

Sign Align – A System to Track and Translate Sign Language

Daniel Nichol

Department of Computer Science

University of Oxford

A thesis submitted for the degree of

Master of Mathematics and Computer Science

Acknowledgements

Abstract

This is where you write your abstract ...

Contents

Contents	iv
List of Figures	vi
Nomenclature	vi
1 Introduction	1
1.1 Sign Language	1
1.2 Hardware	3
1.3 Outline	3
2 Background and Model Selection	4
2.1 Model Selection	4
2.2 Hidden Markov Models	4
2.2.1 The Forward and Backward Variables	6
2.2.2 Forward-Backward Algorithm	8
2.2.3 Limitations of HMMs	10
2.3 Previous Work	11
3 Implementing a Classifier for Isolated Signs	12
3.1 Implementing a Hidden Markov Model Class	12
3.1.1 Scaling	12
3.1.2 Training From Multiple Observation Sequences	14
3.1.3 Observation Vector Quantization	15
3.2 Implementing the signModel class	17
3.2.1 Determining The Initial Topology and Parameters	19

3.2.2	Determining the Weightings	22
3.3	Building a Sign Classifier	22
4	Interfacing With The Kinect Sensor	24
4.1	Recording Training Data	24
4.1.1	The Training Data	26
4.1.2	Limitations in the Training Data	26
5	Testing and Experimentation	28
5.1	Determining the Classifier Parameters	28
5.1.1	The Cluster Number	28
5.1.2	The Acceptance Threshold	28
5.2	Sign Recognition Accuracy	29
6	Analysis and Conclusions	30
7	Additional Work	31
	Appdx A	33
	Appdx B	34
	References	35

List of Figures

2.1	A Hidden Markov Model. The underlying Markov Chain contains 4 states (depicted as circles) which emit two possible observations (a or b).	5
3.1	A UML skeleton of the DHMM class	16
3.2	The K-Means method for vector quantization	17
3.3	A UML diagram for the signModel class	19
3.4	Different Markov Chain Topologies	20
3.5	The signClassifier in UML	22
4.1	A UML diagram of the gesture recording classes	25
4.2	An example of the broken hand tracking with the hands placed together. The yellow joints are inferred by the sensor incorrectly, the actual location of the hands is just below the chin.	27

Chapter 1

Introduction

British Sign Language is the preferred first language of over 125,000 in the United Kingdom, in addition to an estimated 20,000 children and thousands of hearing friends, relatives and interpreters. There are many more people worldwide communicating in an estimated 200 different signed languages, with a huge variation in grammar and pronunciation. As such, a system which can recognise signed languages and convert them into text in real time would be a valuable tool for many people.

The aim of this project is to implement a system to recognise gestures of a signed language from a continuous data stream provided from the Microsoft Kinect camera. The problem of continuous gesture recognition has been studied for over 20 years and has seen solutions which often involve specific gloves or expensive hardware and which have been particularly sensitive to lighting conditions and camera placement. Using the Kinect sensor, which is available to purchase off the shelf and is designed to be used without gloves and in a range of conditions, we aim to build a robust sign recognition system which does not suffer from these limitations.

1.1 Sign Language

Signed languages are natural languages which have evolved independently of the spoken languages of the areas in which they are used and often independently of

one another. For example, British Sign Language (BSL) is not simply oral English transcribed but rather a distinct language with its own sentence structure which differs significantly from English. Further despite the regions sharing a common language American Sign Language (ASL) is a separate language from BSL.

Despite the large variation between regional sign languages the means of communication remain the same. The main component of the sign is the movement of the body and each sign is determined by the movement and location of the hands and arms as well as the hand shape and palm direction. In addition to the manual component of the signs the signer will often use posture, facial expressions, mouth shape or eye movements to convey meaning.

The non-manual components of sign language will be ignored for the purposes of this project. It should be noted that eye tracking and facial expression are essentially independent problems from that of gesture recognition and hence if solutions to the facial expression problem exist they may be combined with the work of this project to produce a more substantial sign recognition system.

We can further break the manual part of a sign in to two components. The first is the movement of the hand, arms and body over time and the second is the orientation of the hand throughout this movement. This is a sensible distinction to make as it helps reduce the number of distinct signs to differentiate between. As the same hand shape might be used with two different arm movements to produce different words, we can reduce the number of distinct gestures we wish to detect by detecting hand movement gestures and hand shapes and then combining the two to identify a sign.

In the next section we will see there is a limitation in our choice of hardware that makes hand shape detection a problem which is outside the scope of this project. However the above observation suggests that the work in this project to detect arm and body gestures can be combined with future work to produce a system which will detect a signed language. Furthermore, a major task of this project is to implement a library of Hidden Markov Model sequence classifiers to detect gestures and this library can be used to detect hand shape gestures just as we use it here to detect body gestures.

1.2 Hardware

The Kinect is a motion sensing camera designed by Microsoft and released in November 2010. The Kinect combines an RGB camera and an infrared depth sensor to detect objects in 3D space. The raw images of these sensors are combined using proprietary software to detect a human body as a skeleton, in particular the Kinect is capable of detecting the positions of the hands, elbows, shoulders and head as a point in 3D space.

Microsoft released the Kinect software development kit (SDK) for public use on January 16, 2011. This SDK allows developers build applications which interact with the proprietary skeleton tracking software using C#, C++ or Visual Basic. In this project we will use this SDK to build the gesture recognition framework.

The Kinect camera was originally designed with the ability to detect hand and finger positions. In fact the original patent claims the device would be able to detect American Sign Language ([LATTA et al., 2010](#)), however this functionality was removed from the original release of the Kinect sensor. The SDK was updated to recognise open and closed hands on March 18th, 2013 ([Microsoft, 2013b](#)) and it is believed that the next iteration of the Kinect hardware will be able to fully detect hand gestures.

Third party hand gesture recognition frameworks do exist, however the most successful of these do not make use of the Kinect SDK and so do not make use of its features at all ([I. Oikonomidis, 2013](#)). We aim to build a framework that will be easily extended when the capabilities of the Kinect are improved and for this reason we choose to only use the tools provided within the Kinect SDK.

1.3 Outline

Chapter 2

Background and Model Selection

2.1 Model Selection

2.2 Hidden Markov Models

A *hidden Markov Model* (HMM) is (following the definitions in [Rabiner \(1989\)](#)) a doubly stochastic process which consists of an underlying discrete Markov chain with state set $S = \{S_1, \dots, S_n\}$ and stochastic matrix $A = [a_{i,j}]_{N \times N}$ where

$$a_{i,j} = \mathbb{P}(q_{t+1} = S_i | q_t = S_j) \text{ for each } 1 \leq i, j \leq N$$

which is hidden from an observer in the respect that one cannot directly observe the current state of the Markov chain. Instead we have a collection of M observable symbols, say $V = \{v_1, \dots, v_M\}$, which may be observed with probability dependent on the state underlying Markov chain (Figure: [2.1](#)).

We encode these so-called *emission probabilities* in a matrix $B = [b_j(k)]_{N \times M}$ where

$$b_j(k) = \mathbb{P}[v_k \text{ occurs at time } t | q_t = S_j]$$

Now if we take an initial state distribution $\boldsymbol{\pi} = [\pi_1, \dots, \pi_N]$ for our Markov chain with

$$\pi_i = \mathbb{P}[q_1 = S_i]$$

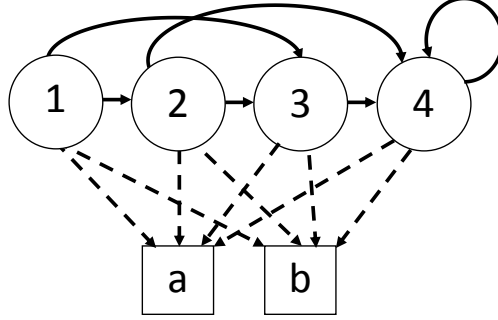


Figure 2.1: A Hidden Markov Model. The underlying Markov Chain contains 4 states (depicted as circles) which emit two possible observations (a or b).

Then the HMM generates a sequence of observations

$$\mathbf{O} = O_1, O_2, \dots, O_T$$

by the following process:

1. Choose an initial state $q_1 = S_i$ stochastically from the initial state distribution $\boldsymbol{\pi}$
2. For $t = 1$ to T
 - i. Choose $O_t = v_k$ according to the distribution $b_i(k)$ of the current state S_i
 - ii. Stochastically transition to a new state S_j from S_i according to A

Note now that a hidden Markov Model is entirely determined by the transition matrix A , the emissions matrix B and the initial distribution $\boldsymbol{\pi}$ (noting that the dimensions N and M are encoded in the dimensions of the matrices) hence for convenience we may denote a HMM by

$$\lambda = (A, B, \boldsymbol{\pi})$$

Hidden Markov Models were first introduced by in a series of papers by L.E Baum and others in the late 1960s (Baum and Petrie, 1966; Baum et al., 1970) and have been since used to study handwriting recognition (Bunke et al., 1995), speech recognition (Jelinek, 1998; Juang and Rabiner, 1991), natural language modelling (Jurafsky et al., 2002; Manning and Schütze, 1999) and biological processes (Durbin et al., 1998; Krogh et al., 1994; Liò and Goldman, 1998).

Given these definitions there exist three basic problems which form the basis of using HMMs as a tool for machine learning

1. Given an observation sequence $\mathbf{O} = O_1, O_2, \dots, O_T$ and a HMM λ , what is $\mathbb{P}(\mathbf{O}|\lambda)$ - the probability that the observation sequence \mathbf{O} was produced by λ ?
2. Given a HMM λ and an observation sequence $\mathbf{O} = O_1, O_2, \dots, O_T$ what is the state sequence of the underlying Markov chain in λ which is most likely to have generated \mathbf{O} ?
3. Given an observation sequence \mathbf{O} , how to do we choose the parameters of $\lambda = (A, B, \boldsymbol{\pi})$ which best optimise $\mathbb{P}(\mathbf{O}|\lambda)$?

In fact the problem of sign language gesture recognition can be seen as instance of problems 3 and 1. First we take for each gesture g a set of training data which are recordings of the joints in 3D space. This training data forms a collection of observation sequences $\mathbf{O}_1, \dots, \mathbf{O}_n$ which are used in a solution to problem 3 to parameterise a hidden Markov Model λ_g to model the gesture.

Then given a fresh observation sequence \mathbf{O} we can use a solution to problem 1 to compute $\mathbb{P}[\mathbf{O}|\lambda_g]$ for each g . We can then take the gesture g for which this probability is maximised and, provided the probability exceeds some threshold, conclude that the gesture g corresponds to the observation sequence \mathbf{O} .

2.2.1 The Forward and Backward Variables

Suppose we are given an observation sequence $\mathbf{O} = O_1, O_2, \dots, O_T$ and a HMM λ and we wish to solve the evaluation problem (problem 1). Rabiner (1989)

notes that a direct computation of $\mathbb{P}[\mathbf{O}|\lambda]$ will take time order $\mathcal{O}(2TN^T)$ which is infeasible even for moderate values of N and T . Instead we use a dynamic programming approach, the forward algorithm, introduced in [Baum and Sell \(1968\)](#) and [Baum et al. \(1970\)](#).

Define the *forward variables* $\alpha_t(i)$ for each $1 \leq t \leq T$ and $1 \leq i \leq N$ by

$$\alpha_t(i) = \mathbb{P}[O_1, \dots, O_t, q_t = S_i | \lambda]$$

Then note these can be inductively computed by the following procedure

$$\begin{aligned} \alpha_1(i) &= \pi_i b_i(O_1) \\ \alpha_{t+1}(j) &= \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1}) \quad \text{for } 1 \leq t \leq T-1 \text{ and } 1 \leq j \leq N \end{aligned}$$

Finally we have

$$\mathbb{P}[\mathbf{O}|\lambda] = \sum_{i=1}^N \mathbb{P}[\mathbf{O}, q_T = S_i | \lambda] = \sum_{i=1}^N \alpha_T(i)$$

The set of forward variables can be computed by dynamic programming in $\mathcal{O}(N^2T)$ time and hence we can compute $\mathbb{P}[\mathbf{O}|\lambda]$ in this time - a significant improvement over direct computation. This algorithm is the *forward algorithm* for evaluation in HMMs.

Symmetrically to the forward variables we can define a collection of *backward variables* by

$$\beta_t(i) = \mathbf{P}[O_{t+1}, O_{t+2}, \dots, O_T | q_t = S_i, \lambda]$$

Which gives the probability of a partial observation sequence from a time t given the state of the HMM λ at time t . These too can be computed iteratively as

$$\begin{aligned} \beta_T(i) &= 1 \quad \text{for each } 1 \leq i \leq N \\ \beta_t(i) &= \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \beta_{t+1}(j) \quad \text{for } t = T-1, T-2, \dots, 1 \text{ and } 1 \leq i \leq N \end{aligned}$$

and these too can be computed by dynamic programming in $\mathcal{O}(N^2T)$ time. The backward variables are not needed in the solution to the evaluation problem but are used in the following section to re-parameterise hidden Markov Models.

2.2.2 Forward-Backward Algorithm

The problem of parameterising a hidden Markov Model λ is considerably more difficult than the other two problems of HMMs. In fact it is known that there is no analytic solution which provides a model λ to maximise the probability of some observation sequence \mathbf{O} . Instead we must use local optimisation methods which given an initial parameterisation λ_0 iteratively improve it until $\mathbb{P}[\mathbf{O}|\lambda]$ reaches a local maximum.

We will use a modified version of the *Baum-Welch algorithm* introduced by [Baum et al. \(1970\)](#) called the *forward-backward algorithm* ([Rabiner, 1989](#)) which is reproduced below. This algorithm is an instance of an expectation-maximization algorithm [Moon \(1996\)](#) which is a general technique used to determine maximum likelihood estimators in a number of machine learning models [Bishop et al. \(2006\)](#).

For $t \in \{1, \dots, T-1\}$ and $i, j \in \{1, \dots, N\}$ define

$$\gamma_t(i, j) = \mathbb{P}[q_t = S_i, q_{t+1} = S_j | \mathbf{O}, \lambda]$$

such that $\gamma_t(i, j)$ is the probability of being in state S_i at time t and transitioning to S_j at the next time step given the HMM λ and the observation sequence \mathbf{O} . The by definition of the forward and backward variables we have

$$\gamma_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta(j)}{\mathbb{P}[\mathbf{O}|\lambda]}$$

We can define for each $1 \leq t \leq T-1$ and each $1 \leq i \leq N$ the probability of being in state i at time t by

$$\gamma_t(i) = \mathbb{P}[q_t = S_i | \mathbf{O}, \lambda] = \frac{\alpha_t(i) \beta_t(i)}{\mathbb{P}[\mathbf{O}|\lambda]}$$

and then we have

$$\gamma_t(i) = \sum_{j=1}^N \gamma_t(i, j)$$

Now using these equations we can compute

$$\begin{aligned} \sum_{t=1}^{T-1} \gamma_t(i) &= \text{The expected number of transitions from } S_i \\ \sum_{t=1}^{T-1} \gamma_t(i, j) &= \text{The expected number of transitions from } S_i \text{ to } S_j \end{aligned}$$

which can be used to reparameterise the model $\lambda = (A, B, \boldsymbol{\pi})$ as follows, set

$$\begin{aligned} \bar{\pi} &= \text{the expected number of times in state } S_i \text{ at time 1} = \gamma_1(i) \\ \bar{a}_{ij} &= \frac{\text{the expected number of transitions from } S_i \text{ to } S_j}{\text{expected number of transitions from state } S_i} \\ &= \frac{\sum_{t=1}^{T-1} \gamma_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \\ \bar{b}_j(k) &= \frac{\text{the expected number of times in state } S_j \text{ observing symbol } v_k}{\text{the expected number of times in state } S_j} \\ &= \frac{\sum_{t=1}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)} \end{aligned}$$

Then if denote $\bar{\lambda} = (\bar{A}, \bar{B}, \bar{\boldsymbol{\pi}})$ then it has been proven ([Baum and Sell, 1968](#); [Levinson et al., 1983](#)) that either

1. $\mathbb{P}[\mathbf{O}|\bar{\lambda}] > \mathbb{P}[\mathbf{O}|\lambda]$ or
2. λ is locally maximized with respect to $\mathbb{P}[\mathbf{O}|\lambda]$ and $\bar{\lambda} = \lambda$

It follows that given an initial HMM λ we may improve it to a locally optimal model for some observation sequence \mathbf{O} via the following local search:

1. Whilst $\mathbb{P}[\mathbf{O}|\lambda]$ increases:
 - i. Compute the $\alpha_t(i)$, $\beta_t(i)$, $\gamma_t(i, j)$ and $\gamma_t(i)$

-
- ii. Compute the new parameters $\bar{A} = [\bar{a}_{ij}]$, $\bar{B} = [\bar{b}_j(k)]$ and $\bar{\pi} = [\bar{\pi}_1, \dots, \bar{\pi}_N]$
 - iii. Set $\lambda := \bar{\lambda}$

2. return λ

Note that this algorithm may not actually terminate. In practice we threshold the increase of $\mathbb{P}[\mathbf{O}|\lambda]$ to ensure termination in a reasonable time. Further, in practice we will not have single observation sequence but rather a collection, perhaps from a variety of signers of different heights, genders, ages or dialects. In the implementation of our HMM framework we will modify this algorithm to work for multiple observation sequences. This will again increase the time complexity of the algorithm, however this problem is not so significant as in practice we will use this algorithm once per sign and save the parameterisations.

2.2.3 Limitations of HMMs

A significant limitation of hidden Markov Models is that they can have only a finite set V of possible observations. In practice this can prevent us using a HMM to model a specific gesture exactly. Take for example the problem of detecting the gesture of a circle drawn on a 2D plane. We might create a hidden Markov Model in which the states of the Markov chain represent some specific points on the plane and want our matrix B to be such that

$$b_j(\mathbf{p}) = \mathbb{P}[\text{the pen is at point } \mathbf{p} \text{ at time } t \mid q_t = S_j] \quad \text{for each point } \mathbf{p} \in \mathbb{R}^2$$

However as the plane is continuous and V is finite this is not possible. One solution is to discretise the plane and have only a finite (but large) set of observation symbols. This method has been used with some success to distinguish between gestures with large variations, for example different tennis strokes (Yamato et al., 1992). It could be argued that this discretisation of the observation space could cause us to lose some subtlety in the input and hence prevent us from distinguishing between similar signs. We attempt to compensate for this by taking into account the full torso in determining a sign and using the elbows, shoulders and head to help distinguish signs where the hands are not sufficient alone.

A second solution is to utilise *Continuous Distribution Hidden Markov Models* which extend HMMs by replacing the emissions matrix B with a probability density function. This method proved to be outside of the scope of this project but is one that has proven fruitful in previous pattern detection systems. As such we have built our sign recognition system to allow the discrete hidden markov models to be replaced by a continuous counterpart provided it implements a specific interface.

2.3 Previous Work

Motivated by the success of hidden Markov Models for speech in the mid 1970s (Baker, 1975; Jelinek et al., 1975) these models were adopted for gesture detection. Yamato et al. (1992) used discrete HMMs with observations from a 2D camera to detect tennis strokes and later Starner and Pentland (1995) implemented a system to recognise American Sign Language in real-time using by continuous density hidden Markov Models. This system relied on a single camera and required the user to wear a pair of special gloves. By imposing a restricted grammar to only allow sentences of a particular structure they were able to achieve an accuracy of 97.0% on an independent training set.

More recently the SignSpeak project (Dreuw et al., 2010), which aims to use a 2D video camera and an extension of HMM techniques (Dreuw et al., 2009), has received EU funding. This project aims to solve the problems of image extraction, sign recognition and language translation to provide a complete video-to-text system for sign language translation.

Chapter 3

Implementing a Classifier for Isolated Signs

3.1 Implementing a Hidden Markov Model Class

Our first task in implementing a classifier for signed languages was to implement a general purpose class `DHMM` (figure: 3.1) which can be instantiated to model a Hidden Markov Model $\lambda = (A, B, \pi)$ by providing appropriate parameters. In this class we implemented solutions to the evaluation and training problems and provided methods for saving and loading the trained parameters. Extending the methods described in the previous section we tailored this implementation for use learning from a collection of data recorded from the Kinect Sensor.

3.1.1 Scaling

In the previous section we introduced methods for solving the evaluation and re-estimation problems for hidden markov models. These methods, whilst mathematically sound, introduced certain problems during implementation. In particular these solutions involve the products of a large number of probabilities, especially for long observation sequence and as a consequence when implemented in C# (which lacks arbitrary precision floats) can cause errors due to underflow. This problem is a common one in the implementation of Hidden Markov Models and to solve it we used the methods used by Rabiner (1989) (and corrected

by [Rahimi \(2000\)](#)) and scale the intermediate α , β and γ variables to prevent underflow. To do this we normalize the $\alpha_t(i)$ by setting

$$\bar{\alpha}_0(i) = \alpha_0(i) \text{ for each } 0 \leq i < N$$

and inductively define for each $0 \leq t < T$

$$\begin{aligned} c_t &= \frac{1}{\sum_{i=0}^{N-1} \bar{\alpha}_t(i)} \\ \hat{\alpha}_t(i) &= c_t \bar{\alpha}_t(i) \text{ for each } i \\ \bar{\alpha}_{t+1}(i) &= \sum_{j=0}^{N-1} \hat{\alpha}_t(i) a_{ij} b_j(O_{t+1}) \text{ for each } i \end{aligned}$$

Which prevents underflow in the intermediate calculation of the α variables. Further we then have ([Rabiner, 1989](#)) that

$$\begin{aligned} 1 &= \sum_{i=0}^N \hat{\alpha}_{T-1}(i) = c_0 c_1 \dots c_{T-1} \sum_{i=0}^N \alpha_{T-i}(j) \\ &= c_0 c_1 \dots c_{T-1} \mathbb{P}(\mathbf{O}|\lambda) \end{aligned}$$

and hence we can compute the log-probability of an observation sequence \mathbf{O} by

$$\log(\mathbb{P}[\mathbf{O}|\lambda]) = - \sum_{t=0}^{T-1} \log(c_t)$$

which allows us to determine probabilities which might have otherwise been lost due to underflow. We then use the same scale factors on the β variables by

inductively defining

$$\begin{aligned}\bar{\beta}_{T-1}(i) &= \beta_{T-1}(i) \\ \hat{\beta}_t(i) &= c_t \bar{\beta}_t(i) \\ \bar{\beta}_{t+1}(i) &= \sum_{j=0}^{N-1} a_{ij} b_j(\mathbf{O}_{t+1}) \hat{\beta}_{t+1}(i)\end{aligned}$$

and redefine the γ variables as

$$\begin{aligned}\gamma_t(i, j) &= \hat{\alpha}(i) a_{ij} b_j(\mathbf{O}_{t+1}) \beta_{t+1}(j) \\ \gamma_t(i) &= \hat{\alpha}_t(i) \hat{\beta}_t(i) \frac{1}{c_t}\end{aligned}$$

Then if we use these variables in the re-estimation it is still the case that the probability $\mathbb{P}[\mathbf{O}|\lambda]$ increases with each iteration (Rabiner, 1989) and hence we can perform the iterative reestimation procedure for HMMs without introducing underflow errors.

3.1.2 Training From Multiple Observation Sequences

A second restriction of the training procedure presented in chapter 2 is that it restricts the training data of the HMM to be a single observation sequence. As we are going to train our HMMs on data which is derived from the Kinect Sensor data of a person performing a sign this is too great of a restriction. For example, the way in which a person performs a sign will be unique each time and will certainly differ between signers and so to train a recogniser on a single observation sequence of a sign would not fully capture the range of movements which might be interpreted as a sign. Fortunately the reestimation procedure can be extended (Li et al., 2000; Rabiner, 1989; van Oosten, 2010) to allow training from a set of multiple observation sequences $\mathcal{O} = \{\mathbf{O}_0, \dots, \mathbf{O}_{K-1}\}$ to maximise

$$\mathbb{P}[\mathcal{O}|\lambda] = \prod_{i=0}^{K-1} \mathbb{P}[\mathbf{O}_i|\lambda]$$

by computing for each $0 \leq i < K$ the set of scaled $\hat{\alpha}$, $\hat{\beta}$ and $\hat{\gamma}$ variables associated with the observation sequence \mathbf{O}_i and the associated scales. Using these scaled variables we can reestimate the parameters of our HMM as

$$\begin{aligned}\bar{\pi}_i &= \frac{\sum_{k=0}^{K-1} \gamma_0^{(k)}(i)}{\sum_{j=0}^{N-1} \sum_{k=0}^{K-1} \gamma_0^{(k)}(j)} \\ \bar{a}_{ij} &= \frac{\sum_{k=0}^{K-1} \sum_{t=0}^{T_k-2} \hat{\alpha}_t^{(k)}(i) a_{ij} b_j(O_{t+1}^{(k)}) \hat{\beta}_{t+1}^{(k)}(j)}{\sum_{k=0}^{K-1} \sum_{t=0}^{T_k-2} \hat{\alpha}_t^{(k)}(i) \hat{\beta}_t^{(k)}(i) \frac{1}{c_t^{(k)}}} \\ \bar{b}_j(l) &= \frac{\sum_{k=0}^{K-1} \sum_{t \in \{0, \dots, T_k-2\} \text{ and } O_t = v_l} \hat{\alpha}_t^{(k)}(i) \hat{\beta}_t^{(k)}(i) \frac{1}{c_t^{(k)}}}{\sum_{k=0}^{K-1} \sum_{t=0}^{T_k-1} \hat{\alpha}_t^{(k)}(i) \hat{\beta}_t^{(k)}(i) \frac{1}{c_t^{(k)}}}\end{aligned}$$

and this will ensure that $\mathbb{P}[\mathcal{O}|\bar{\lambda}] \geq \mathbb{P}[\mathcal{O}|\lambda]$. Using this procedure we implemented the `Reestimate` method of our `DHMM` class.

3.1.3 Observation Vector Quantization

In order to use discrete observation hidden markov models for the problem of gesture recognition we must restrict the continuous observation space to a discrete set of observation symbols. This set of symbols should not be too large (say no more than 30 symbols) to ensure that the problems of training and evaluation can be solved quickly enough for use in a real-time recognition system. Further, an increase in the number of symbols will cause a decrease in the evaluated probability of a given observation sequence being generated by some HMM. If the symbol count is too high this can cause underflow in spite of the effects of scaling. Of course the symbol set should not be too small either, as otherwise some subtle aspects of signs will be lost similar signs will be indistinguishable.

To quantize the training data for a single discrete HMM instance we merged all of the points of all of the training data into a single 3D point cloud and partitioned this point cloud into k clusters. To do this we implemented Lloyd's

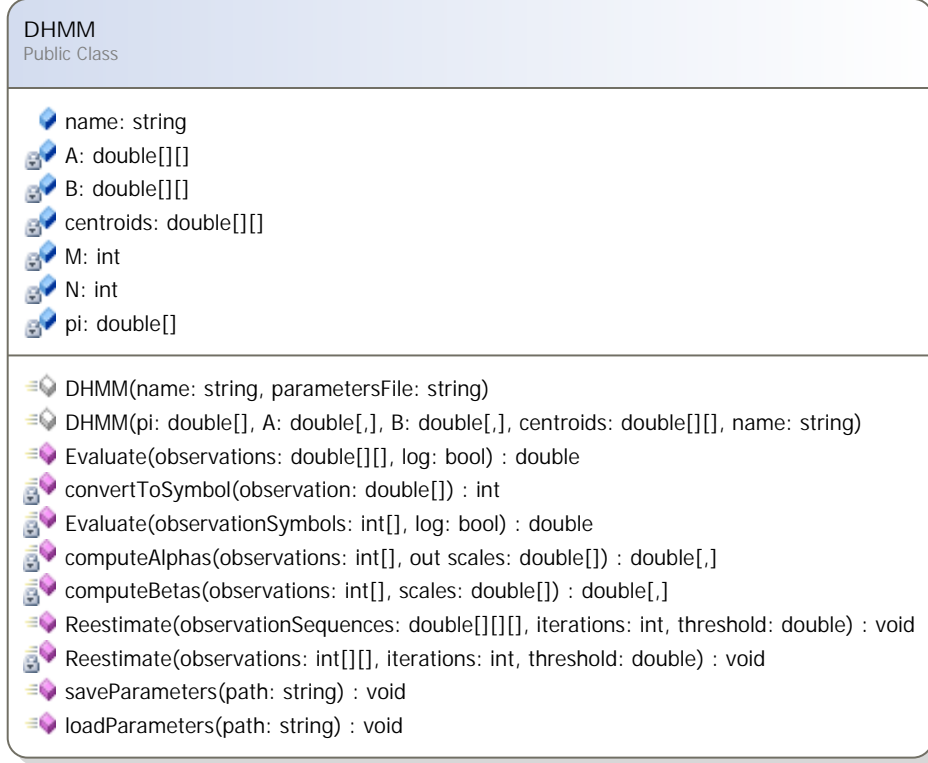


Figure 3.1: A UML skeleton of the DHMM class

algorithm (Lloyd, 1982) to solve the k-means clustering problem (Figure: 3.2). Then we defined a translation from the continuous \mathbb{R}^3 observation space provided by the Kinect Sensor to a finite set of $k+1$ observations by assigning to each point the number of the cluster with the centroid closest (via Euclidean distance) to it, or if the distance from the point to each of the means exceeds some threshold assigning the symbol $k+1$. In fact, this quantization method will work just as well to translate arbitrary dimension continuous spaces to a finite sequence set and hence offers us the opportunity to extend our training data sets to include not only spatial positions but other information as well, say velocity or positions relative to some fixed location.

Using this vector quantization method we were able to extend the **Evaluate** and **Reestimate** method of DHMM to take as input streams of (x, y, z) coordinates captured by the Kinect Sensor.

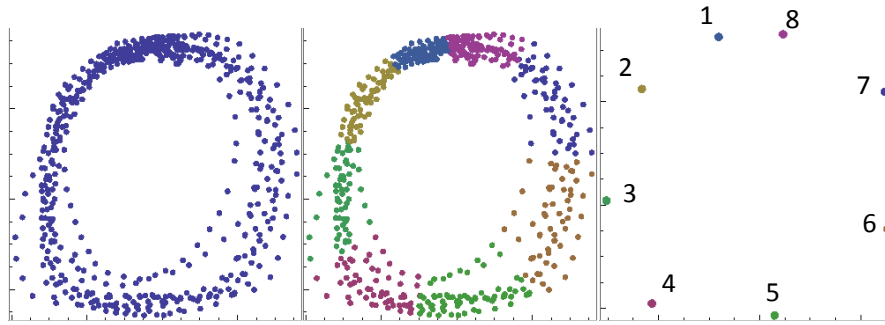


Figure 3.2: The K-Means method for vector quantization

3.2 Implementing the signModel class

The Kinect Sensor provides a stream of real-time positions, as vectors in \mathbb{R}^3 , of a number of joints in the body. A single `DHMM` object can be trained on a collection of observations sequences of one of these joints, the left hand say, and be used to determine the probability that another observation sequence was produced by that HMM. If this probability is high enough we might conclude that this observations sequence matches the sign that produced the training data. A collection of such trained `DHMM` used to distinguish between signs would constitute a sign classifier capable of distinguishing between broadly different gestures of the left hand. However such a classifier will only use observations from the left hand and ignore the information provided by the rest of the body which also gives information about each sign.

One solution to this problem is to create an instance `DHMM` which takes as its input streams vectors of $(\mathbb{R}^3)^J$ where J is the number of joints we choose to track. This will allow us to train the model on full observations of the body and

hence will take into account all of the information available through the Kinect Sensor in computing the probability that an observation sequence corresponds to a specific sign. This method introduces two problems however. The first is purely practical - the computations in the k-means clustering and training algorithms will become more and more expensive as the dimension of the observation vector grows and this will impact negatively on the performance of our sign recognition system.

The second issue is that this method ignores the prior knowledge we have about the joints of that body - that some are more dominant in the expression of a sign than others. For example, we know that the right hand is more dominant than the left and that both are more dominant than the elbows or shoulders. In fact the elbows and shoulders might be only worth considering when the hands alone cannot be used to determine two different signs.

As such we should not let the importance of each joint be determined algorithmically and should specify these parameters with care. For this we create J separate instances of **DHMM**, each trained on a collection of observations sequences of a different joint. Then supposing we have trained a set of Hidden Markov Models $\Lambda = \{\lambda_1, \dots, \lambda_J\}$ for each joint, and observations sequences $\mathcal{O} = \{\mathbf{O}_1, \dots, \mathbf{O}_J\}$ which specify a stream for each joint over the course of one sign, we can compute the log-probability that this sequence corresponds to the sign used to train $\lambda_1, \dots, \lambda_J$ as a weighted sum

$$\log(\mathbf{P}[\mathcal{O}|\Lambda]) = \sum_{j=1}^J c_j \log(\mathbf{P}[\mathbf{O}_j|\lambda_j])$$

Where the weights c_j are used to express the dominance of the j^{th} joint in expressing a sign relative to the other joints.

We implemented the class **SignModel** (Figure: 3.3) to correspond to this definition of Λ . This class contains a collection of **DHMM** objects with each corresponding to a joint of the skeleton provided by the Kinect Sensor, as well as a set of weightings corresponding to the c_1, \dots, c_J . This class contains a method **trainClassifier** which will train an instance of **DHMM** for each training data file in `/signAlign/Data/Training/name` where "name" is the name of the sign. The

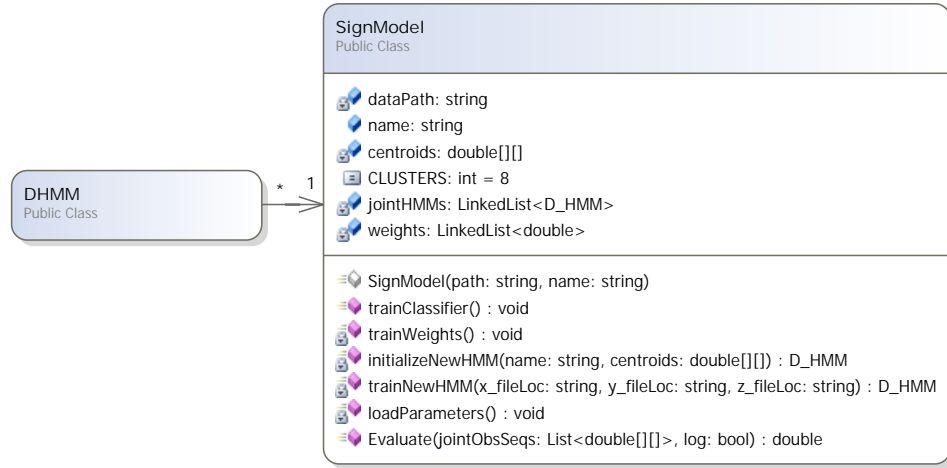


Figure 3.3: A UML diagram for the signModel class

`trainClassifier` class also contains a method to compute the probability that a given collection of joint observation sequences corresponds to this sign using the weighted sum of joints method described above.

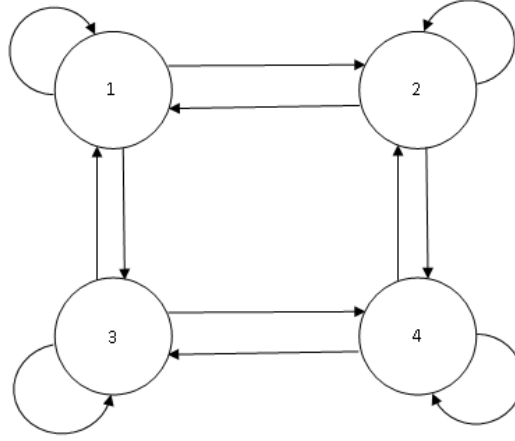
3.2.1 Determining The Initial Topology and Parameters

As our training method utilises a local search to improve our parameterisation its effectiveness will greatly depend on how we initialise our Hidden Markov Models. If the initial parameterisation is poor then even with a large training set we may fail to find a parameterisation which allows us to recognise signs with high enough accuracy.

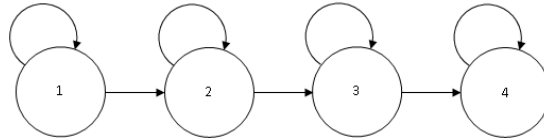
It has been shown that the topology of the underlying Markov Chain can greatly impact the effectiveness of a HMM for a given pattern recognition task (Jelinek, 1998; Rabiner, 1989). The two most common types of HMM (See figure 3.4) are the *ergodic model*, in which each state can transition to any other state in a finite number of steps and the *left-right model* or Bakis model (Bakis, 1976) in which the state sequence is (non-strictly) increasing. The left-right model can be viewed as modelling a signal which changes through time, with each transition representing a movement forward in time. For that reason it is more suitable for modelling gestures or speech than other Markov Chain topologies and is the

topology we choose for our Hidden Markov Models.

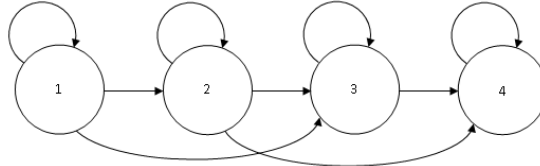
A left-right topology Markov chain is characterized by an upper triangular stochastic matrix A . We can further restrict our model to only allow jumps from a state i to states $i \leq j \leq i + \Delta$ for some Δ . This restriction can prevent the model being pushed, through re-parameterisation, into a trivial Markov chain which remains at the final state N at all times. This can be seen as imposing a restriction on the velocity we can expect a signer's hand to move at - not allowing him to pass by too many states in quick succession.



(a) An Ergodic Markov Chain



(b) A $\Delta = 1$ Left-Right Markov Chain



(c) A $\Delta = 2$ Left-Right Markov Chain

Figure 3.4: Different Markov Chain Topologies

As such we choose to initialize the HMMs in the `signModel` class as left-right

models restricted to $\Delta = 1$ and initialise the stochastic matrix with form

$$A = \begin{pmatrix} a_{11} & a_{12} & 0 & \cdots & 0 \\ 0 & a_{22} & a_{23} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_{(n-1)(n-1)} & a_{(n-1)n} \\ 0 & 0 & \cdots & 0 & 1 \end{pmatrix}$$

and define $\pi = [1, 0, \dots, 0]$. Clearly if we initialise the non-zero a_{ij} of A deterministically then we restrict ourselves to single outcome from the re-estimation procedure regardless of how many times we retry it, as such we choose to introduce some stochasticity by setting

$$a_{ii} = 0.5 - r \qquad a_{i(i+1)} = 0.5 + r$$

for some r chosen uniformly from the interval $[-0.3, +0.3]$. This interval is chosen to ensure that none of the links are set too weakly and the forward transitions broken during re-estimation.

As we expect the system to differentiate a large number of signs it would be a major task to tailor the initialisations of each Hidden Markov Model to the sign it is expected to recognise. Hence, we choose to use the same initial estimation procedure for each HMM and do not make any inferences about the structure of the emissions matrix B given the sign it is expected model. We initialise B almost uniformly with

$$b_{ij} = 1/M + r_{ij} \qquad \text{for each } 1 \leq i \leq N \text{ and } 1 \leq j \leq M$$

where each r_{ij} is some small random number and subject to the condition that the matrix B remains stochastic.

The benefit of introducing stochasticity into the initial parameter estimations is that if the re-estimation procedure does not produce a sufficiently good parameterisation we can re-initialise our HMM and try again. That is to say we may use the forward-backward algorithm as part of a random restart hill climbing algorithm ([Russell et al., 1995](#)) to increase the chances that we will find a good

parameterisation for the model.

3.2.2 Determining the Weightings

3.3 Building a Sign Classifier

A `signModel` object, when trained using suitable set of training data from a single sign, allows us to evaluate how likely a new observation sequence is to be an instance of the sign used to train the model. Hence if we have a trained `signModel` instance m_s for each sign s in some dictionary D and a collection of joint observation sequences $\mathcal{O} = \{\mathbf{O}_1, \dots, \mathbf{O}_J\}$ read from the Kinect Sensor then we can determine which sign model was most likely to have generated \mathcal{O} as

$$m = \operatorname{argmax}_{s \in D} \{\mathbb{P}[\mathcal{O}|m_s]\} \quad (*)$$

where the value $\mathbb{P}[\mathcal{O}|m_s]$ is can be computed by the `Evaluate` method of the `signModel` object. Then if $\mathbb{P}[\mathcal{O}|m]$ is sufficiently large then we can conclude that \mathcal{O} corresponds to the sign used to train m .

We implemented this means of determining a signs as the class `signClassifier` (Figure: 3.5) which contains for each sign a trained `signModel` object and a method `getSign` performs the calculation of (*).

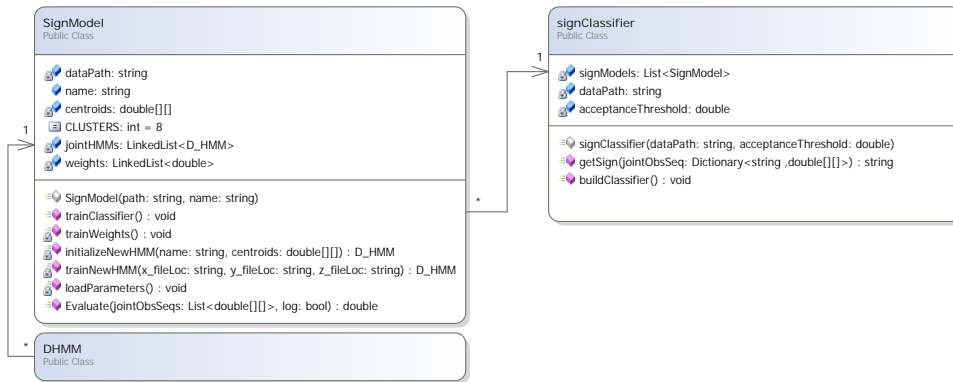


Figure 3.5: The `signClassifier` in UML

A trained instance of `signClassifier` forms the basis of the `SignAlign` and

provides a means of recognising signs using the Kinect Sensor as input. We designed the `buildClassifier` method to search for any parameterisations for `signModels` on disk and load them and then to search for any new training sets and create new `signModel` instances from them (saving their parameters to disk). As such the `SignAlign` dictionary can be extended by simply providing more training data without the need to retrain the rest of the model. This is a particularly desirable feature as there have been over 18 million Kinect Sensors sold worldwide, each providing the same format of joint and skeletal data. Hence the `SignAlign` system will be able to take training data from a wide variety of users if they choose to provide it.

The `signClassifier` and `signModel` classes contain a number of constants which we have not yet given a specified value, for example M the number of observation symbols or the acceptance threshold for detecting a sign. There is no a priori reason we should choose any particular value for these signs and chose determine them through experimentation (Chapter 5) and hence needed to build a catalogue of training data and test cases using the Kinect Sensor.

Chapter 4

Interfacing With The Kinect Sensor

The Kinect SDK provides access to the Kinect Sensor from within C# code by defining a `KinectSensor` object which interfaces directly with the hardware. When this object is instantiated, and provided a Kinect Sensor is connected via the USB port, it provides access to the depth, RGB and skeletal streams of the camera. The `KinectSensor` object also contains an event `AllFramesReady` which fires each time the streams have successfully updated. The streams update at approximately 30 frames per second but this rate can be significantly reduced if the sensor attempts to track too many skeletons. Using this `KinectSensor` object we built a general purpose skeletal tracking controller class `GestureController` (Figure: 4.1) which initialises a `KinectSensor` object to provide a single skeletal stream (for the left most person in the camera shot) and subscribes to the `AllFramesReady` event via the method `KinectAllFramesReady`.

4.1 Recording Training Data

In order to record the training data we implemented a class `GestureRecorder`, which extends `GestureController`, and a class `GestureRecording` which stores a single recording of the skeleton over time (Figure: 4.1). This second class contains a hashmap from `JointType` (an enum listing the different possible joints

tracked by the Kinect Sensor) to a list of 3-vectors representing the readings of that joint through time, and a method `addReading` which takes a `Skeleton` object and adds each joint location to the appropriate list. The `GestureRecorder` class stores a list of completed `GestureRecording` objects, along with one which is the current recording, and extends the `KinectAllFramesReady` method to pass the `Skeleton` provided by the sensor to the `currentRecording` object. This class also provides methods to start and stop the current recording (adding it to the list of recordings when stopped) and to save the list of recordings to file. This `GestureRecorder` class allows us to record the movements of the hands, wrists, elbows, shoulders and head over time and save them as training data. Further we implemented two methods for storing the data - as positions relative to the camera or as positions relative to the signer's head. This second method of storing positions prevents the same sign creating different readings when performed at different distances from the camera, in a different environment or by a different signer. It is these head relative distances that we use to train the sign models.

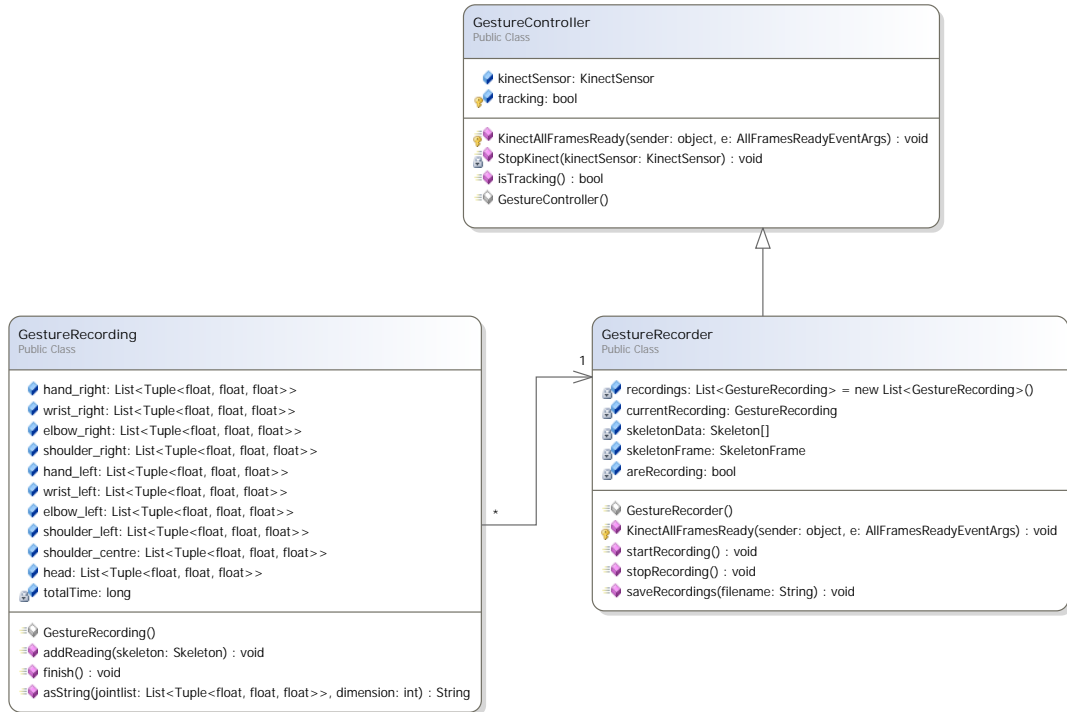


Figure 4.1: A UML diagram of the gesture recording classes

Using this controller we implemented a small program which can be used to record a collection of training data by repeating the sign the desired number of times and giving a start/stop signal between each. For convenience we defined this signal to be when hands are raised a certain distance above the waist, a signal consistent with the pose a signer takes whilst pausing between sequences of signs.

4.1.1 The Training Data

Using this controller we recorded for each of a collection of 30 British Sign Language (See Figure:) signs a set of 40 training examples and a further set of 10 testing examples. Each sign lasted between 1 and 2 seconds and was recorded at the maximum skeletal frame rate - 30 frames per second.

4.1.2 Limitations in the Training Data

The quality of the training sets was degraded by two factors. The first is that the signer (myself) that performed the training sets is not fully fluent in sign language and hence does not reproduce the signs consistently, this introduced inconsistencies into the training set and caused it to contain instances of signs that a fluent signer might not ever make. The second factor is a limitation of the Kinect Sensor. The Kinect can track joints well when they are isolated but performs poorly when two parts of the body, for example the hands, are in contact [4.2](#) and as many common signs require the hands to be placed together this issue caused a significant amount of the noise in the training data for those signs. The hand tracking issue is expect to be fixed in an upcoming SDK release but the contact issue could still caused noise in certain signs, for example those that require the signer to touch their elbows, shoulders or face.



Figure 4.2: An example of the broken hand tracking with the hands placed together. The yellow joints are inferred by the sensor incorrectly, the actual location of the hands is just below the chin.

Chapter 5

Testing and Experimentation

5.1 Determining the Classifier Parameters

5.1.1 The Cluster Number

5.1.2 The Acceptance Threshold

Given a collection of joint observations \mathcal{O} , the sign classifier determines the sign model m which is most likely to have generated it and if $\mathbb{P}[\mathcal{O}|m]$ is sufficiently large, say larger than some threshold τ , concludes that \mathcal{O} corresponds to the sign that trained m . However there is no *a priori* reason to choose a given value for probability threshold τ . If τ is chosen too small then the classifier is more likely to associate a non-sign observation sequence with a sign, resulting in a false positive. Conversely, if τ is chosen too great then the classifier may fail to associate the observation sequence of a sign with the appropriate sign model, resulting in a false negative. To determine the appropriate value of τ we trained the classifier on the training set and tested the classifier on half of the test data for values of $\log(\tau)$ from -4000 up to 0 in increments of 100 . We recorded for each value of $\log(\tau)$ we recorded the number of false positives and false negatives and plotted them together (See Figure:), from this we determined that $\log(\tau) = -610$ to be an appropriate choice. It should be noted that we cannot test the values of τ on the full testing set as this constitutes training the data on the test set and could overfitting in our model.

5.2 Sign Recognition Accuracy

Once parameterised we

Chapter 6

Analysis and Conclusions

Chapter 7

Additional Work

An immediate improvement that could be made to the sign recognition system is to gather training data from a user fluent in British Sign Language. This would likely improve the recognition accuracy of our system as the training sets would be more consistent. As we have implemented an efficient way to gather training sets through the Kinect Sensor this could be achieved fairly easily for a small training set, however it would still be prohibitively time consuming for a single user to generate the training sets for a large dictionary of signs. As the Kinect Sensor is widely available this problem might be solved by crowd-sourcing the training database and allowing users of the system to submit their own training data for new signs.

An extension of the system to track hands and factor their positions into predicting signs would likely yield significant improvements to sign recognition. For example, the signs for "you" and "your" are essentially the same if the hand is ignored and are differentiated by a pointing finger or the hand forming a fist and so cannot be accurately detected by the current system. At present the Kinect SDK does not provide information about the locations of the fingers or shapes of the hands and so the system cannot be extended to use hand shapes without using third party software or implementing our own finger tracking, however it is expected that the Kinect will soon be updated to provide hand tracking and so we may be able to use hand tracking in the near future using only the Kinect

SDK.

If hand tracking is implemented then we will be able to significantly extend the functionality of our sign recognition system to detect finger spelling. Finger spelling is commonly used in sign language to communicate words with no sign equivalent and is a key part of communicating in sign language. As such the ability to accurately finger spell is integral to a full sign language translation system and it further can be used to overcome the problem of not having training data for a large dictionary - by allowing a user to finger spell those words not yet available. The finger spelling problem consists of detecting between 26 distinct letter signs and could be solved using a similar Hidden Markov Model technique to the one presented in this project; by replacing joint tracking with finger and knuckle tracking. As we have obtained a good accuracy in classifying a small dictionary of signs it suggests that this method might produce good results in classifying the different letter signs of finger spelling.

Another feature of sign language which we have ignore in this project is facial expressions. Facial expressions are often used to modify the signs being communicated by the hands. For example the sign "you" can be converted to "are you ... " by raising the eyebrows. The Kinect SDK has been used to track facial expressions ([Microsoft, 2013a](#)) and this work could be intergrated with the sign recognition to produce a more robust sign recognition system.

Appdx A

Appdx B

References

- James Baker. The dragon system—an overview. *Acoustics, Speech and Signal Processing, IEEE Transactions on*, 23(1):24–29, 1975. [11](#)
- Raimo Bakis. Continuous speech recognition via centisecond acoustic states. *The Journal of the Acoustical Society of America*, 59:S97, 1976. [19](#)
- Leonard E Baum and Ted Petrie. Statistical inference for probabilistic functions of finite state markov chains. *The Annals of Mathematical Statistics*, 37(6):1554–1563, 1966. [6](#)
- Leonard E Baum and George R Sell. Growth transformations for functions on manifolds. *Pacific Journal of Mathematics*, 27(2):211–227, 1968. [7](#), [9](#)
- Leonard E Baum, Ted Petrie, George Soules, and Norman Weiss. A maximization technique occurring in the statistical analysis of probabilistic functions of markov chains. *The annals of mathematical statistics*, pages 164–171, 1970. [6](#), [7](#), [8](#)
- Christopher M Bishop et al. *Pattern recognition and machine learning*, volume 4. springer New York, 2006. [8](#)
- Horst Bunke, Markus Roth, and Ernst Günter Schukat-Talamazzini. Off-line cursive handwriting recognition using hidden markov models. *Pattern recognition*, 28(9):1399–1413, 1995. [6](#)
- Philippe Dreuw, Pascal Steingrube, Thomas Deselaers, and Hermann Ney. Smoothed disparity maps for continuous american sign language recognition. *Pattern Recognition and Image Analysis*, pages 24–31, 2009. [11](#)

REFERENCES

- Philippe Dreuw, Jens Forster, Yannick Gweth, Daniel Stein, Hermann Ney, Gregorio Martinez, Jaume Verges Llahi, Onno Crasborn, Ellen Ormel, Wei Du, et al. Signspeak—understanding, recognition, and translation of sign languages. In *Proceedings of 4th Workshop on the Representation and Processing of Sign Languages: Corpora and Sign Language Technologies*, pages 22–23, 2010. 11
- Richard Durbin, Sean R Eddy, Anders Krogh, and Graeme Mitchison. *Biological sequence analysis: probabilistic models of proteins and nucleic acids*. Cambridge university press, 1998. 6
- A. A. Argyros I. Oikonomidis, N. Kyriazis. Kinect 3d Hand Tracking, March 2013. URL <http://cvrlcode.ics.forth.gr/handtracking/>. 3
- Frederick Jelinek. *Statistical methods for speech recognition*. MIT press, 1998. 6, 19
- Frederick Jelinek, L Bahl, and R Mercer. Design of a linguistic statistical decoder for the recognition of continuous speech. *Information Theory, IEEE Transactions on*, 21(3):250–256, 1975. 11
- Biing-Hwang Juang and Lawrence R Rabiner. Hidden markov models for speech recognition. *Technometrics*, 33(3):251–272, 1991. 6
- Dan Jurafsky, James H Martin, and Andrew Kehler. *Speech and language processing: an introduction to natural language processing, computational linguistics and speech recognition*, volume 2. MIT Press, 2002. 6
- Anders Krogh, Michael Brown, I Saira Mian, Kimmen Sjolander, and David Haussler. Hidden markov models in computational biology: Applications to protein modeling. *Journal of molecular biology*, 235(5):1501–1531, 1994. 6
- Stephen G. LATTA, Kudo TSUNODA, Kevin GEISNER, Relja MARKOVIC, Darren Alexander BENNETT, and Kathryn Stone PEREZ. Gesture keyboarding, 08 2010. URL http://www.patentlens.net/patentlens/patent/US_2010_0199228_A1/en/. 3

REFERENCES

- Stephen E Levinson, Lawrence R Rabiner, and Man Mohan Sondhi. An introduction to the application of the theory of probabilistic functions of a markov process to automatic speech recognition. *Bell Syst. Tech. J*, 62(4):1035–1074, 1983. 9
- Xiaolin Li, Marc Parizeau, and Réjean Plamondon. Training hidden markov models with multiple observations-a combinatorial method. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 22(4):371–377, 2000. 14
- Pietro Liò and Nick Goldman. Models of molecular evolution and phylogeny. *Genome research*, 8(12):1233–1244, 1998. 6
- Stuart Lloyd. Least squares quantization in pcm. *Information Theory, IEEE Transactions on*, 28(2):129–137, 1982. 16
- Christopher D Manning and Hinrich Schütze. *Foundations of statistical natural language processing*. MIT press, 1999. 6
- Microsoft. Face Tracking, March 2013a. URL <http://msdn.microsoft.com/en-us/library/jj130970.aspx>. 32
- Microsoft. What’s New - Kinect SDK 1.7, March 2013b. URL <http://www.microsoft.com/en-us/kinectforwindows/Develop/New.aspx>. 3
- Todd K Moon. The expectation-maximization algorithm. *Signal Processing Magazine, IEEE*, 13(6):47–60, 1996. 8
- Lawrence R Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989. 4, 6, 8, 12, 13, 14, 19
- Ali Rahimi. An erratum for ‘a tutorial on hidden markov models and selected applications in speech recognition’, April 2000. URL <http://xenia.media.mit.edu/~rahimi/rabiner/rabiner-errata/rabiner-errata.html>. 13
- Stuart Jonathan Russell, Peter Norvig, John F Canny, Jitendra M Malik, and Douglas D Edwards. *Artificial intelligence: a modern approach*, volume 74. Prentice hall Englewood Cliffs, 1995. 21

REFERENCES

- Thad Starner and Alex Pentland. Real-time american sign language recognition from video using hidden markov models. In *Computer Vision, 1995. Proceedings., International Symposium on*, pages 265–270. IEEE, 1995. [11](#)
- Jean-Paul van Oosten. *Can markov properties be learned by hidden markov modelling algorithms*. PhD thesis, Masters thesis, University of Groningen, The Netherlands, 2010. [14](#)
- Junji Yamato, Jun Ohya, and Kenichiro Ishii. Recognizing human action in time-sequential images using hidden markov model. In *Computer Vision and Pattern Recognition, 1992. Proceedings CVPR'92., 1992 IEEE Computer Society Conference on*, pages 379–385. IEEE, 1992. [10](#), [11](#)