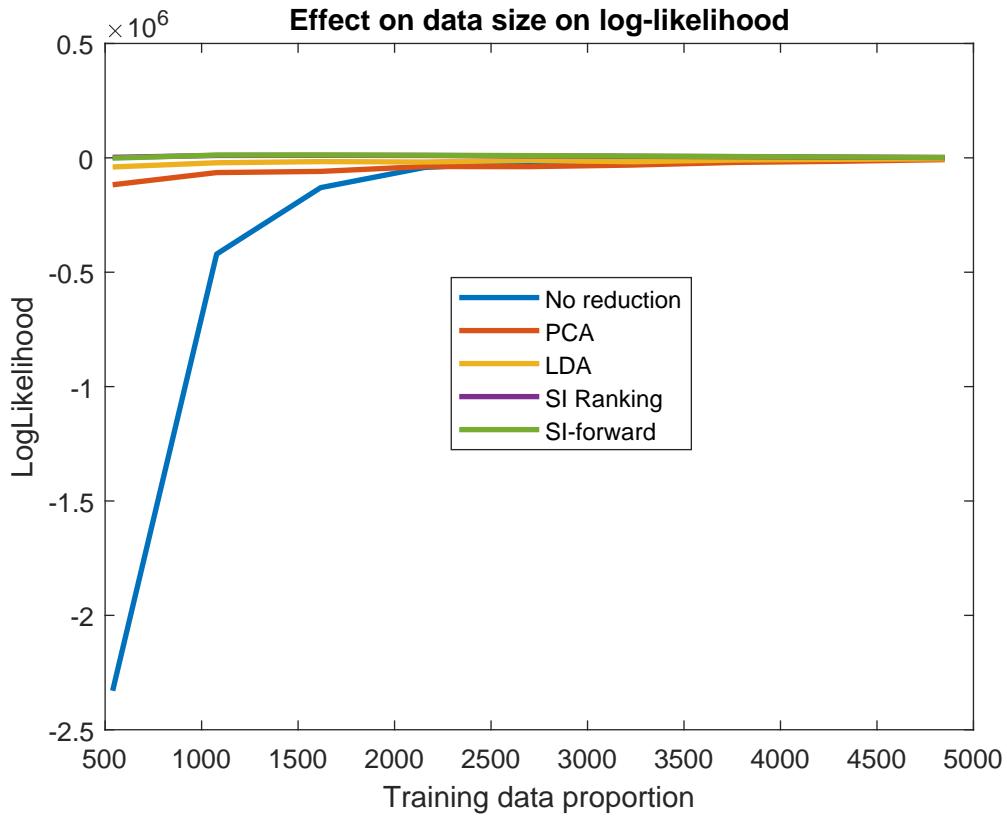


Figure 4.8: The effect of training datasize on the log likelihood

confirmation of the above discussion.

From these similar results, we can conclude that the very fast feature ranking 3.5.3 method can be used in lieu of the slower forward feature selection method 3.5.4. This result confirms the efficacy of the 'feature classifiability' 3 as feature ranking criterion for classification applications.

No reduction vs PCA and LDA

The two feature extraction methods decreased the performance of the models for all data sizes. In addition, their performances do not consistently increase with the increase in data size. The nature of IMU measurements is not suited to these two feature extraction methods. The experiment could have been better with other successful feature extraction methods found in literature such as Fast Fourier Transform (FFT) [3], time-frequency domain analysis [4].

4.2. EXPERIMENT 2: THE IMPACT OF DIMENSIONALITY REDUCTION ON PERFORMANCE

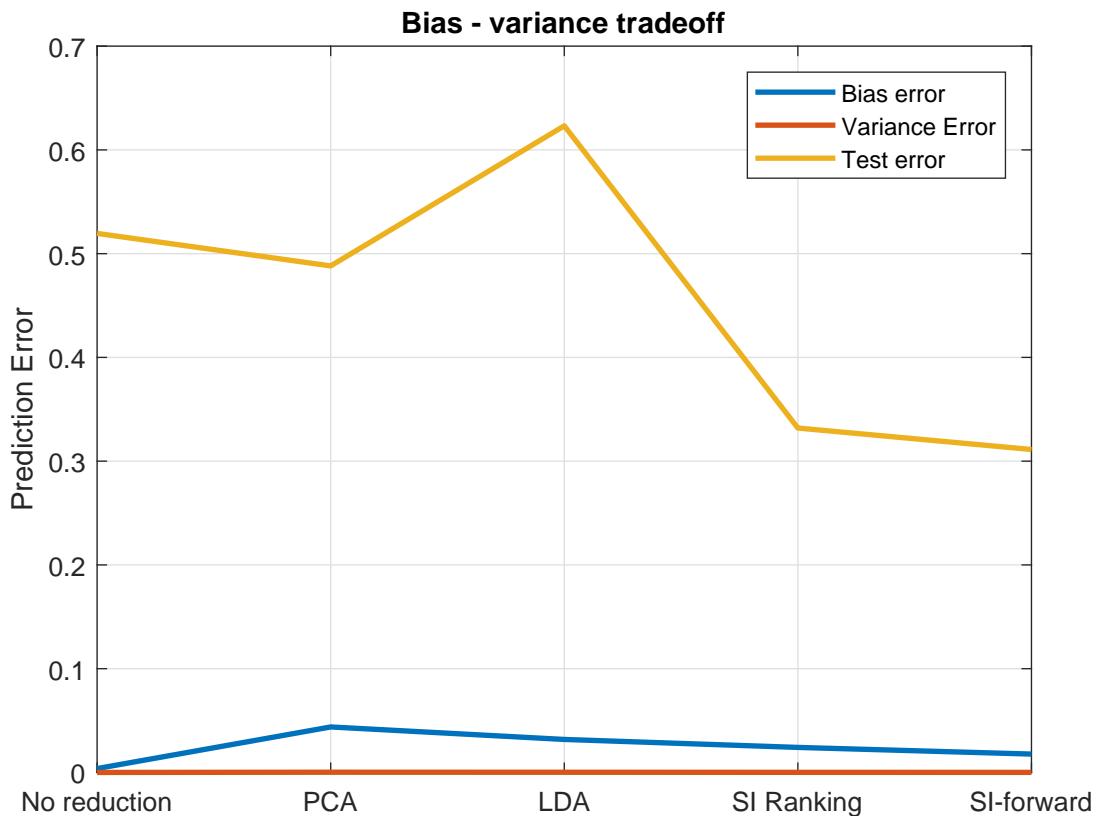


Figure 4.9: Bias-Variance tradeoff analysis

Bias - variance error analysis

The presented results in 4.9 is for a dataset with 539 samples. The bias error for all five models is relatively low. This because the test data was similar to the training data. Interestingly, the variance in the training errors over the 100 differet runs are all very close to zero. This is not very surprising because the model's parameters, particularly the mean and the co-variance matrices were designed to have low variance. By looking at the test errors and the bias errors, the best model appears to be the model built with the forward feature selection method, $\lambda_{SI-forward}$ when the dataset size is small. The next comparable model would be the model built with the feature ranking metho, i.e, λ_{SI} . So, with a very strict constraint on the model's speed, λ_{SI} should be favoured because, it is significantly faster than $\lambda_{SI-forward}$.

4.3. EXPERIMENT 3: THE NECESSITY OF COMBINING THE FRONT AND BACK IMU SENSOR MEASUREMENTS

4.2.6 Conclusions and recommendations of the experiment

Considering the above findings, we can effectively confirm that dimensionality reduction can effectively improve the performance of an IMU based CHMM when the dataset is not large enough. However, this does not apply to every dimensionality reduction method. A large range of technique should be explored to select the most appropriate such as the forward feature selection and the feature ranking methods.

4.3 Experiment 3: The necessity of combining the front and back IMU sensor measurements

4.3.1 Aim of the experiment

The objective of this experiment is to determine whether or not it is necessary to consider both the front and back IMU measurements when predicting only the front or back footfalls. Thus, the following hypothesis was constructed:

When predicting only the front or back footfalls of the dog, building the model with the aggregated front and back IMU measurements can improve its performance.

4.3.2 Experiment apparatus

To perform this experiment, the following materials are required:

- λ_f and λ_b , two continuous Hidden Markov Models for the front and back limbs respectively specified by $\lambda_f = (A_f, \beta_{jm_f}, \mu_{jm_f}, \Sigma_{jm_f}, \pi_f)$ and $\lambda_b = (A_b, \beta_{jm_b}, \mu_{jm_b}, \Sigma_{jm_b}, \pi_b)$. These models will be further split into two.
- Four data samples: two distinct training sets and two distinct test sets for the two models, where each contains both the front and back IMU measurements.
- A measure to evaluate the performance of the CHMM model.
- Finally, a way to visualise the results of the experiments

4.3. EXPERIMENT 3: THE NECESSITY OF COMBINING THE FRONT AND BACK IMU SENSOR MEASUREMENTS

4.3.3 Experiment procedure

To verify our hypothesis, similar experiments to predict the back and front footfalls of the dogs were performed. Further details are given next.

First of all, two front and back models: $\lambda_{f_{combined}}$ and $\lambda_{b_{combined}}$ were constructed by using the combined IMU data from two sensor sets.

Furthermore, two other front and back models: $\lambda_{f_{just-front}}$ and $\lambda_{b_{just-back}}$ were built using only the front IMU and back IMU measurements, respectively.

Thus, we have two couples of models to be compared: ($\lambda_{f_{combined}} vs. \lambda_{f_{just-front}}$) and ($\lambda_{b_{combined}} vs. \lambda_{b_{just-back}}$)

Secondly, the models to be compared were trained using the same training dataset, where the back or front IMU measurements were removed in the respective case of $\lambda_{f_{just-front}}$ and $\lambda_{b_{just-back}}$.

Finally, the models were tested with their corresponding training datasets. This was done purposefully since the relative comparison of the models with 'combined' and 'separate' IMU data is the subject matter, not the individual performances.

The experiment was repeatedly performed while varying the size of the training dataset. The state decoding accuracies and the corresponding log-likelihood were recorded. They are presented in the following sub-section.

4.3.4 Experiment results

The results of the current experiment are presented in 4.10, 4.11, 4.12 and 4.13. In the same order, they represent the front footfalls decoding accuracy, the front sequence log-likelihood, back footfalls decoding accuracy and the back sequence log-likelihood.

4.3.5 Analysis and discussion of results

Let's proceed by analysing the front and back footfalls separately before putting them together with conclusive thoughts.

4.3. EXPERIMENT 3: THE NECESSITY OF COMBINING THE FRONT AND BACK IMU SENSOR MEASUREMENTS

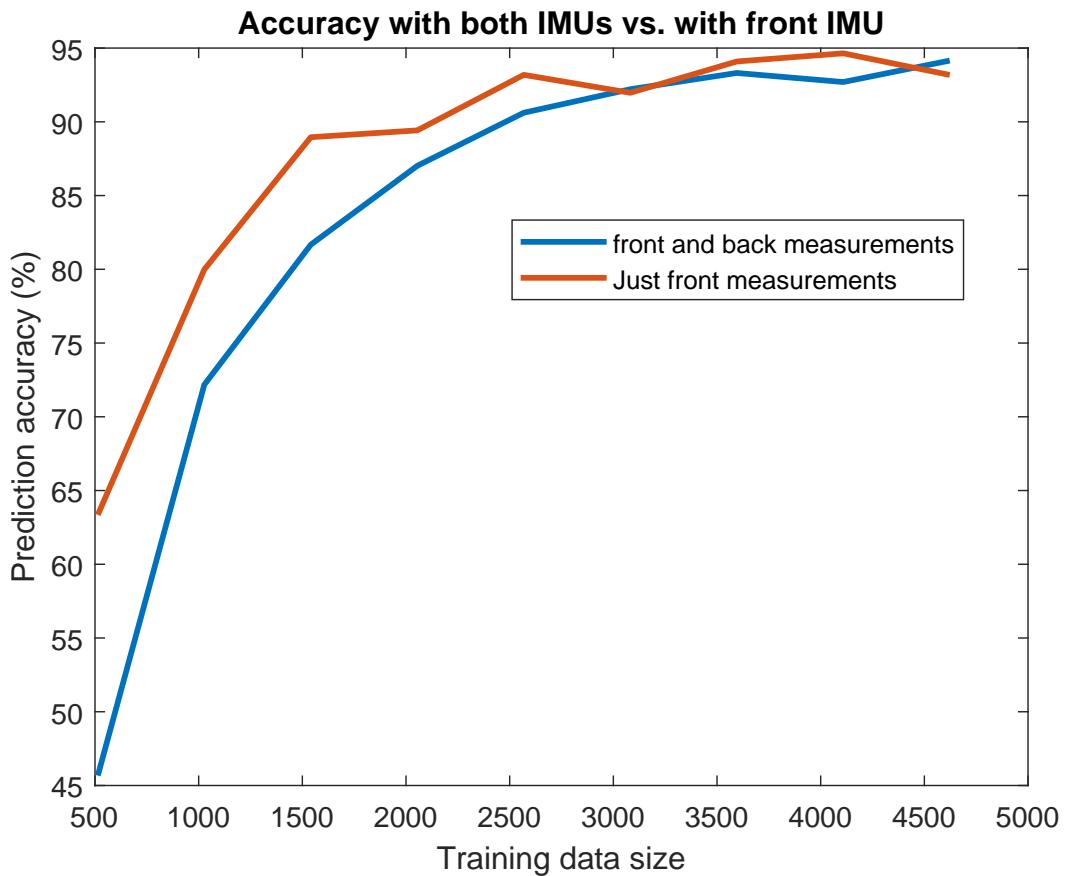


Figure 4.10: Front footfalls prediction accuracy of both IMUs vs with only the front IMU

Front footfalls

The experiment revealed that for both the combined and just the front IMU data, the classification accuracy increases with the increase in the training data size as testified by 4.10. Overall, building the model with just front sensor measurements achieves better accuracy except with about 3000 and 4500 observation samples.

The difference in performance is more significant with smaller data sizes. For instance, with 539 samples, the combination resulted in a very poor performance, below 50% accuracy, whereas, taking just the front measurements resulted in about 63%. This is at least 35% improvement. This difference is due to the dimensionality of the observation sequences as demonstrated in 4.1. With just the front measurements, we are dealing with 9 features against 18 features when the front and back measurements are considered. After about 3000 samples, the two models prediction accuracies become comparable. This is because there is sufficiently enough training data to cater for the high dimensionality.

4.3. EXPERIMENT 3: THE NECESSITY OF COMBINING THE FRONT AND BACK IMU SENSOR MEASUREMENTS

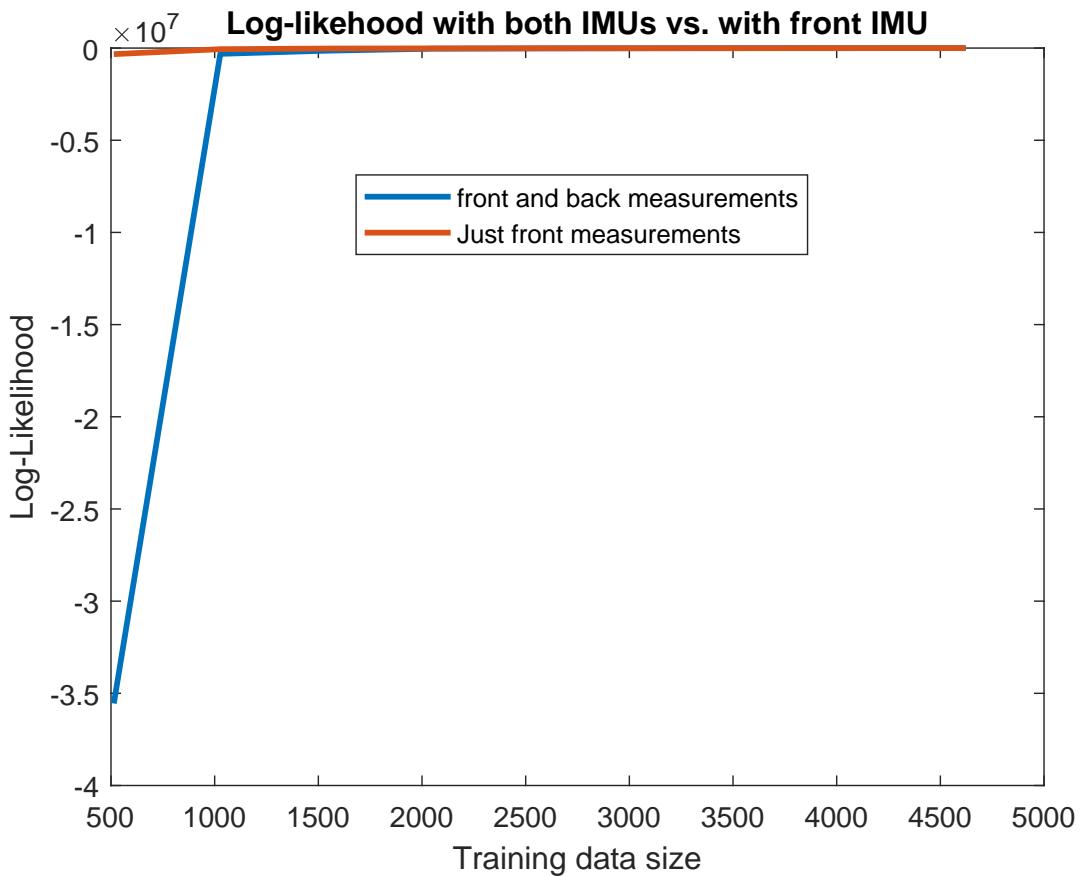


Figure 4.11: Front footfalls prediction log-likelihood with both IMUs vs with only the front IMU

The ability of the two models to recognise the IMU measurements generated by the dog are very comparable from 1000 samples onwards as shown in 4.11. Before 1000 samples, the model with just the front sensor measurements gives a higher likelihood value, constantly close to zero. This confirms the trend in the prediction accuracy.

Back footfalls

The back footfalls prediction with just the back sensor data and both the front and back sensor data show similar results to the front footfall case. However, the model with the combined sensors measurements catches up quicker at around 1000 samples. It even shows slightly better performance in accuracy after 2500 samples.

The common observation about both the front and back footfalls is the following. Only, the front and the back sensor measurements can better predict or recognise the front and back footfalls respectively, when the available observation sequence is relatively small,

4.3. EXPERIMENT 3: THE NECESSITY OF COMBINING THE FRONT AND BACK IMU SENSOR MEASUREMENTS

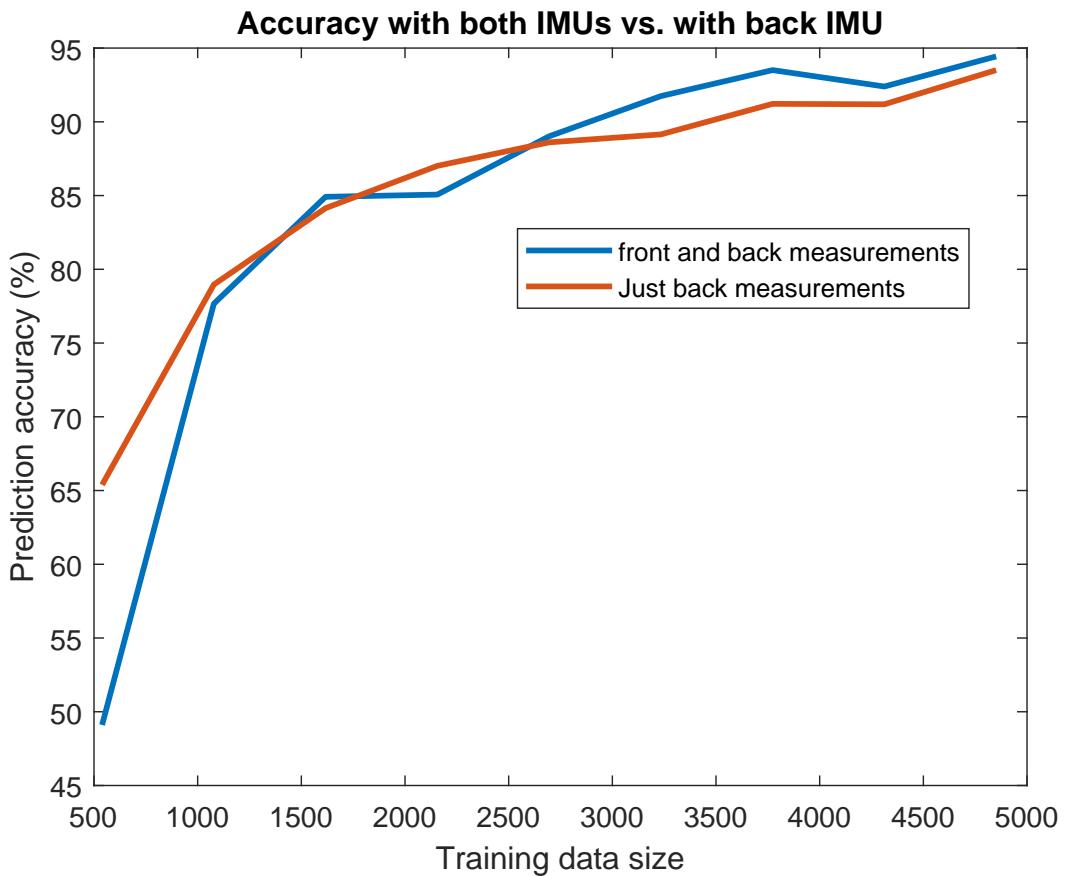


Figure 4.12: Back footfalls prediction accuracy with both IMUs vs with only the front IMU

below 1000. This finding is advantageous for several reasons. In fact, this realisation can be used to reduce the dimension of the data from 18 down to 9 when dealing with a small dataset. On one hand, this is not only good in terms of the model's accuracy. On the other hand, when dealing with just the front or the back legs, the required logistics such as the memory size, the unnecessary back or front sensors themselves.

4.3.6 Conclusions and recommendations of the experiment

In light of the finding above, it can be concluded that: using all the front and back IMU measurements to predict just the front and back footfalls does not necessarily improves the performance. It degrades the performance when there is not enough training data because of the high dimensionality. An improvement to this experiment would be to explore how the dimensionality reduction of the different datasets would influence the outcome. Moreover, a futher investigation how well can the back sensors measurement be used to predict the front footfalls. Similarly, how well can the front sensors be used

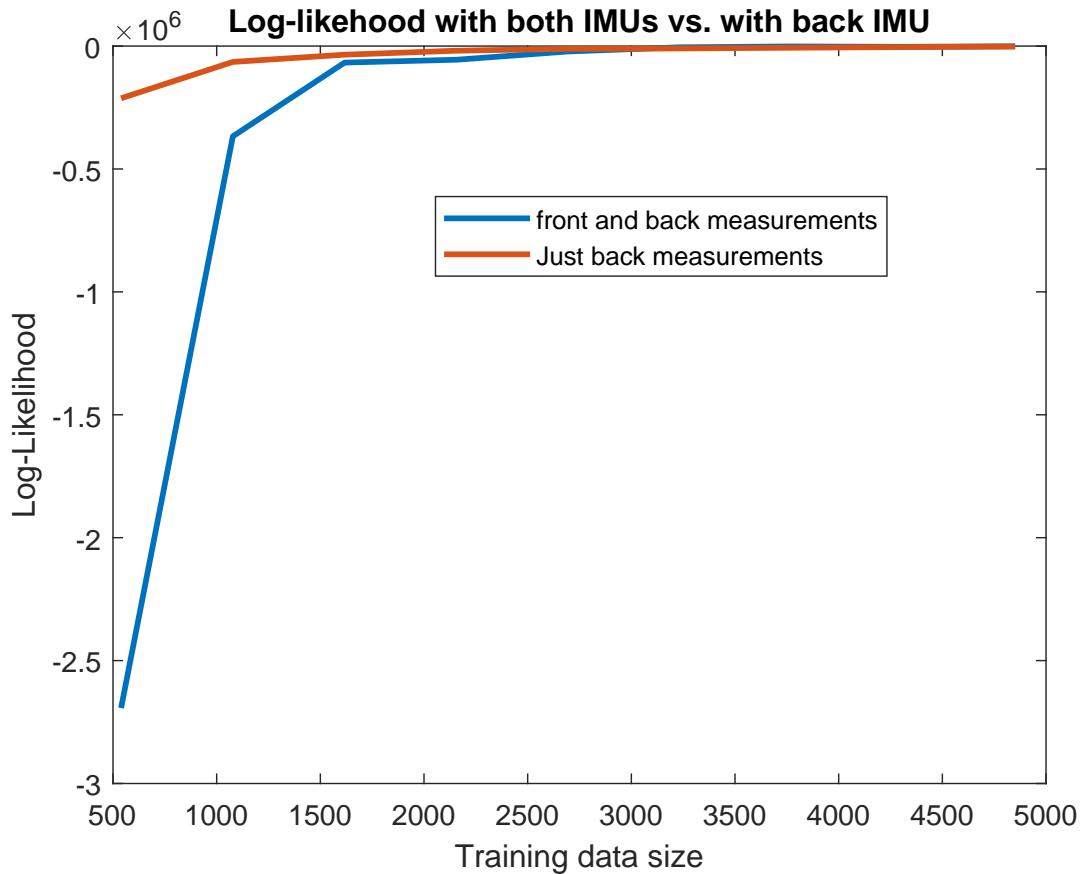


Figure 4.13: Back footfalls prediction log-likelihood with both IMUs vs with only the front IMU

to predict the back footfalls.

4.4 Experiment 4: Motion type recognition

4.4.1 Aim of the experiment

This experiment is about using IMU based HMM to identify a particular action of the dog. The three actions considered are: running, walking, and trotting. The hypothesis under investigation is, therefore, the following: ***IMU based HMMs can be used to successfully perform action recognition.***

4.4.2 Experiment apparatus

The main tools required to perform the experiment are listed next. Note that the experiment was performed for the front limbs only.

- Three distinct CHMM models: $\lambda_{running}$, $\lambda_{walking}$ and $\lambda_{trotting}$ to model the dog's run, walk, and trot.
- Three separate IMU datasets for the three action denoted by D_r , D_w , D_t respectively for running, walking and trotting.
- An action recognition criterion. In this experience, the prediction accuracies confusion matrix and the log-likelihood values were used.

4.4.3 Experiment procedure

The following steps were followed to run the experiment.

1. Extracting the three distinct datasets: This step is very similar to the pre-processing described here 3.4. The only difference is, the data for the dog's run, walk and trot were separated to obtain the three distinct sets. The initial sequences were mirrored three times for running walking motions and four times for the trotting action, to obtained large enough samples (please refer to this section 3 for more details on the mirroring process).
2. Constructing and training the models: 65% of the three distinct datasets were used to build and estimate the parameters of $\lambda_{running}$, $\lambda_{walking}$ and $\lambda_{trotting}$.
3. Testing and evaluating the models: Each model was tested three times with the remaining 35% of the three distinct datasets. The prediction accuracies and the log-likelihoods were recorded and tabulated as presented in the results sub-section. The log-likelihood value 4.2 measure the likelihood of the IMU measurement sequence having been emitted by the CHMM model in question. Thus, for a given sequence of measurements, can be attributed to the model that outputs the highest log-likelihood value [2]. The same purpose can be achieved by using the prediction accuracy confusion matrix, where the best candidate model is the model with the highest accuracy for given an observation sequence.

4.4.4 Experiment results

Table 4.2: the hidden state decoding accuracy and 4.1: the log-likelihood values, summary the results of the experiment.

	D_r	D_w	D_t
$\lambda_{running}$	91.16%	2.06%	0.22%
$\lambda_{walking}$	21.06%	100.00%	75.53%
$\lambda_{trotting}$	27.40%	45.72%	100.00%

Table 4.1: Footfall prediction confusion matrix (% accuracy)

	D_r	D_w	D_t
$\lambda_{running}$	0.00	-0.00×10^{14}	-0.00×10^{14}
$\lambda_{walking}$	-0.00×10^{14}	0.00	-0.00×10^{14}
$\lambda_{trotting}$	-1.44×10^{12}	-0.1302	0.00

Table 4.2: Footfall sequence log-likelihood

4.4.5 Analysis of results

By inspection, it can be observed that, the principal diagonal of the confusion matrix generates the highest accuracies in 4.1. It can, therefore, be concluded that the different models effectively recognise their corresponding motion type. This fact is confirmed by the log-likelihood table with the highest values on the main diagonal. Please note that the values -0.00×10^{14} are very tiny non-zero quantities.

It is worth mentioning the 100% accuracies when the unseen walking and trotting dataset were used to test their corresponding models. Strictly speaking, the 100% accuracies might have been obtained because the transition errors were ignored (ref. 3.7).

Regardless, this is a better performance on unseen dataset compared to the performance seen thus far when the dataset from three actions are aggregated (ref. 2). Naturally, the non-aggregated training dataset better models a particular action than building a common model that fits all the three action types.

4.4.6 Conclusions of the experiment

On the account of the findings discussed in the above subsection, the following conclusions can be drawn. Indeed, IMU based CHMM can be used to successfully perform motion type recognition.

Chapter 5

Discussion

Different experiments were performed and presented in the results section 4. This section further discusses the findings under the following headings.

5.1 Data dimensionality and bias analysis

This section is based on experiment 1 4.1 and 2 4.2. From the two experiments, it is clear that with a limited-size training data, building the model will all 18 features results in a poor performance. Thus, dimensionality reduction needs to be used to increase the model's precision. The best technique to do this is the combination of forward feature selection and feature ranking using separability index. More than 80% classification accuracy was obtained with four selected features using this filter technique.

However, up to 95% accuracy was obtained without dimensionality reduction when the training data size is relatively large. Therefore, it may not be necessary to reduce the feature size with enough data set, around 2000 samples. The robustness of the model without dimensionality reduction backs this statement. In fact, the variance error was extremely low both models: with and without dimensionality reduction. This is not surprising for two different reasons. Inherently, HMMs are very robust because they are probabilistic models. Moreover, the estimated means and the co-variances are the two potential sources of bias in the model. But, the mean is an unbiased estimator as proven in 5.1

$$E[\mu^*] = E\left[\frac{1}{N} \sum_{k=1}^N X^k\right] = \frac{1}{N} \sum_{k=1}^N E[X^k] = \mu \quad (5.1)$$

Furthermore, although variance is a biased estimator, they were normalised with $N-1$. Thus, with a large N value, the bias tends to zero as shown here in 5.2 and 5.3. So, with a large dataset, the gait sequence estimator is unbiased when the observation sequence is long and biased, otherwise. In fact, with small training data sequence, the bias increases with the feature size. This explains the very poor performance of the model built with 539 18-dimensional observations.

$$\text{Biased } E[\sigma^2] = E\left[\frac{1}{N} \sum_{k=1}^N (X^k - \mu^*)^2\right] = \frac{N-1}{N} \sigma^2 \quad (5.2)$$

$$\text{with a large } N, \quad E[\sigma^2] = \sigma^2 \text{ becomes unbiased} \quad (5.3)$$

5.2 Data aggregation and mirroring

Overall, the different models performed better with the increase in training data size. This fact shows that the data aggregation and the mirroring techniques worked very well. Nevertheless, the aggregation of the data from all the action types led to a loss of specificity as expected, therefore, decreasing the model's precision. This fact was uncovered by the motion recognition experiment 4.4. In this experiment, when specific models were built for each action, the prediction accuracies in table 4.1 were very close to 100%.

Even though the reverse sequences were appended, the models still performed very well when tested with the correct data. This testifies that the increase of data by mirroring does not reduce performance, contrary to the aggregation method.

5.3 State duration time

Chapter 6

Conclusions

The present report consolidates the work done on gait sequence modelling and estimation using Hidden Markov Models (HMMs). We went from the continuous IMU measurements of a dog moving, to gait state identification using HMM with Gaussian mixture density functions.

Given limitation imposed by the small size of the dataset, dimensionality reduction techniques were explored. In addition, data aggregation and mirroring techniques were employed to increase the data. After successfully implementing the various models, experiments were designed and performed to test the hypotheses formulated in the introduction. These hypotheses are repeated here as a reminder.

1. HMMs can successfully model gait sequence dynamics using IMU data, in the absence of huge training data.
2. Data aggregation and mirroring techniques can be used to overcome training data size limitations.
3. Dimensionality reduction can increase the performance of an HMM, in the absence of a large training set.
4. HMMs can be used to successfully perform gait activity recognition.

In light of the above results and the discussions, these hypotheses can be confidently accepted.

Indeed, continuous HMMs with Gaussian mixture can be used to build robust gait sequence estimators from IMU measurements with up to 95% precision. This accuracy is

obtainable by aggregating the measurements of the different gait actions and/or by using considering the reverse gait sequence as a valid sequence.

Furthermore, when dealing with a small training data, dimensionality reduction can increase the model's performance by up to 78%. However, with a large dataset, it might not be necessary to perform any reduction with just 18-dimensions. This last conclusion should be taken with some degree of reservation. Finally, the last experiment showed that IMU based HMM can be used for gait type classification.

The objectives of this project can, therefore, be considered met. In the next section recommendations are made for possible improvements and future works.

Chapter 7

Recommendations

In this chapter are presented the recommendations that follow from the completion of the project. These recommendations can serve as an extension of the project for improvement and future work.

- *Online real-time testing:* Even though the model was tested with unseen data, an online real-time implementation and testing on a moving dog would testify of the reliability of the model built.
- *Gather more data from different dogs:* The HMM model was formulated and the experiments were performed based on measurements from a single subject. More data should be gathered from multiple dogs to test the present and subsequently tune the model's parameters based on the outcome of the tests.
- *Incorporate more domain knowledge:* Given the limited amount of time, domain-specific knowledge was not thoroughly explored. Making more use of available knowledge on quadruped movement may improve the model's complexity and performance.
- *Compare the HMM models to other algorithms:* The performance of the constructed models should be compared to other classification and pattern recognition methods.
- *Combine the front and back HMMs:* In this design, separate models were built for the front and the back legs of the dog. Although this made the algorithm applicable to bipedal gait estimation, it would be useful to combine the two models and evaluate the performance of the holistic 16-states HMM.
- *Investigate more dimensionality reduction techniques:* The feature extraction methods performed did not perform very well, more dimensionality reduction methods such

as Fast Fourier Transform (FFT) [3], time-frequency analysis [4] should be investigated.

Bibliography

- [1] L. R. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition", *Proceedings of the IEEE*, vol. 77, pp. 257-286, 1989, ISSN 0018-9456.
- [2] Ghazaleh Panahandeh, Nasser Mohammadiha, Arne Leijon, Peter Handel, "Continuous Hidden Markov Model for Pedestrian Activity Classification and Gait Analysis", *Instrumentation and Measurement IEEE Transactions on*, vol. 62, pp. 1073-1083, 2013, ISSN 0890-6955, DOI 10.1109/TIM.2013.2247032
- [3] Guangyi Shi and Yuexian Zou and Yufeng Jin and Xiaole Cui and W. J. Li, "Towards HMM based human motion recognition using MEMS inertial sensors", *2008 IEEE International Conference on Robotics and Biomimetics*, pp. 1762-1766, 2009, DOI 10.1109/ROBIO.2009.54913268.
- [4] G. Panahandeh and N. Mohammadiha and A. Leijon and P. Handel, "Chest-mounted inertial measurement unit for pedestrian motion classification using continuous hidden Markov model", *2012 IEEE International Instrumentation and Measurement Technology Conference Proceedings*, pp. 991-995, 2012, ISSN 1091-5281, DOI 10.1109/I2MTC.2012.6229380.
- [5] I. P. I. Pappas, M. R. Popovic, T. Keller, V. Dietz, and M. Morari, A reliable gait phase detection system, *IEEE Trans. Neural Syst. Rehabil. Eng.*, vol. 9, no. 2, pp. 113-125, Jun. 2001.
- [6] C. Senanayake and S. Senanayake, Computational intelligent gait-phase detection system to identify pathological gait, *IEEE Trans. Inf. Technol. Biomed.*, vol. 14, no. 5, pp. 1173-1179, Sep. 2010.
- [7] R. Durbin and S. Eddy and A. Krogh and G. Mitchison, Markov chains and hidden Markov models", in *Biological sequence analysis*, UK: Cambridge University Press, 1998, pp. 46-79.

- [8] G. Pfundstein, "Hidden Markov Models with Generalised Emission Distribution for Analysis of High-Dimensional, Non-Euclidean Data," M.S. thesis, Dept. Stat. and Gene Center, Munich Univ., Munich, Germany, 2011.
- [9] Huseyin M. Ertunc and Kenneth A. Loparo and Hasan Ocak, "Tool wear condition monitoring in drilling operations using hidden Markov models (HMMs)", *International Journal of Machine Tools and Manufacture*, vol. 41, no. 9, pp. 1363-1384, 2001, ISSN 0018-9456, DOI [https://doi.org/10.1016/S0890-6955\(00\)00112-7](https://doi.org/10.1016/S0890-6955(00)00112-7). [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0890695500001127>
- [10] Mingxia Liu and Daoqiang Zhang, "Feature selection with effective distance", *Neurocomputing*, vol. 215, no. Supplement C, pp. 100-109, 2016, note SI: Stereo Data, ISSN 0925-2312, DOI <https://doi.org/10.1016/j.neucom.2015.07.155>. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0925231216306336>.
- [11] Ali Mahvash and Annie Ross, "Two-phase flow pattern identification using continuous hidden Markov model", *International Journal of Multiphase Flow*, vol. 34, no. 3, pp. 303-311, 2008, ISSN 0301-9322, DOI <https://doi.org/10.1016/j.ijmultiphaseflow.2007.08.006>. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0301932207001462>.
- [12] Saul Solorio-Fernandez and J. Ariel Carrasco-Ochoa and Jos Fco. Martnez-Trinidad, "A new hybrid filterwrapper feature selection method for clustering based on ranking", *Neurocomputing*, vol. 214, no. Supplement C, pp. 866-880, 2016, ISSN 0925-2312, DOI <https://doi.org/10.1016/j.neucom.2016.07.026>. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0925231216307718>.
- [13] Hui-Huang Hsu and Cheng-Wei Hsieh and Ming-Da Lu, "Hybrid feature selection by combining filters and wrappers", *Expert Systems with Applications*, vol. 38, no. 7, pp. 8144-8150, 2011, ISSN 0957-4174, DOI <https://doi.org/10.1016/j.eswa.2010.12.156>. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0957417410015198>.
- [14] Jeong-Su Han and Sang Wan Lee and Zeungnam Bien, "Feature subset selection using separability index matrix", *Information Sciences*, vol. 223, no. Supplement C, pp. 102-118, 2013, ISSN 0020-0255, DOI <https://doi.org/10.1016/j.ins.2012.09.042>. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0020025512006354>.
- [15] Zhao Yongli and Zhang Yungui and Tong Weiming and Chen Hongzhi, "An improved feature selection algorithm based on MAHALANOBIS distance for Network Intrusion

Detection”, *PROCEEDINGS OF 2013 International Conference on Sensor Network Security Technology and Privacy Communication System*, pp. 69-73, 2013, DOI 10.1109/SNS-PCS.2013.6553837.

- [16] Dimension Reduction - Pattern Recognition Tools.[Online]. Available: <http://37steps.com/prtools/examples/dimension-reduction/>
- [17] Youness Aliyari Ghassabeh and Frank Rudzicz and Hamid Abrishami Moghaddam, ”Fast incremental LDA feature extraction”, *Pattern Recognition*, vol. 48, no. 6, pp. 1999-2012, 2015, ISSN 0031-3203, DOI <https://doi.org/10.1016/j.patcog.2014.12.012>. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0031320314005214>.
- [18] Miguel Simo and Pedro Neto and Olivier Gibaru, ”Using data dimensionality reduction for recognition of incomplete dynamic gestures”, *Pattern Recognition Letters*, vol. 99, no. Supplement C, pp. 32-38, 2017, ISSN 0167-8655, DOI <https://doi.org/10.1016/j.patrec.2017.01.003>. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S016786551730003X>.
- [19] Timothy M. Griffin and Russell P. Main and Claire T. Farley, ”Biomechanics of quadrupedal walking: how do four-legged animals achieve inverted pendulum-like movements?”, *The Journal of Experimental Biology*, vol. 207, pp. 3545-3558, 2004, DOI 10.1242/jeb.01177.
- [20] Feature Selection - MathWorks.[Online]. Available: <https://www.mathworks.com/help/stats/feature-selection.html>
- [21] K-Nearest Neighbor Classifier - Pattern Recognition Tools.[Online]. Available: <http://www.37steps.com/prhtml/prtools/knnc.html>
- [22] hmmestimate - MathWorks.[Online]. Available: <https://www.mathworks.com/help/stats/hmmestimate.html>
- [23] McLachlan, G., and D. Peel, Finite Mixture Models. Hoboken, NJ: *John Wiley & Sons*, Inc., 2000.
- [24] fitgmdist - MathWorks.[Online]. Available: <https://www.mathworks.com/help/stats/fitgmdist.html>
- [25] Ljung, L. System Identification: Theory for the User, Upper Saddle River, NJ, Prentice-Hall PTR, 1999
- [26] Akaike’s Information Criterion for estimated model - MATLAB aic.[Online]. Available: <https://www.mathworks.com/help/ident/ref/aic.html>
- [27] Bayes Information Criterion - MATLAB - MathWorks.[Online]. Available: <https://www.mathworks.com/help/stats/gmdistribution.bic.html>

BIBLIOGRAPHY

- [28] Qiuqiang Kong, "matlab-hmm".[Online]. Available:
<https://github.com/qiuqiangkong/matlab-hmm>
- [29] Trainable mapping for forward feature selection - Pattern Recognition Tools. [Online]
. <http://www.37steps.com/prhtml/prtools/featself.html>
- [30] Mahalanobis distance - Pattern Recognition Tools. [Online]
<http://www.37steps.com/prhtml/prtools/distmaha.html>
- [31] Principal component analysis (PCA or MCA on overall covariance matrix) - Pattern Recognition Tools. [Online] <http://www.37steps.com/prhtml/prtools/pcam.html>
- [32] Optimal discrimination linear mapping (Fisher mapping, LDA) - Pattern Recognition Tools. [Online]. <http://www.37steps.com/prhtml/prtools/fisherm.html>
- [33] Create cross validation partition for data - MATLAB - MathWorks.[Online].
<https://www.mathworks.com/help/stats/cvpartition.html>
- [34] Random sampling of datasets for training and testing - Pattern Recognition Tools.
[Online]. <http://www.37steps.com/prhtml/prtools/gendat.html>
- [35] PRTools - A Matlab toolbox for pattern recognition.[Online]. <http://prtools.org/>
- [36] Statistics and Machine Learning Toolbox - MATLAB - MathWorks.[Online].
<https://www.mathworks.com/products/statistics.html>
- [37] K. Fukunaga, "Introduction to statistical pattern recognition", 2nd ed., *Academic Press*, New York, 1990.
- [38] C. Liu and H. Wechsler, "Robust Coding Schemes for Indexing and Retrieval from Large Face Databases", *IEEE Transactions on Image Processing*, vol. 9, no. 1, 2000, 132-136.

Appendix A

Additional Files and Schematics

A.1 Implementation and experiment code

The implementation can be accessed from this github repository: <https://github.com/h-kouame/gait-sequence-modelling-and-estimation-v2>

A.2 Additional results

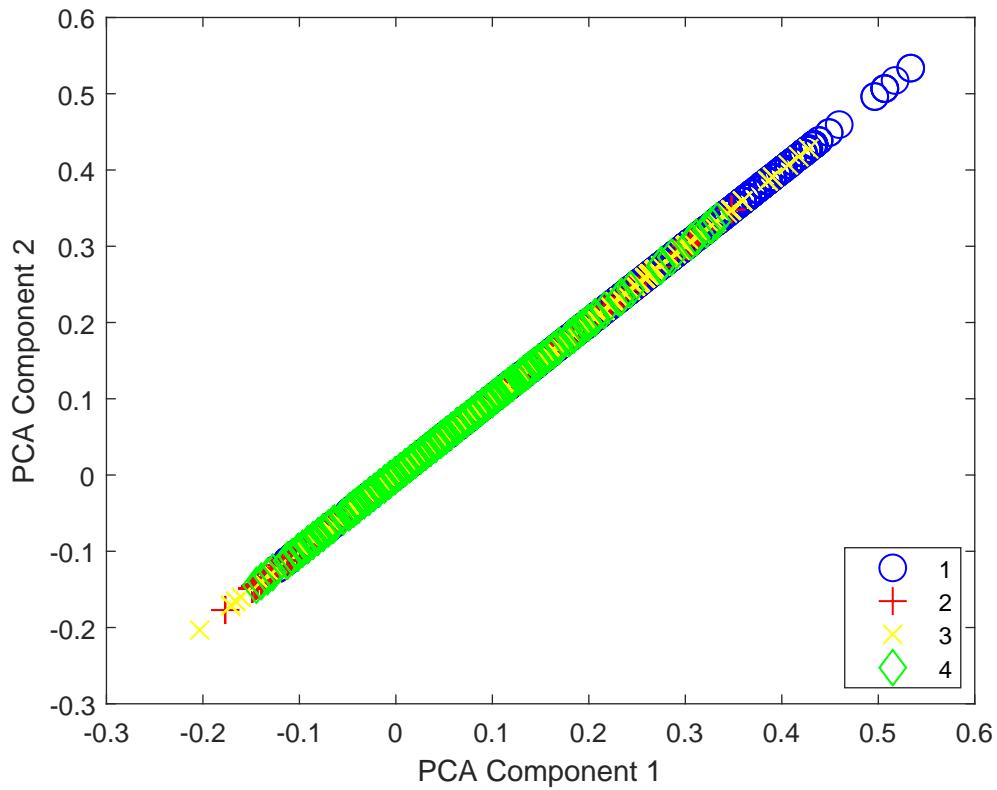


Figure A.1: Scatter plot of 1-dimensional observations grouped according the ground-truth state sequence

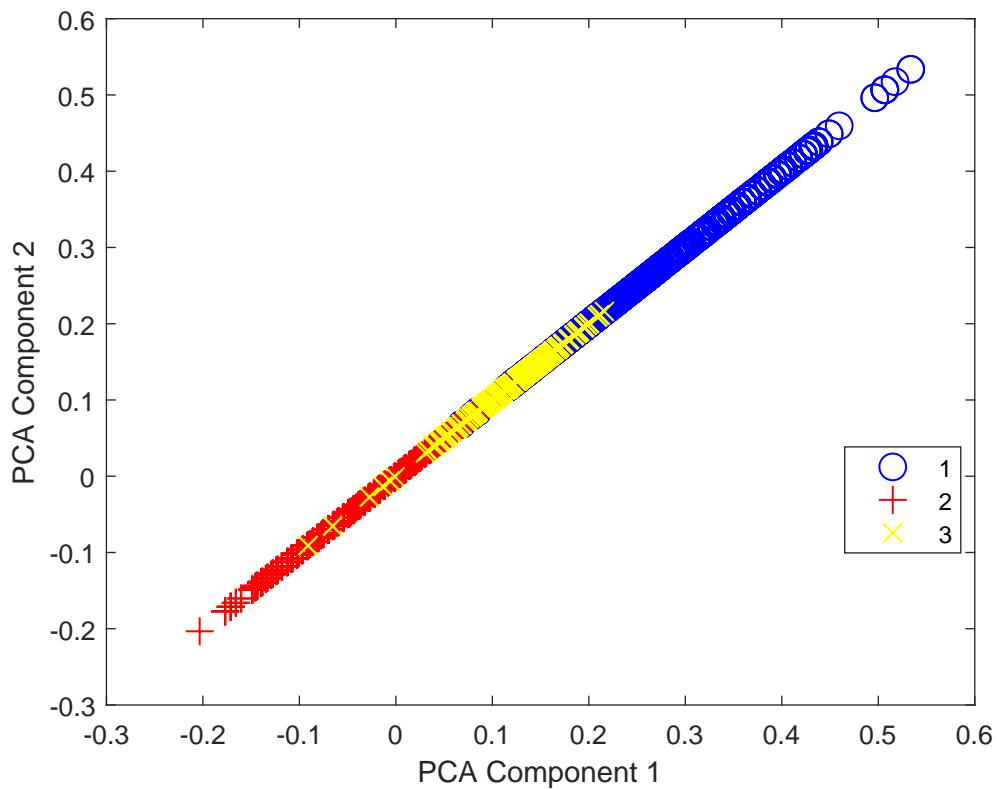


Figure A.2: Scatter plot of 1-dimensional observations grouped according to the estimated state sequence

A.2. ADDITIONAL RESULTS

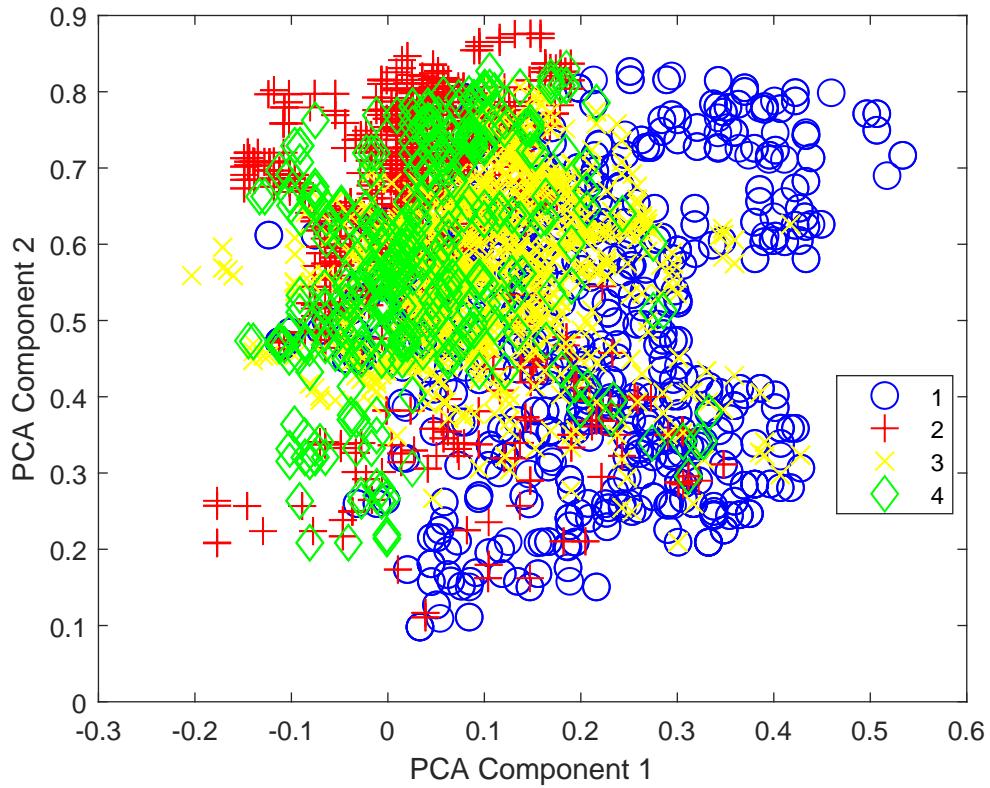


Figure A.3: Scatter plot of 2-dimensional observations grouped according to the ground-truth state sequence

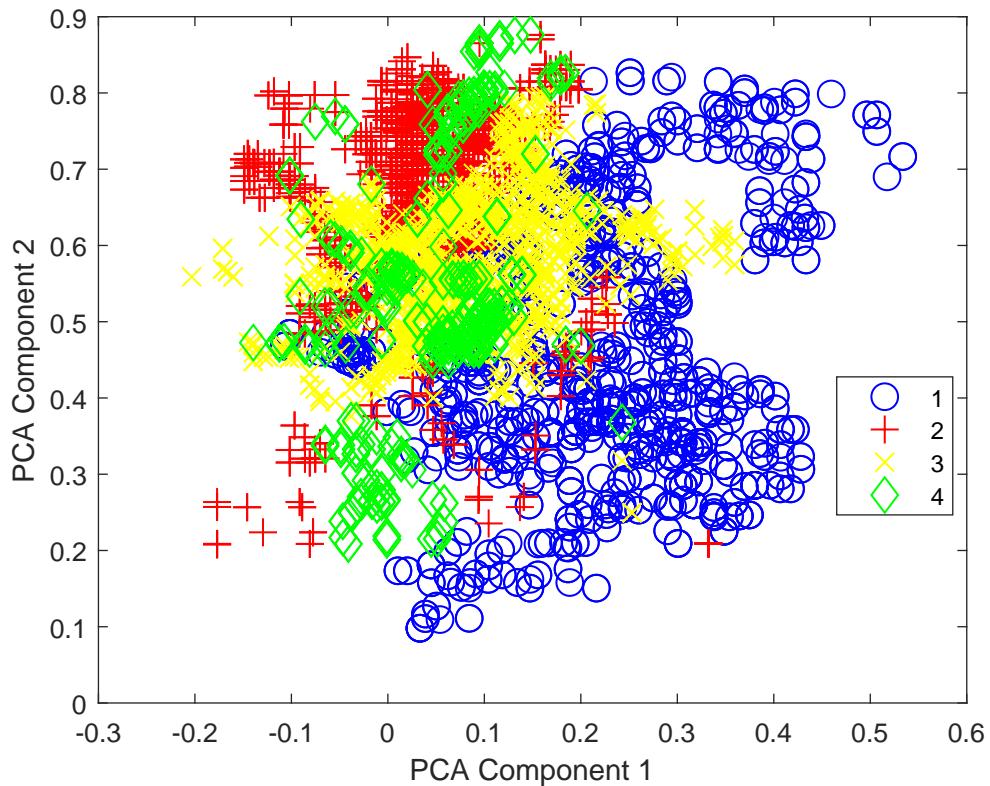


Figure A.4: Scatter plot of 2-dimensional observations grouped according to the estimated state sequence

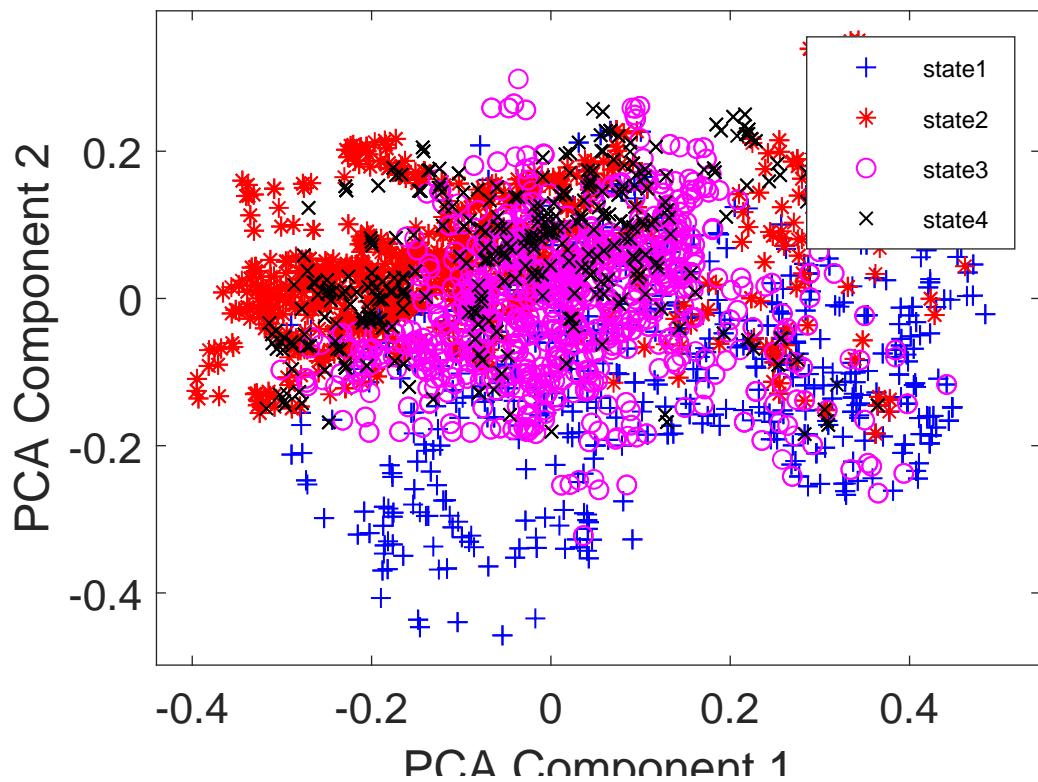


Figure A.5: Scatter plot of 3-dimensional observations grouped according the ground-truth state sequence

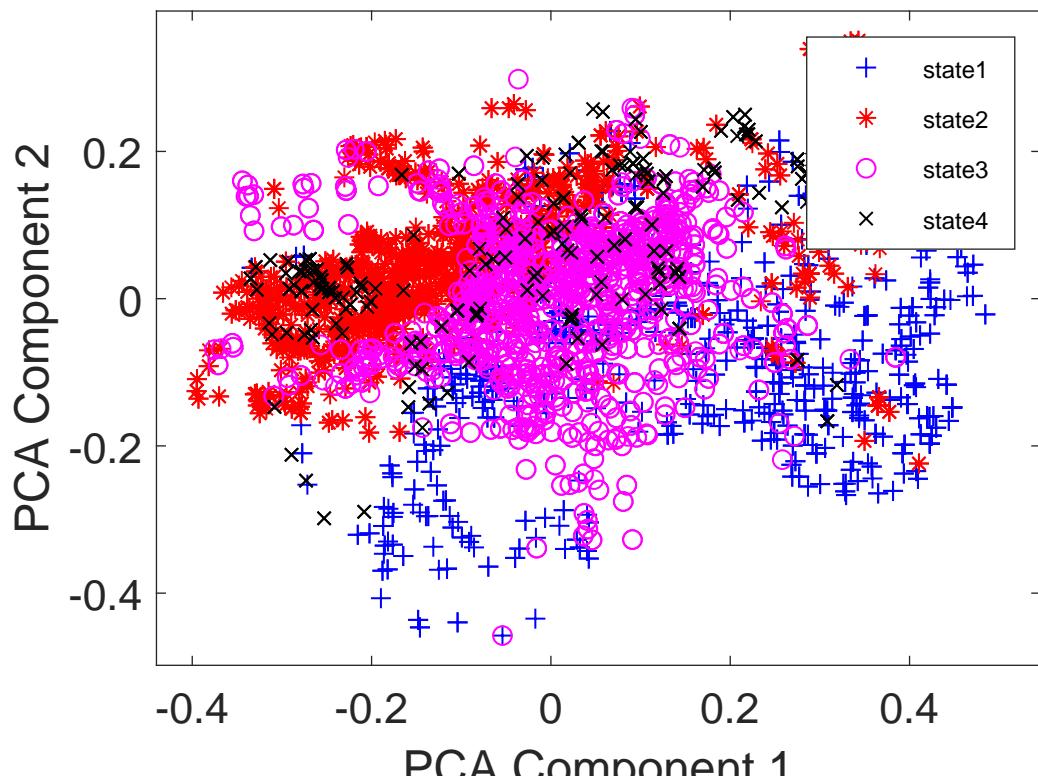


Figure A.6: Scatter plot of 3-dimensional observations grouped according to the estimated state sequence

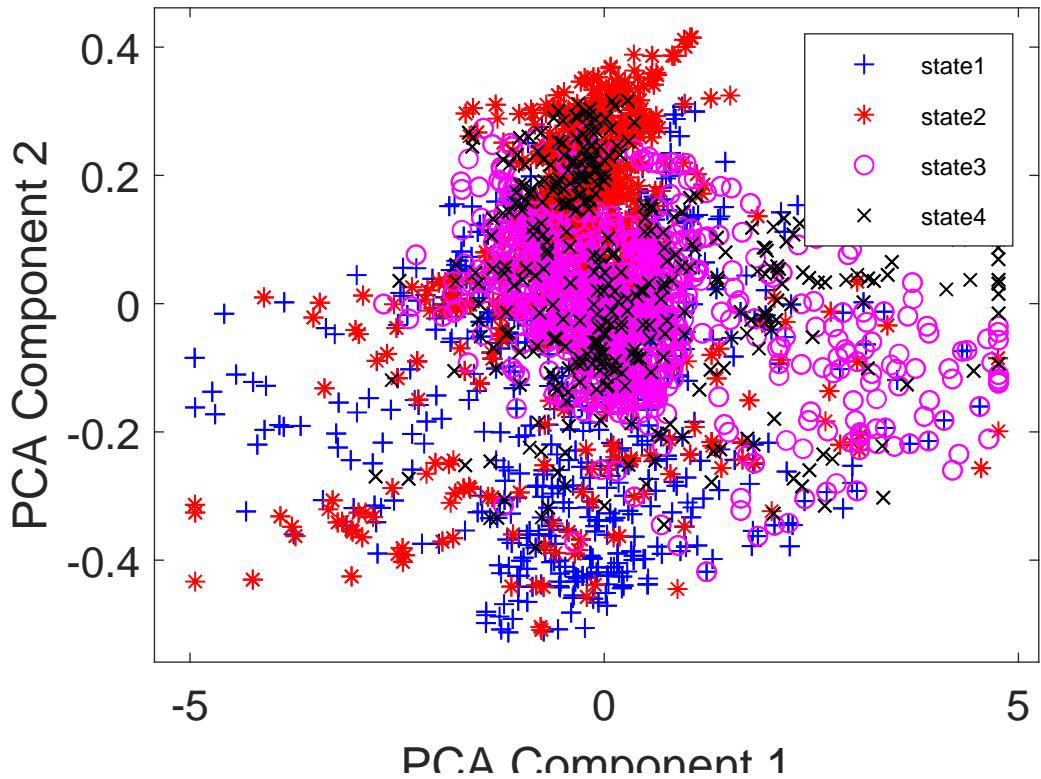


Figure A.7: Scatter plot of 4-dimensional observations grouped according to the ground-truth state sequence

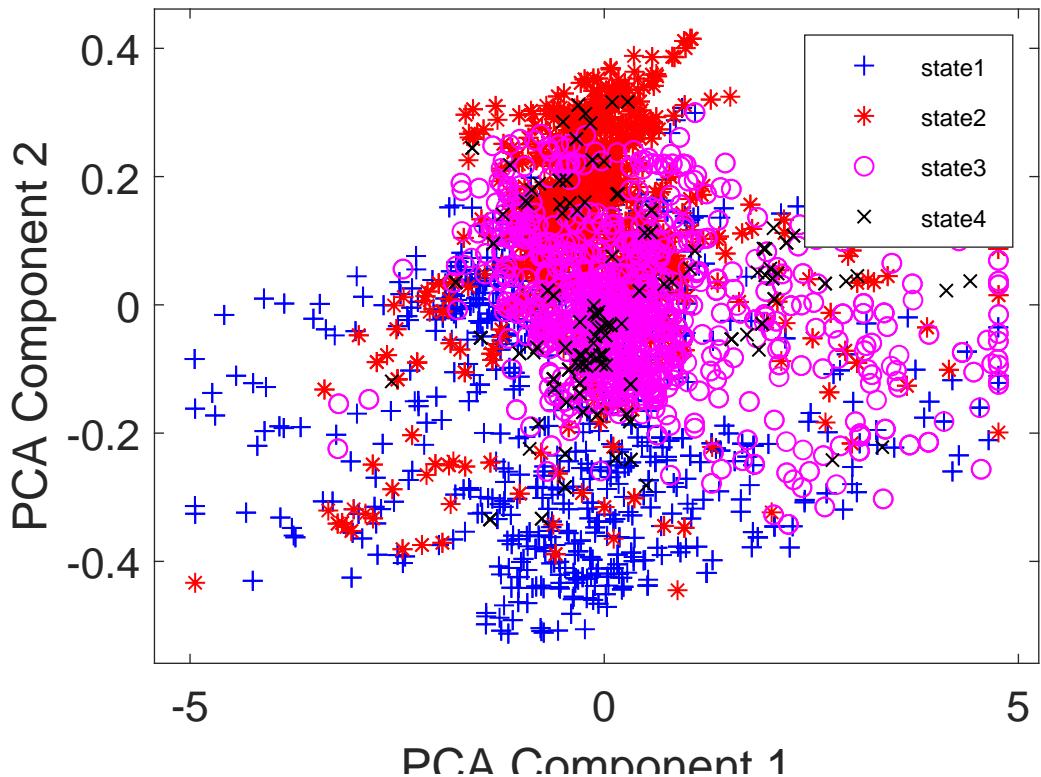


Figure A.8: Scatter plot of 4-dimensional observations grouped according to the estimated state sequence

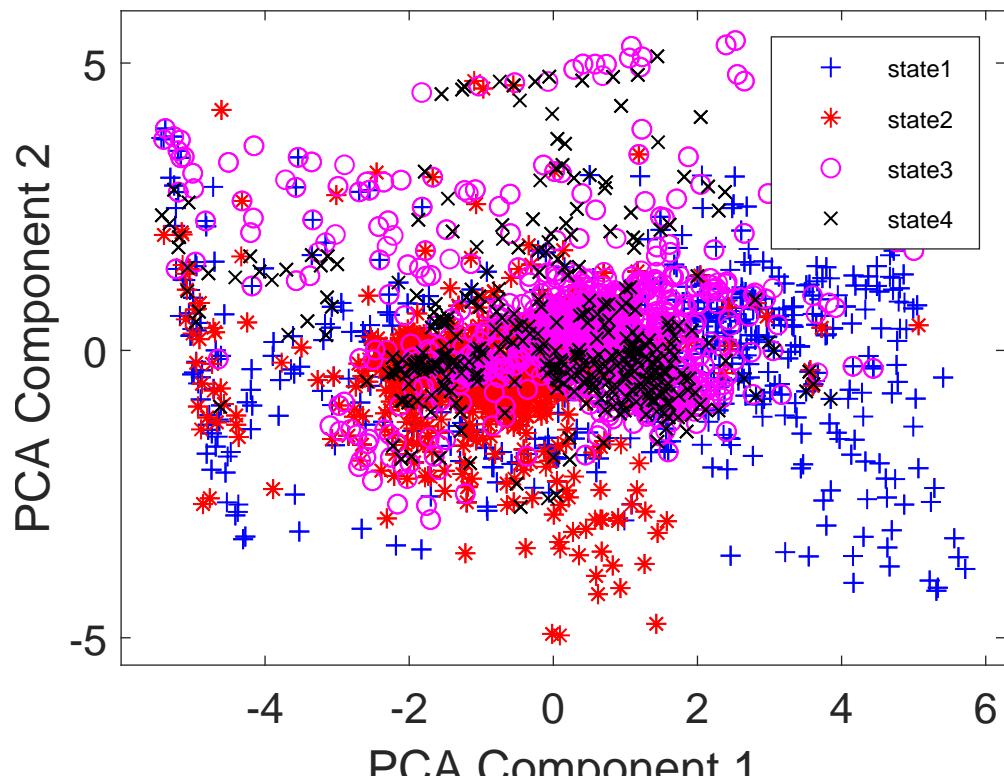


Figure A.9: Scatter plot of 6-dimensional observations grouped according to the ground-truth state sequence

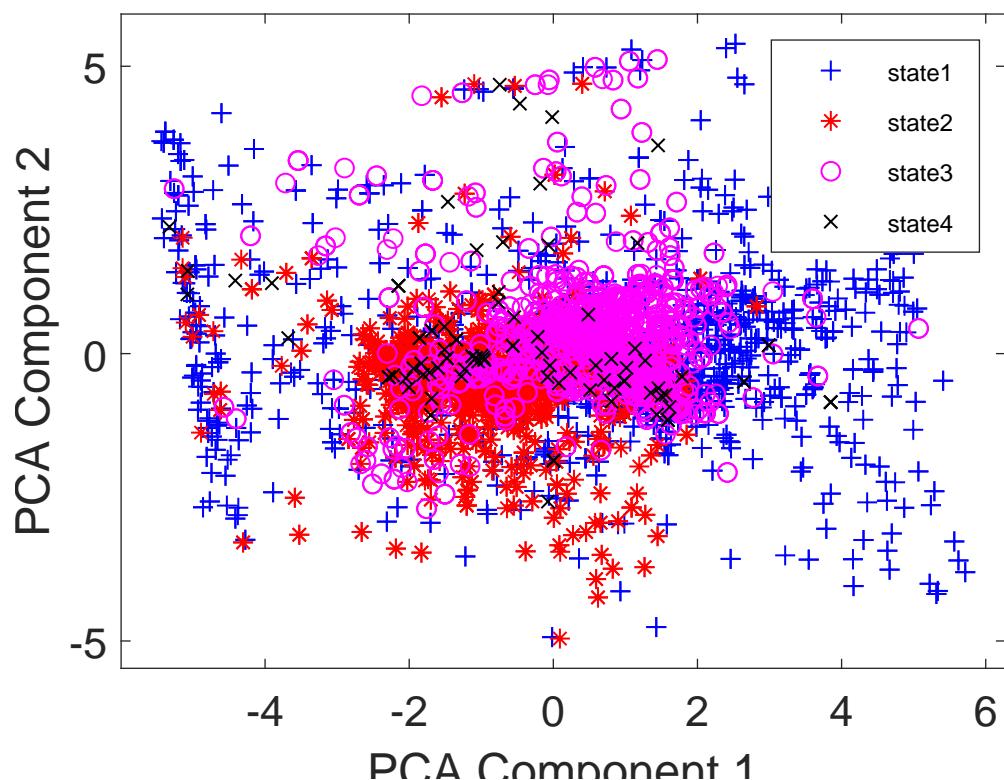


Figure A.10: Scatter plot of 6-dimensional observations grouped according to the estimated state sequence

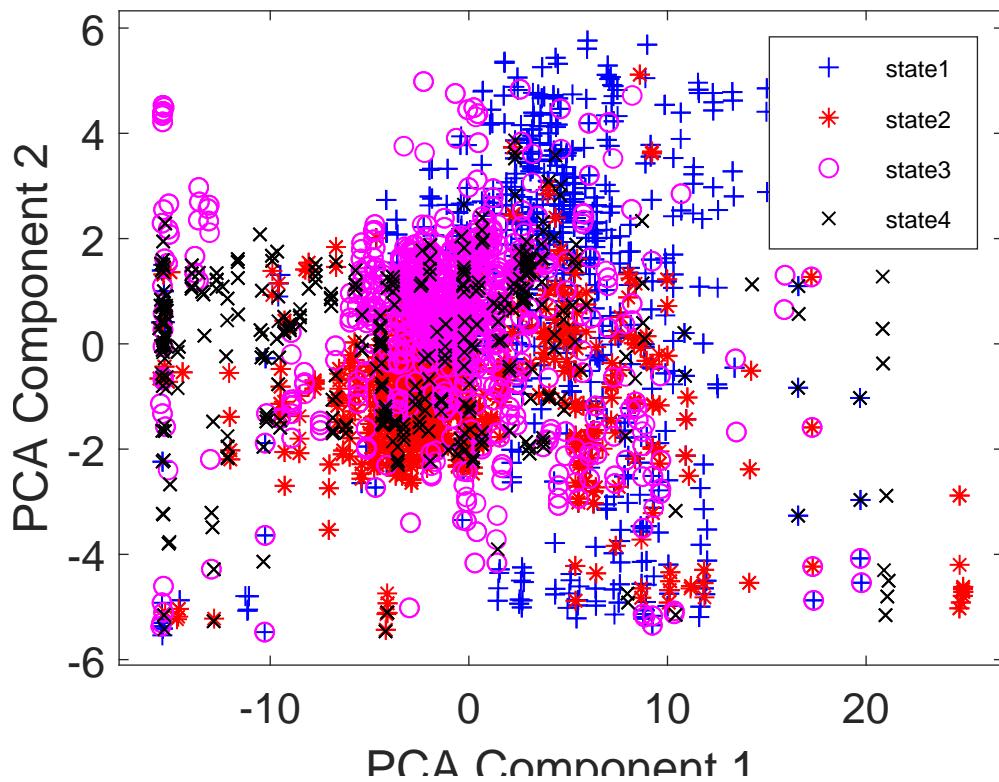


Figure A.11: Scatter plot of 7-dimensional observations grouped according to the ground-truth state sequence

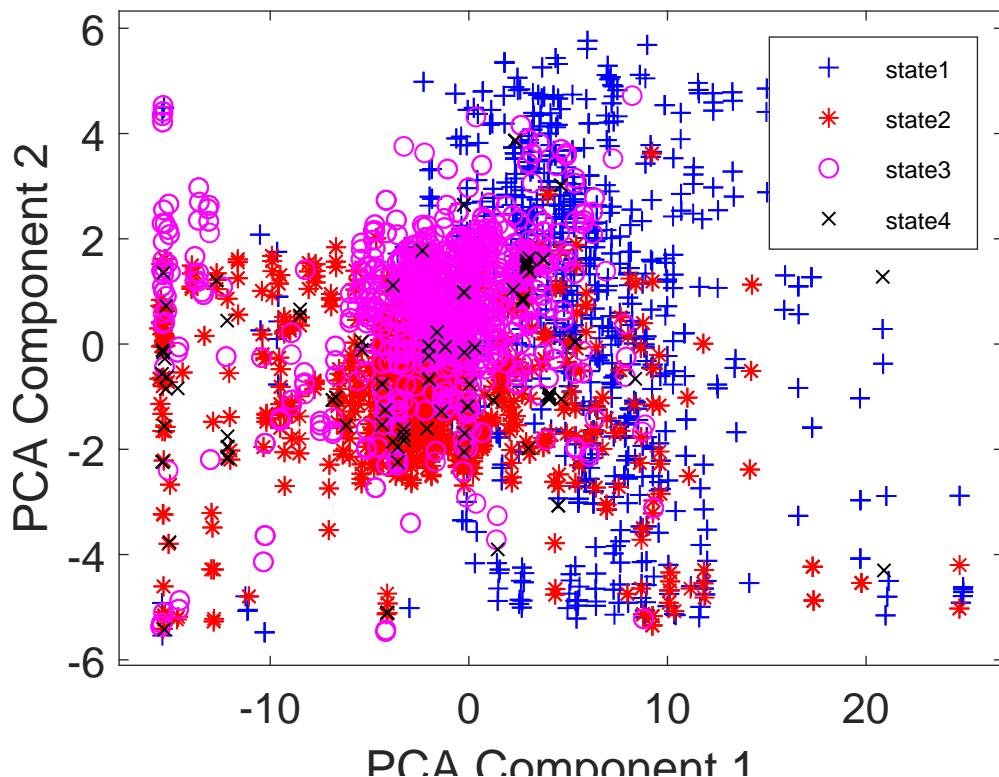


Figure A.12: Scatter plot of 7-dimensional observations grouped according to the estimated state sequence

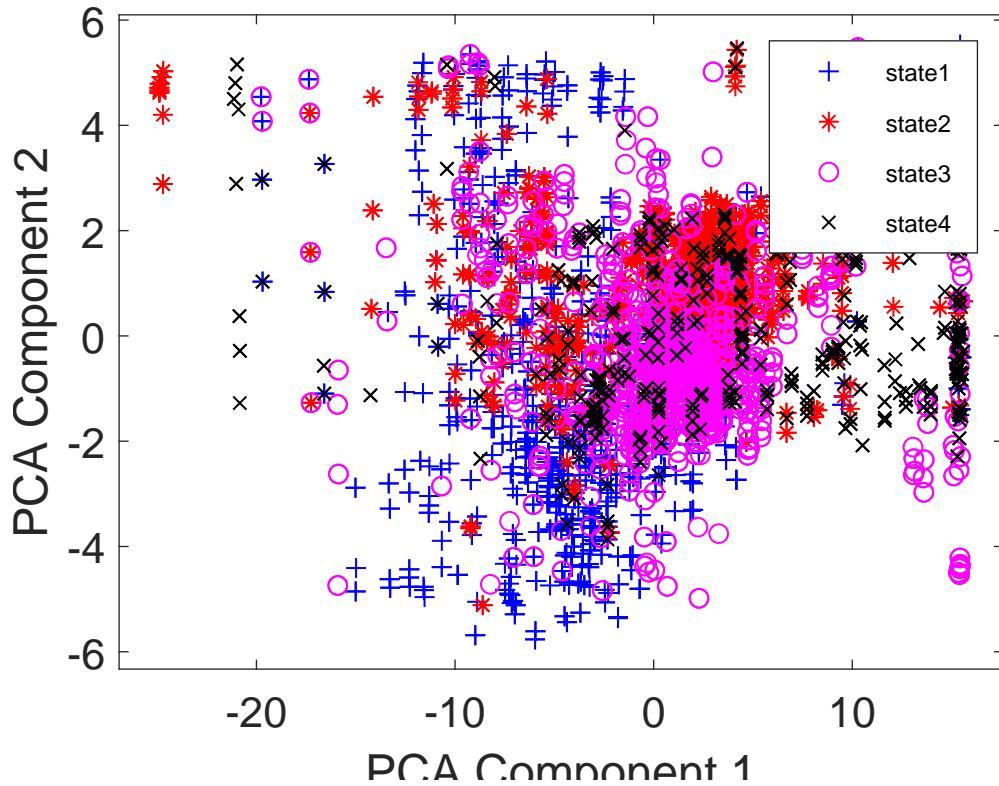


Figure A.13: Scatter plot of 8-dimensional observations grouped according to the ground-truth state sequence

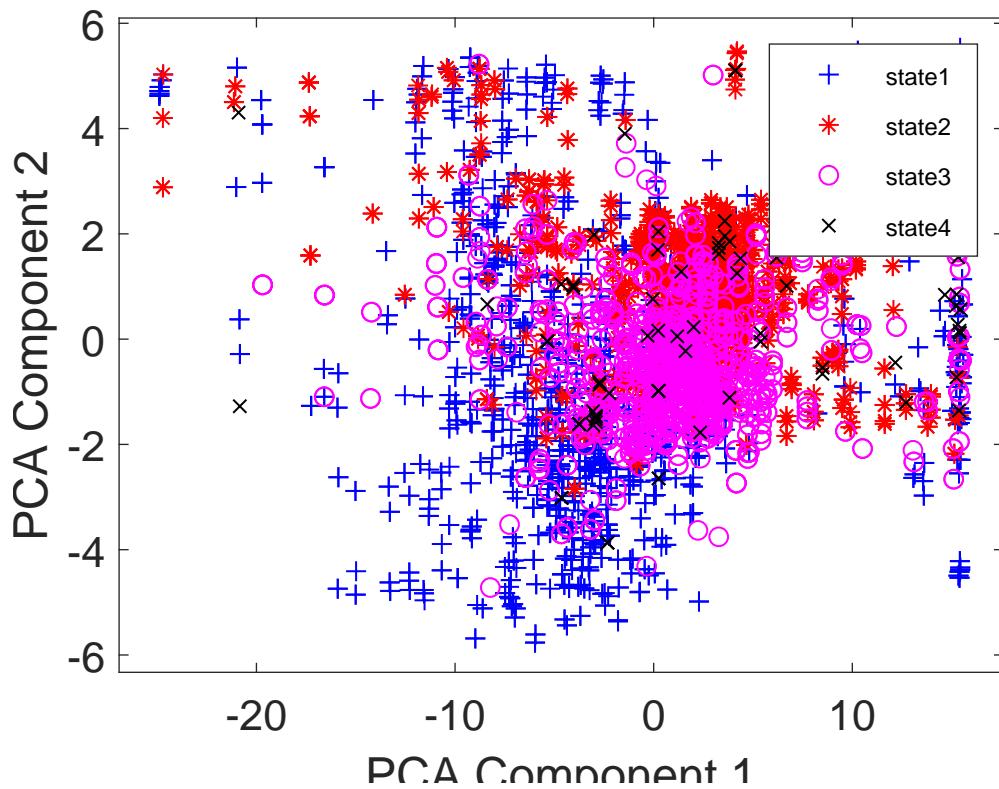


Figure A.14: Scatter plot of 8-dimensional observations grouped according to the estimated state sequence

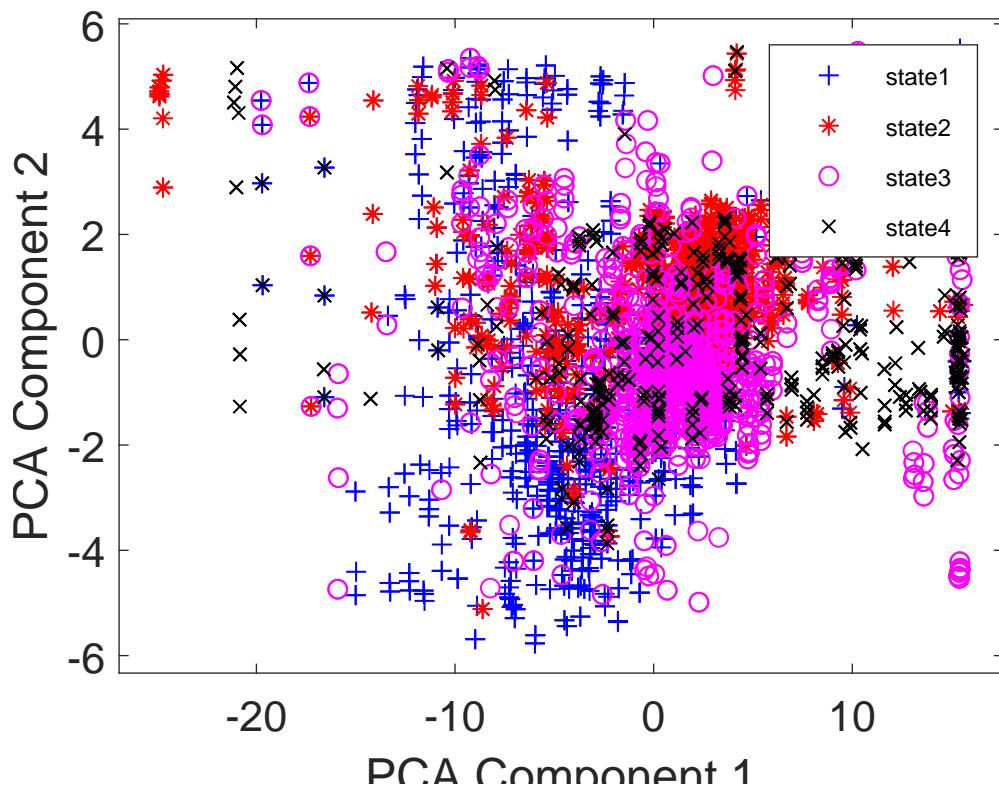


Figure A.15: Scatter plot of 9-dimensional observations grouped according to the ground-truth state sequence

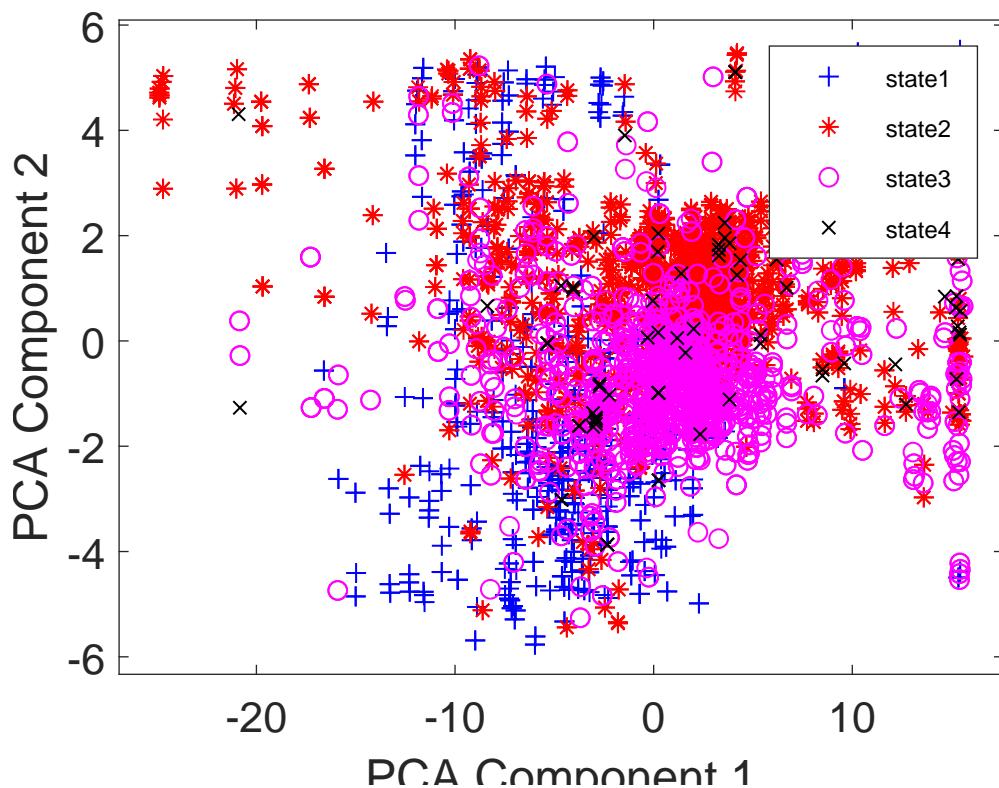


Figure A.16: Scatter plot of 9-dimensional observations grouped according to the estimated state sequence

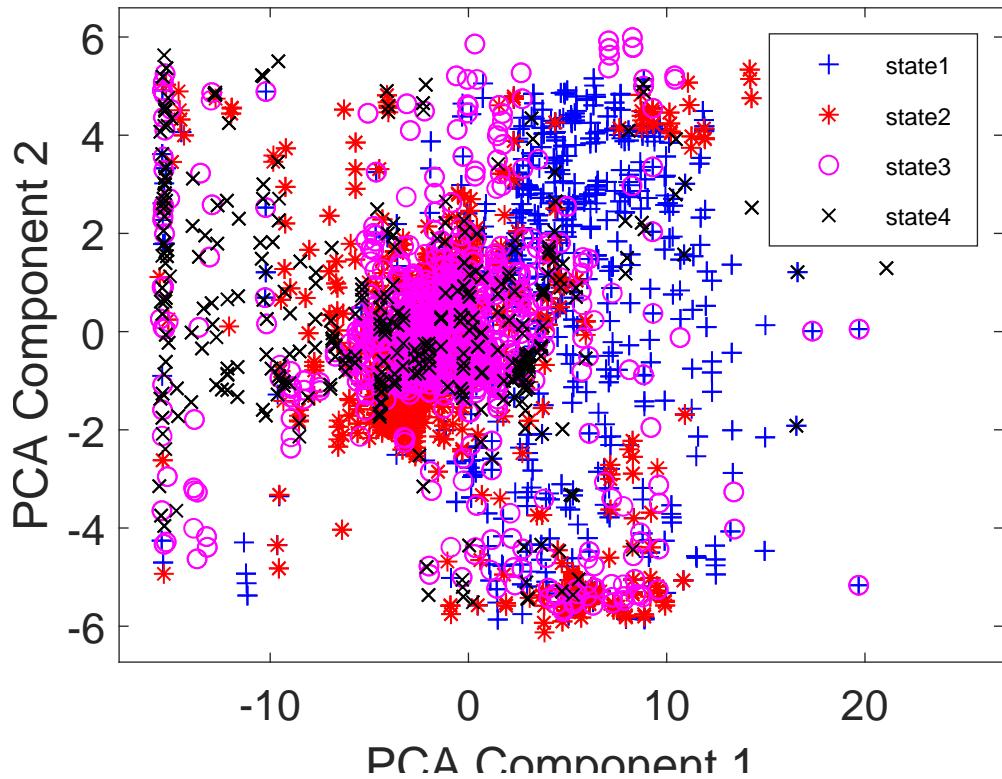


Figure A.17: Scatter plot of 10-dimensional observations grouped according to the ground-truth state sequence

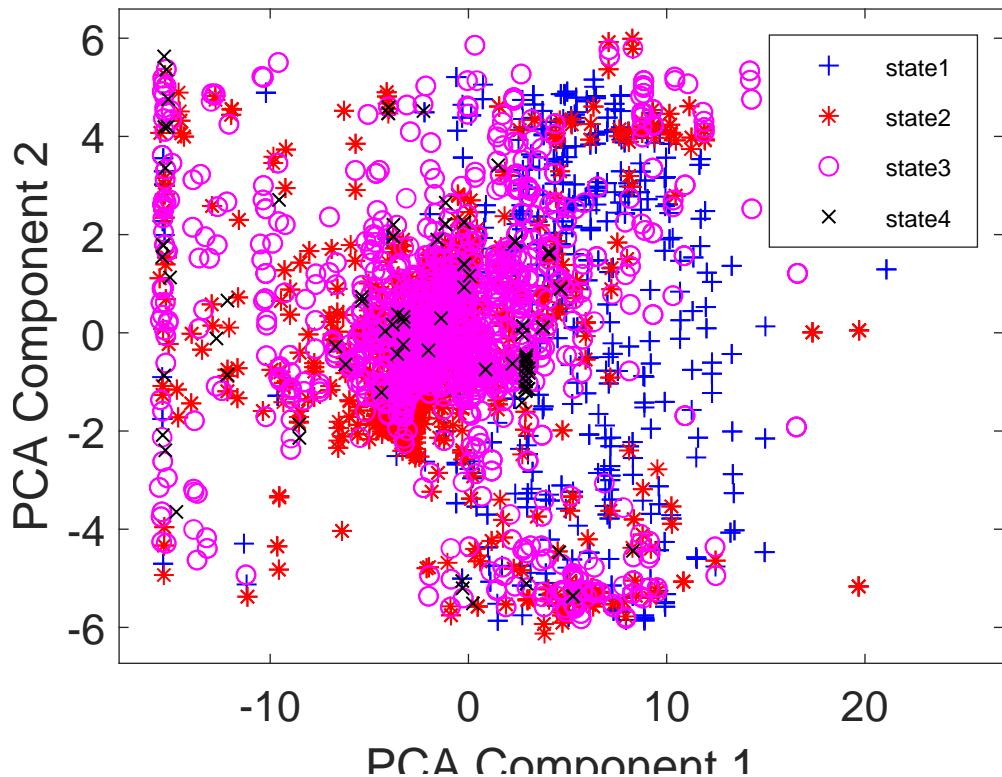


Figure A.18: Scatter plot of 10-dimensional observations grouped according to the estimated state sequence

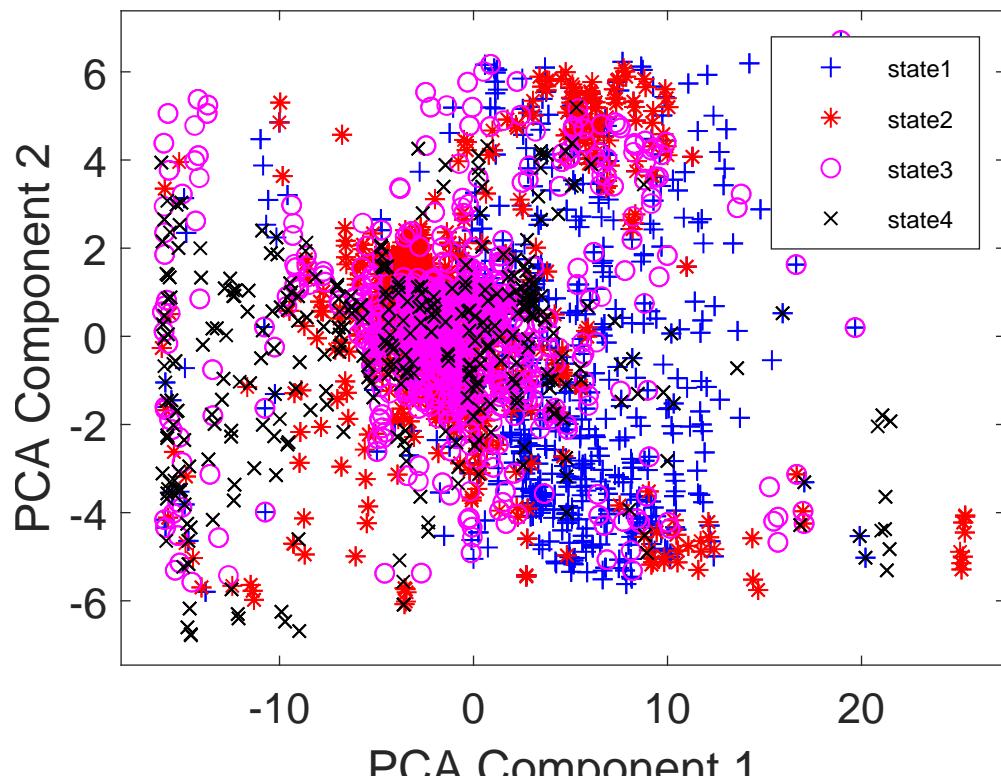


Figure A.19: Scatter plot of 11-dimensional observations grouped according to the ground-truth state sequence

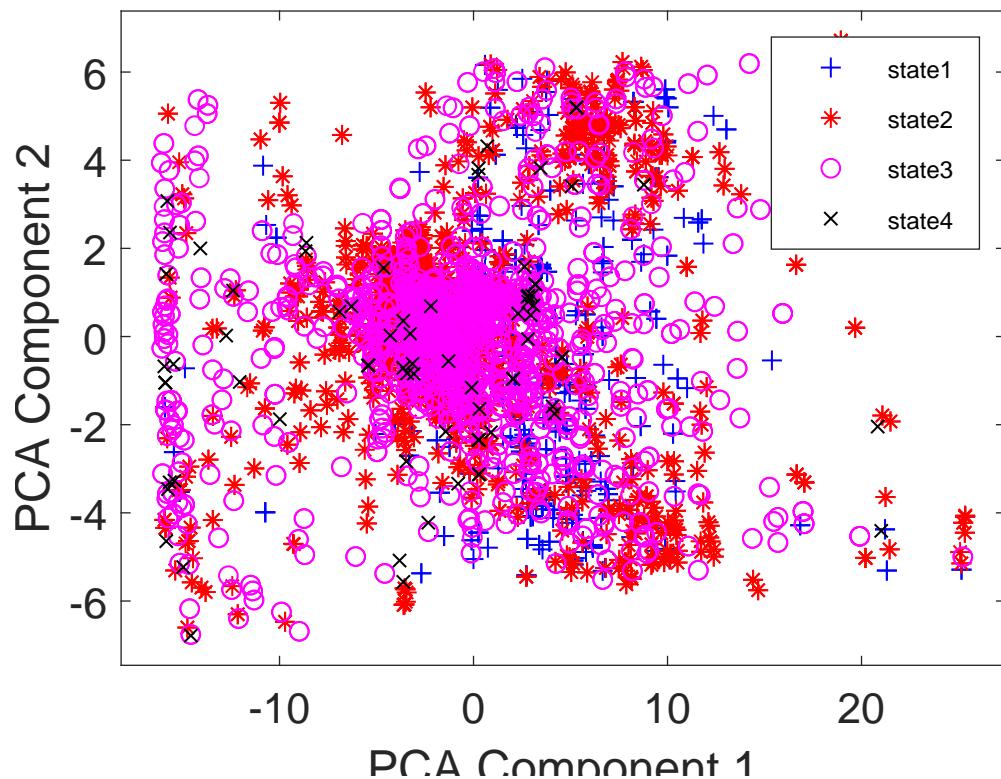


Figure A.20: Scatter plot of 11-dimensional observations grouped according to the estimated state sequence

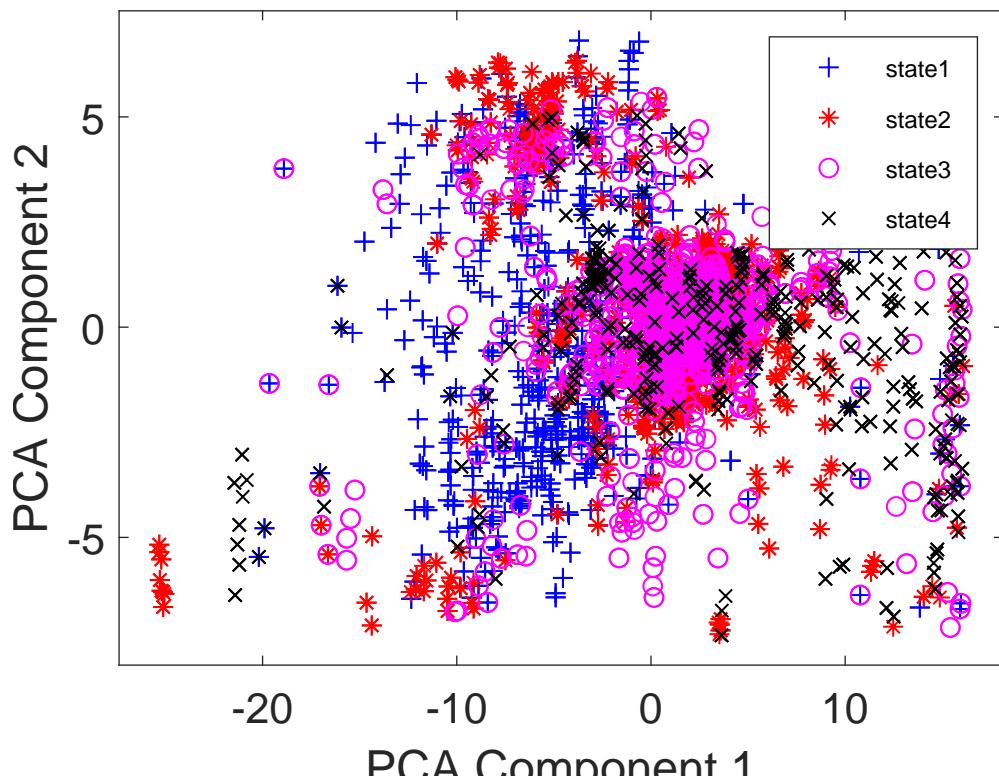


Figure A.21: Scatter plot of 12-dimensional observations grouped according to the ground-truth state sequence

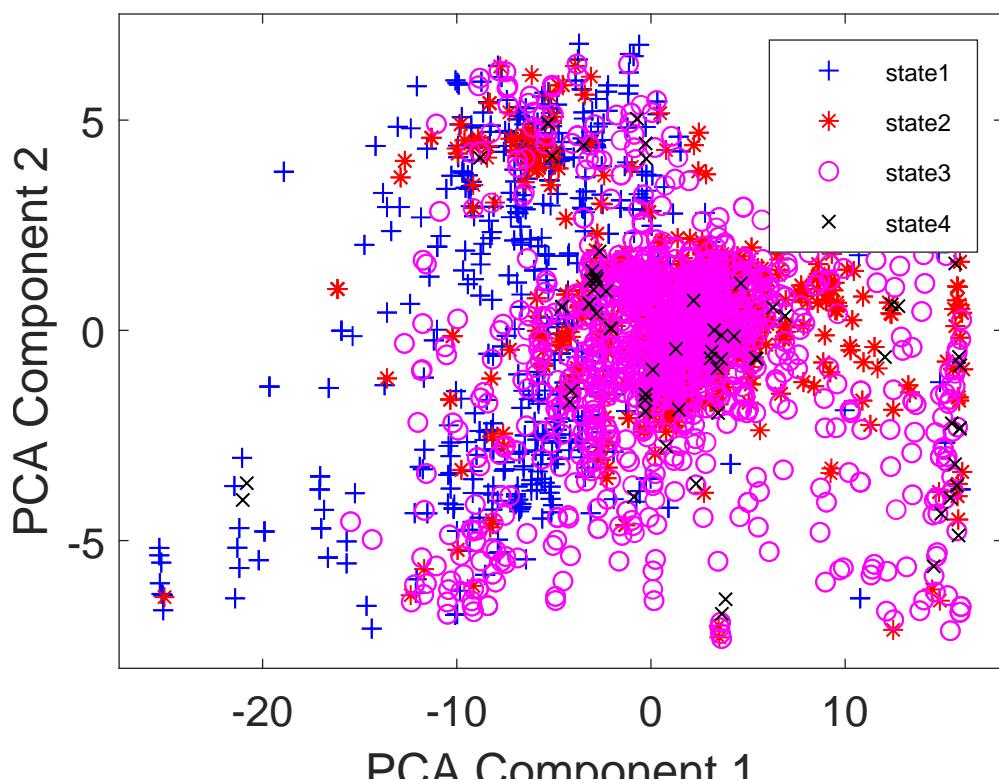


Figure A.22: Scatter plot of 12-dimensional observations grouped according to the estimated state sequence

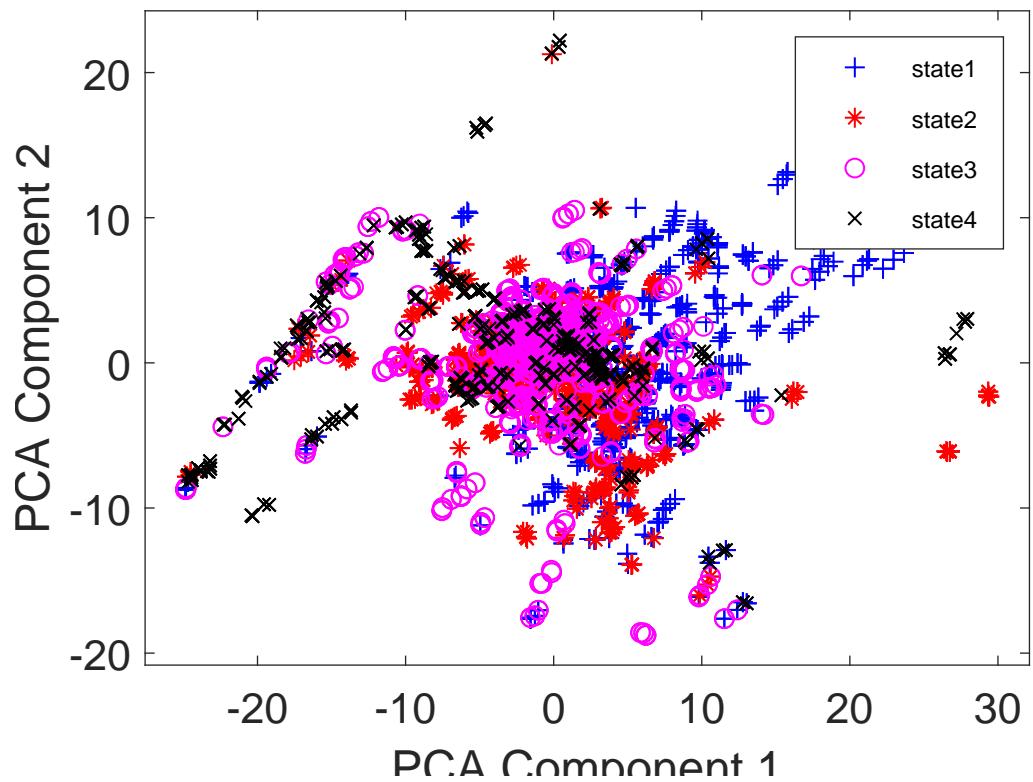


Figure A.23: Scatter plot of 13-dimensional observations grouped according to the ground-truth state sequence

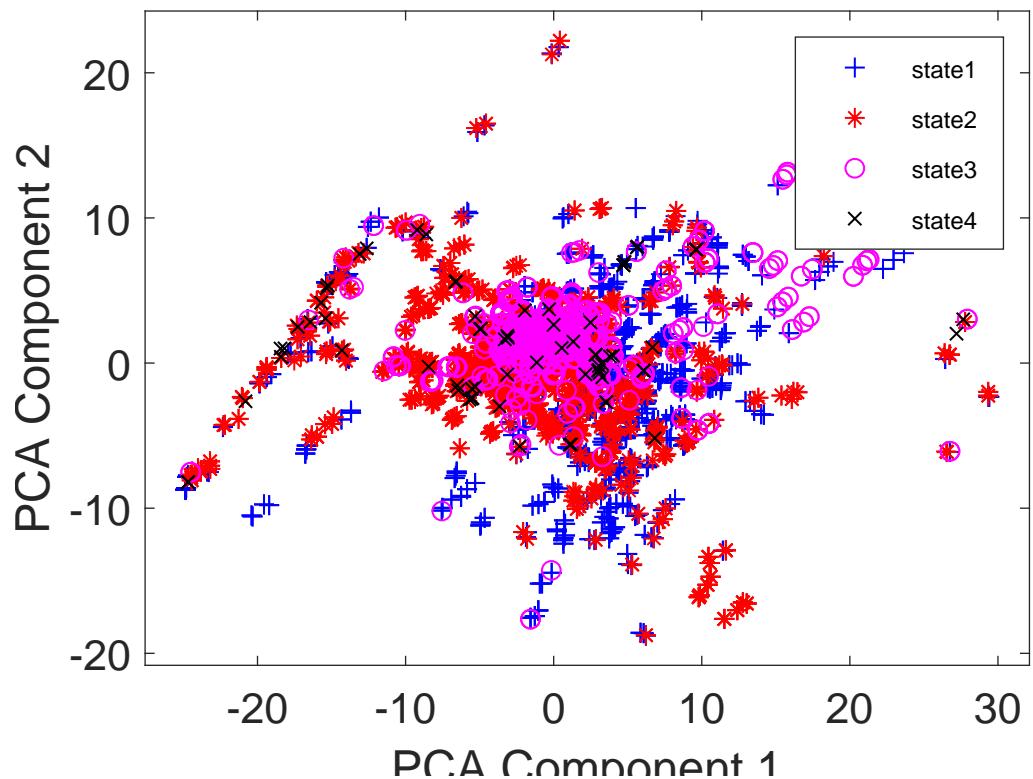


Figure A.24: Scatter plot of 13-dimensional observations grouped according to the estimated state sequence

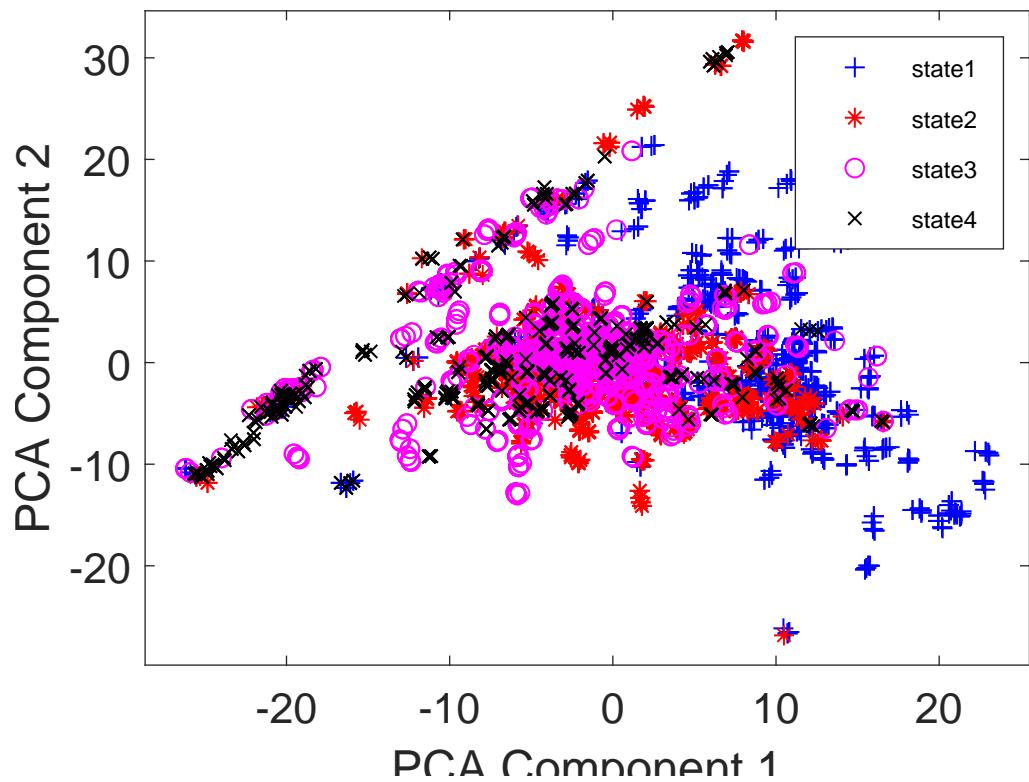


Figure A.25: Scatter plot of 14-dimensional observations grouped according to the ground-truth state sequence

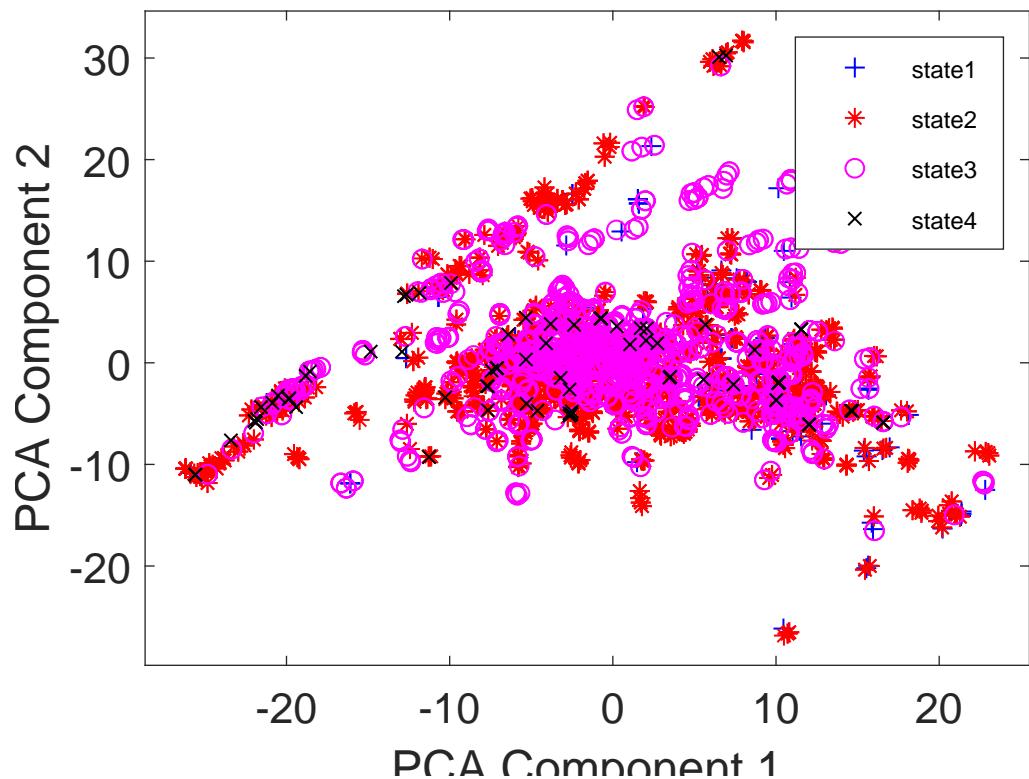


Figure A.26: Scatter plot of 14-dimensional observations grouped according to the estimated state sequence

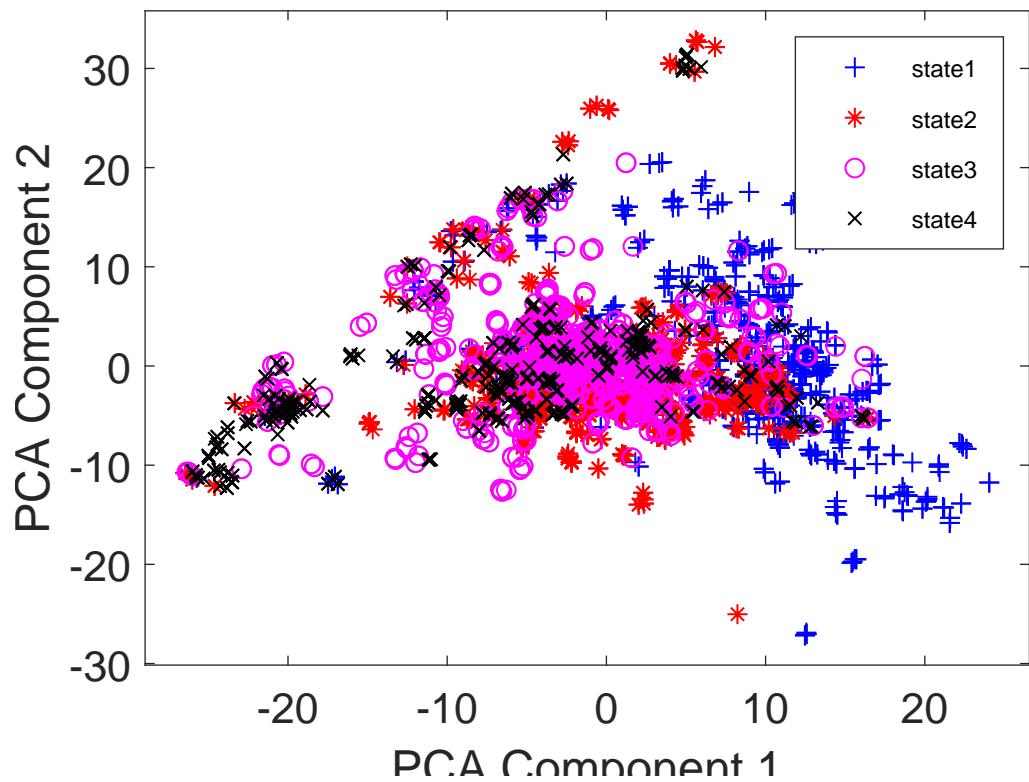


Figure A.27: Scatter plot of 15-dimensional observations grouped according to the ground-truth state sequence

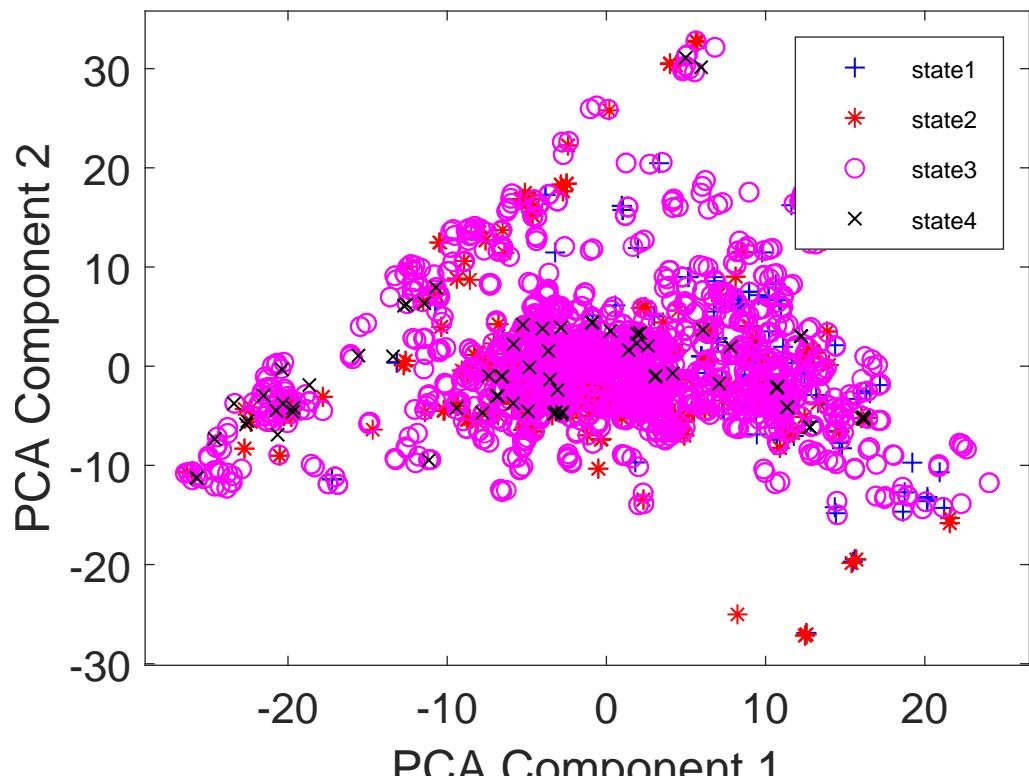


Figure A.28: Scatter plot of 15-dimensional observations grouped according to the estimated state sequence

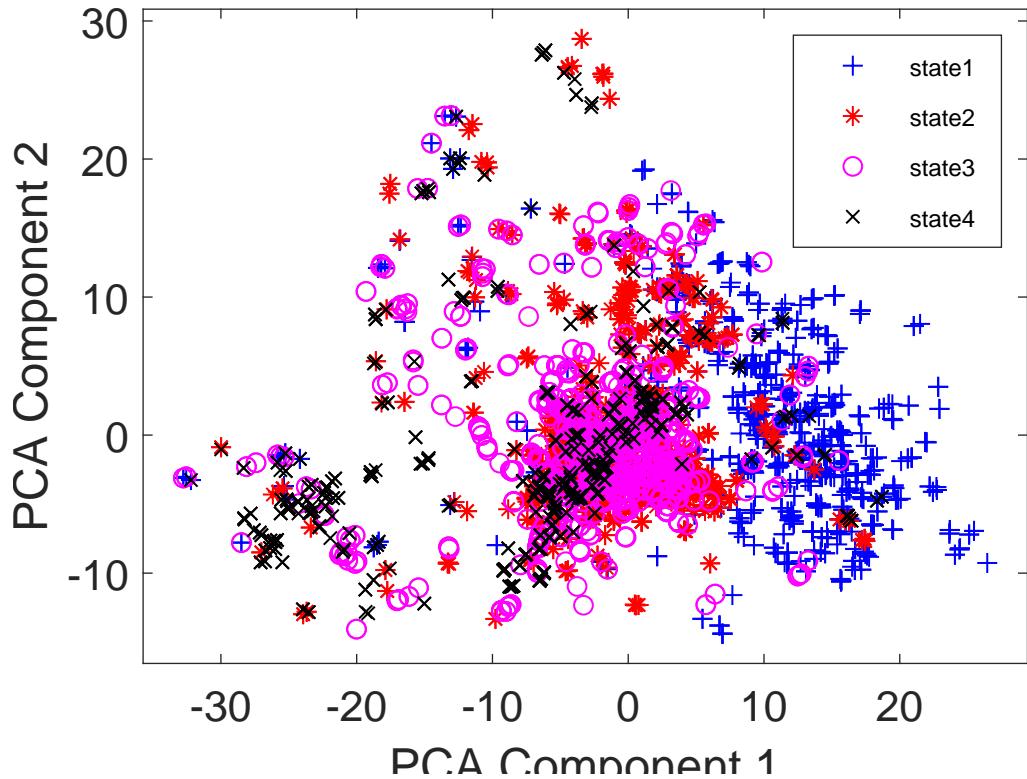


Figure A.29: Scatter plot of 16-dimensional observations grouped according to the ground-truth state sequence

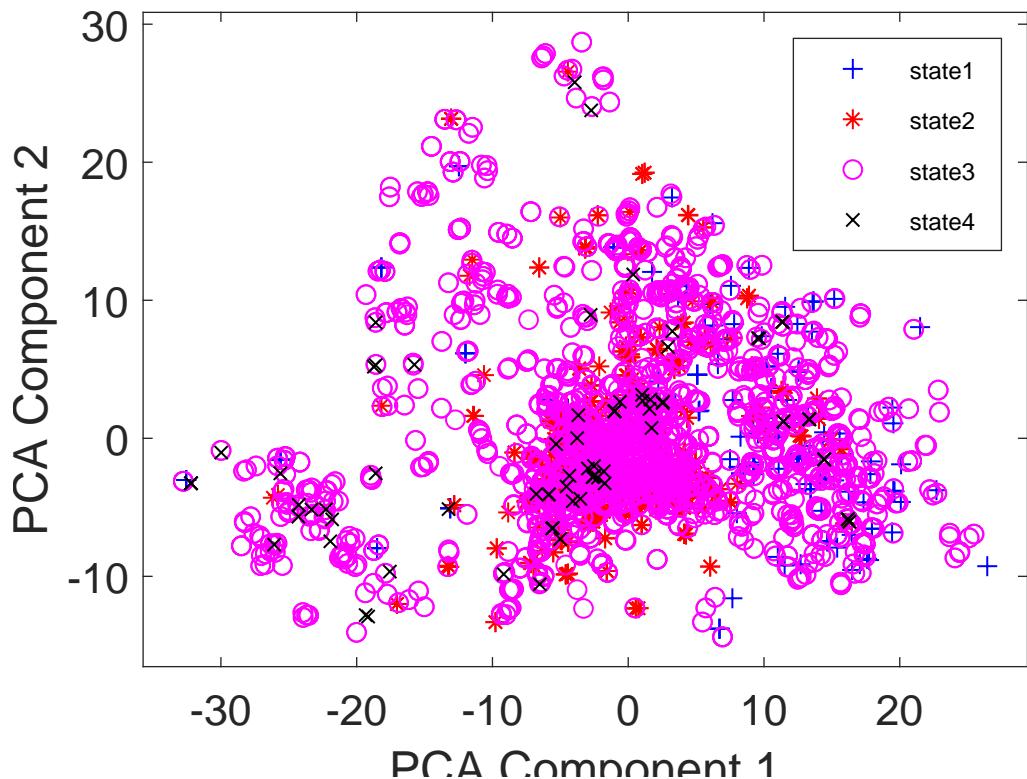


Figure A.30: Scatter plot of 16-dimensional observations grouped according to the estimated state sequence

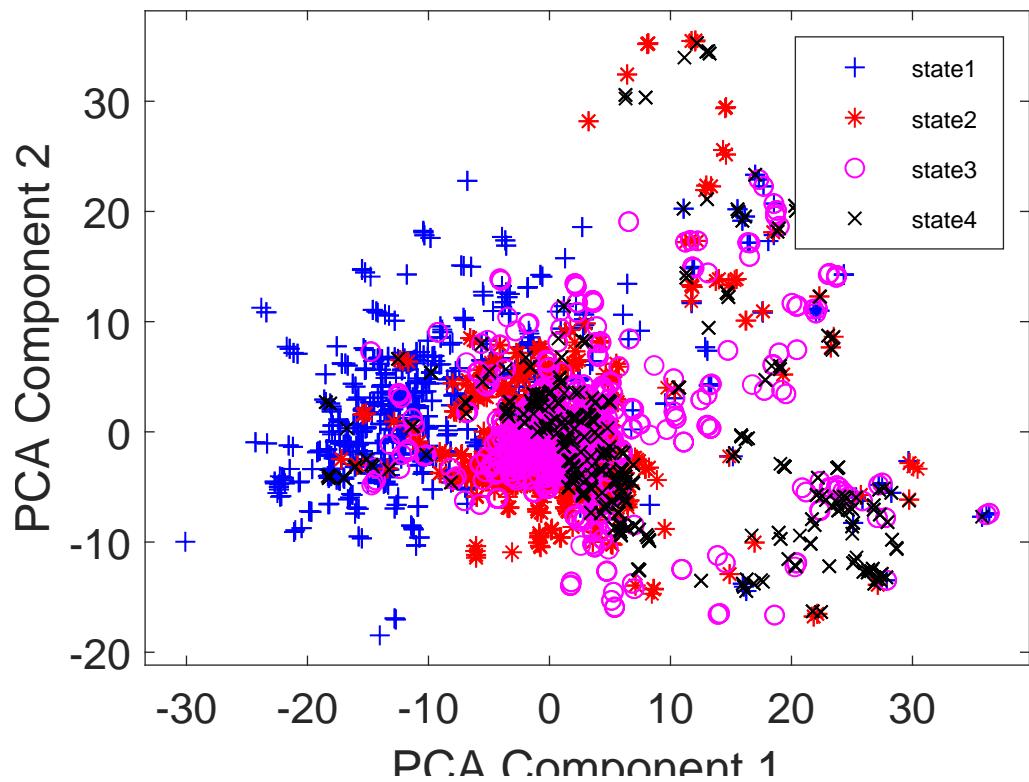


Figure A.31: Scatter plot of 17-dimensional observations grouped according to the ground-truth state sequence

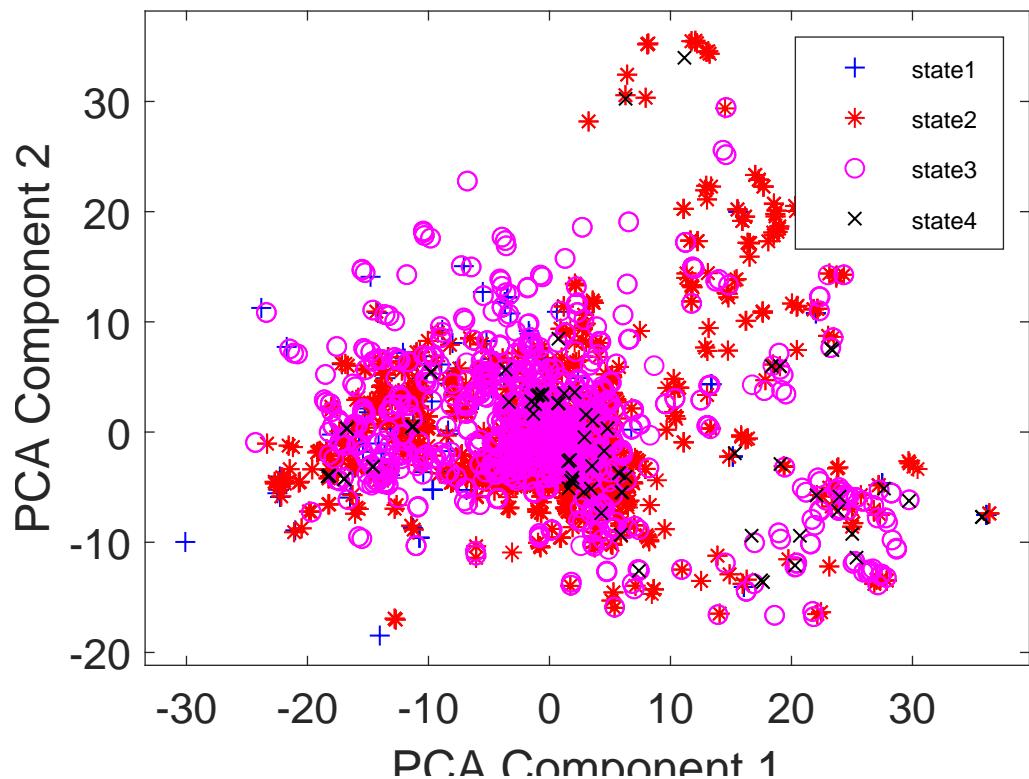


Figure A.32: Scatter plot of 17-dimensional observations grouped according to the estimated state sequence

Appendix B

Addenda

B.1 Ethics Forms

B.1. ETHICS FORMS

Application for Approval of Ethics in Research (EiR) Projects
Faculty of Engineering and the Built Environment, University of Cape Town

APPLICATION FORM

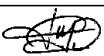
Please Note:

Any person planning to undertake research in the Faculty of Engineering and the Built Environment (EBE) at the University of Cape Town is required to complete this form **before** collecting or analysing data. The objective of submitting this application *prior* to embarking on research is to ensure that the highest ethical standards in research, conducted under the auspices of the EBE Faculty, are met. Please ensure that you have read, and understood the **EBE Ethics in Research Handbook**(available from the UCT EBE, Research Ethics website) prior to completing this application form: <http://www.ebe.uct.ac.za/usr/ebe/research/ethics.pdf>

APPLICANT'S DETAILS		
Name of principal researcher, student or external applicant		Kouame Hermann Kouassi
Department		Electrical Engineering
Preferred email address of applicant:		Ksskou001@myuct.ac.za
If a Student	Your Degree: e.g., MSc, PhD, etc.,	Bsc
	Name of Supervisor (if supervised):	Fred Nicolls
If this is a research contract, indicate the source of funding/sponsorship		
Project Title		Gait sequence modelling and estimation

I hereby undertake to carry out my research in such a way that:

- there is no apparent legal objection to the nature or the method of research; and
- the research will not compromise staff or students or the other responsibilities of the University;
- the stated objective will be achieved, and the findings will have a high degree of validity;
- limitations and alternative interpretations will be considered;
- the findings could be subject to peer review and publicly available; and
- I will comply with the conventions of copyright and avoid any practice that would constitute plagiarism.

SIGNED BY	Full name	Signature	Date
Principal Researcher/ Student/External applicant	Kouame Kouassi		23 Aug 2017

APPLICATION APPROVED BY	Full name	Signature	Date
Supervisor (where applicable)	Fred Nicolls		25/08/2017
HOD (or delegated nominee) Final authority for all applicants who have answered NO to all questions in Section1; and for all Undergraduate research (Including Honours).			
Chair : Faculty EIR Committee For applicants other than undergraduate students who have answered YES to any of the above questions.			