

Poznan University of Technology
Faculty of Computing
Institute of Computing Science

Bachelor's thesis

**GESTURE RECOGNITION LIBRARY FOR LEAP MOTION
CONTROLLER**

Michał Nowicki, 95883
Olgierd Pilarczyk, 100449
Jakub Wąsikowski, 101560
Katarzyna Zjawin, 98826

Supervisor
dr inż. Wojciech Jaśkowski

Poznań, 2014

Tutaj przychodzi karta pracy dyplomowej;
oryginał wstawiamy do wersji dla archiwum PP, w pozostałych kopiach wstawiamy ksero.

Contents

1 Abstract	1
2 Introduction	2
2.1 Scope of work	2
2.2 Motivation	2
2.3 Objectives	3
2.4 Thesis organization	3
3 Introduction to gesture recognition	4
3.1 Classification of gestures proposed in literature	4
3.2 State of the art methods	6
3.2.1 Enabling Technology	6
Non-vision-based Technology	6
Vision-based Technology	6
3.2.2 Gesture representation	7
3d model-based	7
Appearance-based	8
3.2.3 Gesture recognition methods	8
Static gesture recognition	8
Dynamic gesture recognition	10
4 Leap Motion controller	11
5 Gesture recognition for Leap Motion	12
5.1 Classification of gestures	12
5.2 Gesture data representation	14
5.3 Additional processing steps	14
6 Static gestures recognition	15
6.1 Proposed methods	15
6.2 Evaluation methodology	17
6.2.1 Assumptions	17
6.2.2 Recorded datasets	18
6.3 Experiments	19
6.4 Finger differentiation	23
6.4.1 Methods	23
6.4.2 Evaluation methodology	23
Recorded dataset	23

6.4.3	Features	23
6.5	Experiments	24
7	Detection of dynamic gestures	26
7.1	Proposed methods	26
7.1.1	Hidden Markov Model	26
7.1.2	HMM observation from Leap Motion data	29
7.2	Evaluation methodology	31
7.3	Experiments	31
8	LMGesture library dedicated for Leap Motion controller	37
8.1	Architecture	37
8.1.1	Model	37
8.2	Processes	38
8.2.1	The learning process	38
8.2.2	The recognition process	39
8.3	Gesture recorder and visualizer	40
8.3.1	LMR files	40
8.3.2	Visualizer	41
8.3.3	Recorder	42
8.4	Used libraries	42
8.5	Samples of code using the library dedicated to Leap Motion Controller	42
9	Conclusions	44
	Bibliography	45

Chapter 1

Abstract

Nowadays, when computers are ubiquitous, people need to develop more natural and intuitive interface to use computer. People would like to use gestures of everyday life and translate them into a virtual world. This paper presents library for gesture recognition dedicated to Leap Motion Controller called LMGesture. Leap Motion is a new device, which tracks fingers and other objects up to 1/100th of a millimeter. The LMGesture library can be used for various kinds of gestures. This can be done, because of the use of different detection methods for different types of gestures. For static gestures recognition has been used support vector machine (SVM) and for dynamic gestures – hidden Markov model (HMM). In library can be found additional modules: recorder – for recording gestures in a format supported by the library, visualizer – for reviewing recorded gestures, finger differentiation module – for differentiating fingers in a performed gesture. In this paper has been also presented classification of gestures and its modification in the context of Leap Motion Controller. Additionally this work includes descriptions of helper methods used for the gesture recognition, such as pre-processing of data obtained from the device, which gets rid of existing noise or method for fingers differentiating, whereby the obtained results are more accurate.

[dodac opis testow]

Chapter 2

Introduction

2.1 Scope of work

The objective scope of this paper includes:

- designing library architecture,
- selection and implementation of algorithms for gesture recognition,
- implementation (or learning) of gestures built-in library,
- implementation examples of the library,
- tests using Leap Motion Controller.

Following thesis concerns machine learning. In this paper were used two algorithms from this field of science: support vector machine and hidden Markov model. SVM is included in the group of supervised learning, where algorithm knows set of input data and responses to the data, and tries to create a predictor model that generates reasonable predictions for the response to new data. HMM is an example of unsupervised learning, where algorithms are trying to find hidden structure using unlabeled data.

Subjective scope of this work is to examine Leap Motion Controller in gesture recognition. This device is an innovative approach to the computer usage. The scope of this work is to examine this controller and create an interface between it and the user.

Time range of the thesis is October 2013 – January 2014.

2.2 Motivation

Nowadays computer usage is not a natural human behavior. To rotate the object in the virtual world user need to click the mouse and move it in 2D plane. Man rotating objects in the real world has to catch it and turn using hands. To make the usage of computer more intuitive and natural the best solution would be to transfer gestures performed on a daily basis into the virtual world. The latest technological solutions allow to control computer using gestures. Hands are a fundamental tools of every human being. With them people perform hundreds of operations every day. Without this basic manipulator people are not able to cope with the simplest activities. Hands give a large scope of activities and gestures. They are even used for non-verbal communication. Using this tool man at a low cost can do and achieve almost everything. Currently there are several devices that support gesture recognition. An example of such a controller is Kinect, but it is not highly accurate. It is suitable for applications that use the whole body of the user, but for the recognition

of hand gestures may be no useful. Leap Motion Controller is a small device that can be placed in front of the computer. Its operating range is between user and PC, and its accuracy is very high. It is an ideal device for recognizing hand gestures through which people have the opportunity for more natural user interface of the computer.

2.3 Objectives

The main purpose of this paper is creation of library, which recognizes hand gestures using Leap Motion controller. Additionally, below bulleted objectives will also be realized:

- comparison and evaluation of existing methods in the context of hand gesture recognition,
- creating a module to manage gestures,
- creating an initial database of gestures to the library,
- performing hand gesture recognition quality tests,
- presentation of example of library usage.

2.4 Thesis organization

The structure of the paper is as follows. Chapter 5 presents an overview of the literature on gesture recognition. In this part there are descriptions of gestures classification and known methods of gesture recognition. Chapter 6 is a presentation of Leap Motion Controller. Chapter 7 is devoted to gestures recognizing in the context of Leap Motion Controller. There are descriptions of gestures classification, data representation and additional processing steps for hand gestures recognition using the Leap Motion device. Chapter 8 contains proposed methods, evaluation methodology and experiments for static gesture recognition. Chapter 9 presents a similar description for dynamic gesture recognition. In chapter 10 has been described created library - its architecture, processes, additional modules, used libraries for its implementation and an example of its usage. Chapter 11 provides conclusions of the paper.

Michał Nowicki - svm do statycznych, hmm do dynamicznych Olgierd Pilarczyk - projektowanie architektury, implementacja preprocesingu Jakub Wąsikowski - visualizer, recorder, rozróżnianie palców, wyznaczanie liczby klastrów, klasyfikacja gestów Katarzyna Zjawin - visualizer, recorder, rozróżnianie palców, klasyfikacja gestów

Chapter 3

Introduction to gesture recognition

3.1 Classification of gestures proposed in literature

The vast multiplicity of gestures that human can perform makes the number of classes – to which we can divide these gestures – substantial. Therefore the classification can be performed in different ways, taking into account the different characteristics of gestures. Most of presented theories include the knowledge that originates from a variety of science such as anthropology, linguistics, cognitive science and other. In this chapter it is provided review of the most common gesture classifications in Human Computer Interaction (HCI) context. It was focused mainly on gestures that are relate to hand and arm movements.

The basic classification of gestures is the division into static and dynamic gestures. Group of static gesture includes fixed gestures which are not take into account the changes in time. Dynamic gestures is group of time varying gestures.

There is also another general division considered by Kammer et al. [17] due to the type of actions activated by the gesture, to online and offline. The first group includes gestures which are processed during performing. It is a group of direct manipulation gestures that provide additional information about dynamics of gesture. There are often used to manipulate objects in space. The second one is group of action gestures, which are processed at the end of gesture. Most often these are gestures that convey the occurrence of specific meaning.

Karam and schraefel [18] proposed more extensive gesture taxonomy dedicated for Human Computer Interaction. Their classification is based on Quek et al. [29] publication, which provides clarification of gestures taxonomies presented in the past literature. Karam et al. defined following gesture classes: deictic, gesticulation, manipulation, semaphores and sign language. In another publication, Aigner et al. [32] presented classification, which are more tailored for hand gesture recognition purposes, drawing on concepts from Karam and Schraefel work. They distinguished five categories of gestures:

- pointing – used to point object or indicate direction. Formally, it involve pointing to establish the identity or spatial location of an object within the context of the application domain [18]. It applies not only an indication by the index finger but also to any finger and any number of fingers. It is also independent of finger orientation and curvature, while gesture has a indication meaning. Equivalent to deictic gesture in Karam and Schraefel literature.
- semaphoric – group which consists gesture posture and dynamics, which are used to convey specific meanings. Formally, Semaphoric approaches may be referred to as “communicative” in that gestures serve as a universe of symbols to be communicated to the machine [29].

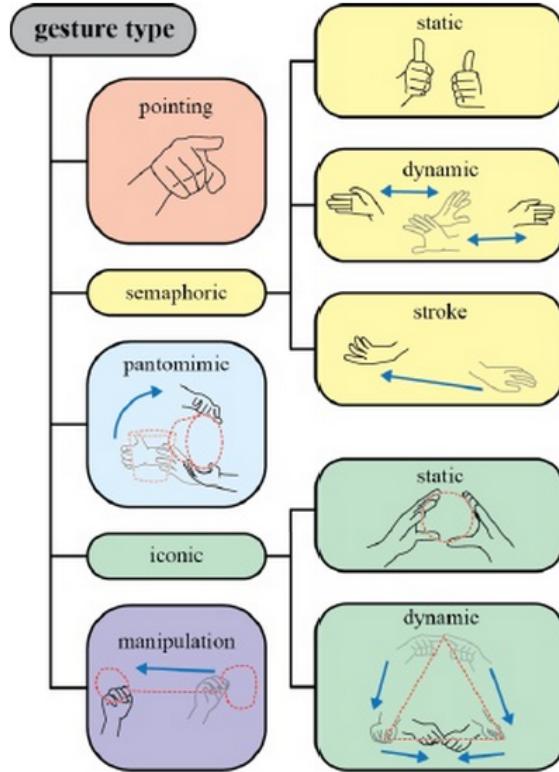


FIGURE 3.1: The classification of gestures proposed by Aigner et al. [32]

Due to the fact that semaphoric are symbolic gestures, their layout can be unrelated to their meaning. Distinguished three semaphoric gesture types of static, dynamic and strokes. The first one concerns specific hand posture, such as thumbs-up meaning approval symbol. Dynamic semaphorics convey their meaning through movement, for example waving of hand to greet somebody. The last one group are similar to dynamic semaphorics gesture, but this represents fast, stroke-like movements, such as swipe gesture.

- **iconic** – used to demonstrate shape, size, curvature of object or entities. In contrast to the semaphoric gestures, their layout or motion path is strictly related to their meaning. Iconic gestures can be divided into static and dynamic. First one are performed by hand postures, such as rectangle formed by the thumb and forefingers of both hands. Dynamic iconic gestures are often used to map edge line of objects by means of motion paths. For instance showing a simplified sine function characteristics with finger movements.
- **pantomimic** – presents imitated perform of specific task or activity without any tools or objects. Pantomimic gesture are characterized by a high variability of posture and movements. An example of this gesture type can be weapon reload or movement of a knife slicing bread.
- **manipulation** – used to control the position, rotation and scale of the object or entity in space. Manipulation gestures constitutes a direct interaction between the manipulated object and hand or tool that performs gesture. It follows, that the movement of the manipulated object must be strictly dependent on the motion gesture.

3.2 State of the art methods

In this chapter presented review of the state-of-the-art in human gesture recognition. The problem of gesture recognition can be divided in two main problems: the gesture representation problem and the decision/inference problem. Therefore, review includes discussion about enabling technology, gesture representations and analysis of recognition methods. Additionally introduced general problems related to the recognition of gestures and their common solutions.

3.2.1 Enabling Technology

In this subsection, overviewed the enabling technology for gesture recognition. The main existing gesture recognition approaches related to type of the devices are as follow:

- non-vision-based devices – tracking devices, instrumented gloves, armbands and others,
- vision-based devices – using one or many cameras.

Non-vision-based Technology

This type of devices uses various technologies to detect motions, such as accelerometers, multi-touch screens, EMG sensors and other, which includes several detectors. There are few categories of non-vision-based devices [16]:

- wearable – these kind of device is in the form of garment, which includes sensors needed to recognize arrangement and motions of examined part of body. Often occur in the form of gloves (CyberGlove®), armband (Myo) or the whole outfit (IGS-190). For instance, CyberGlove® device was used in system developed by Kevin et al. [19], which recognize multi-dimensional gestures using condensation-based trajectory matching algorithm. These devices are often related to biomechanical and inertial technologies,
- biomechanical – type of device, which use biomechanical techniques such as electromyography, to measure parameters of gesture. Example of using this type of device is project developed by Kim et al. [20] for Realtime Biosignal Interfacing based on EMG sensors. Example of these devices is Myo armband, which detects gestures and movements using EMG sensors,
- inertial – these devices measure the variation of the earth magnetic field in order to detect the motion. This kind of devices uses accelerometers [22] and gyroscopes [14] to measurements,
- haptics – various kinds of touch screens. For instance, Webel et al. [39] developed module for dynamic gestures recognition in multi-touch devices,
- electromagnetic – these devices measure the variation of an artificial electromagnetic field deriving from wireless networks, electronic devices or produced by themselves. Example of such devices is WiSee, which leverages ongoing wireless transmissions in the environment (e.g., WiFi) to enable whole-home sensing and recognition of human gestures [28].

Vision-based Technology

Vision-based devices include one or several cameras and provide performed data from the captured video sequences. Processing of frame is based on filtering, analyzing and data interpreting. The following types of vision-based technology can be distinguished [16][40]:

- typical video cameras – gesture recognition techniques based on data derived from monocular camera using detection methods such as color or shape based techniques, learning detectors from pixel values or 3D model-based detection,
- stereocameras – techniques based on captured images from two cameras, which provide an approximation of the recorded data to a 3D model representation,
- active techniques – require the projection of some form of structured light. Examples of this kind devices are Kinect or Leap Motion,
- invasive techniques – systems which require using of body markers such as color gloves [38], LED lights (Play Station Move controller).

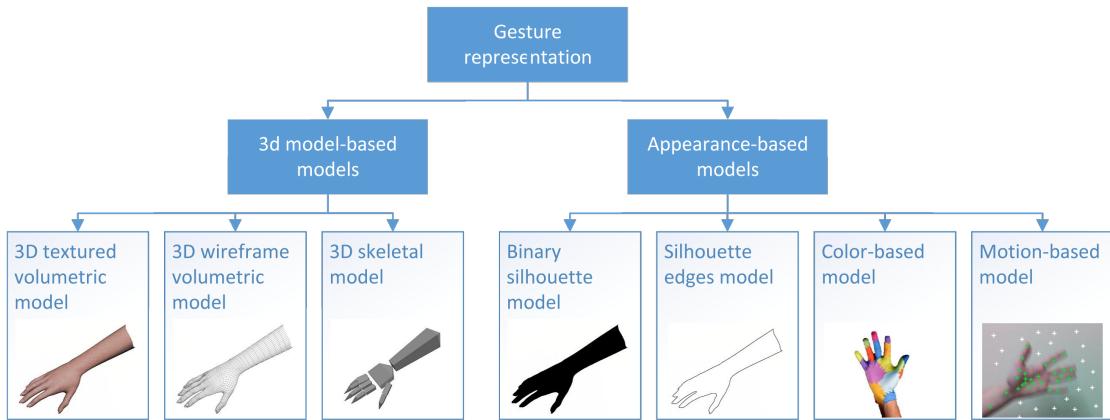


FIGURE 3.2: Diagram of gesture representation

3.2.2 Gesture representation

In this subsection, provides overview the spatial modeling of gestures. In particular, it was focused on one type of spatial gesture representations namely 3d model-based. Depending on the type of the input data, the approach for recognize gesture could be done in different ways. There are following two main types of gesture representation defined in literature [15][25]:

- 3d model-based
- Appearance-based

3d model-based

Defines the 3D spatial decription of the human body parts. They can be classified in two large groups:

- volumetric models,
- skeletal models.

Volumetric models reproduce with high accuracy the shape of the hand or arm. Real object is often interpreted as ordered in space mesh of vertices or NURBS. This model is commonly used for computer vision purposes or for computer animation. The drawback of this approach is that is very demanding computationally and difficult to analyze in real time. With the power of today's computers, simplified models of the hand or arm (such as skeletal models) are more recommended.

As indicated earlier, instead of dealing with all the parameters of volumetric type, model can be reduced to set of equivalent joint angle parameters together with segment lengths. These are known as skeletal models. There are several advantages of using this type:

- Simplified model with the most important parameters allows the detection program to focus on the significant parts of the body.
- Due to the smaller amount of data, processing algorithms are faster.

Appearance-based

The second group of models do not use direct description of the spatial object points, because this model is based on shape of hands or arms in the visual images. The gestures are modeled by relating the appearance of any gesture to the appearance of the set of predefined, template gestures [25]. In this group distinguished a large variety of models. The most used 2d models are:

- color based model – in general, using body markers to track the motion of the body part,
- binary silhouette based model – models based on the geometric properties of the object silhouette,
- deformable gabarit based model – they are generally based on deformable active contours,
- motion-based model – based on the motion of individual pixels or image part description.

3.2.3 Gesture recognition methods

As was written earlier, one of the main issues of gesture recognition is decision problem. Currently several solutions has been proposed, which can be used regardless of device type or data representation for different classes of gestures. Classification of gestures, which should be taken into account while choosing gesture recognition method was being chosen is static and dynamic division. For each of them may be used different tools due to the different properties of these gestures. In the case of static gesture recognition, an important feature is arrangement of object, which performs gesture. In other words how the individual parts of the object are arranged in relation to each other. For dynamic gestures – as described in the previous subsection – a very important feature is the variation in time (dynamics of the gesture or dynamics of individual parts of the object performing the gesture).

To recognize static gestures, general classifier, neural network or template-matcher can be used. Methods which are capable to recognize dynamic gestures have to take into account an aspect of time. The example of this kind of method is Hidden Markov Model.

Static gesture recognition

Different techniques to perform accurate static gesture recognition have been decribed in the literature. The most common methods are neural networks, support vector machines and simple pattern techniques [34].

An neural network (NN) is an information-processing system that has been developed as generalizations of mathematical models of human cognition or neural biology. A neural network is characterized by its pattern of connections between the neurons, method of determining the weights on the connections, and its activation function [11]. They can be used both for static and dynamic gestures.

Hasan et al. [13] presented hand gesture recognition based on shape analysis. Tests were conducted for six static gestures using multi-layer perception of neural network and back-propagation learning algorithm. NN architecture consisted of one hidden layer (100 nodes), 1060 inputs and 6 output for each gesture. They achieved a recognition rate of 86.38% for a training set of 30 images and a testing set of 84 images.

Xu et al. [41] developed virtual training system of Self-Propelled Gun based on static gesture recognition and hand translations and rotations. The input data for algorithms was captured using a 18-sensor DataGlove. To recognize gestures was used feed-forward neural network with 40 nodes in single hidden layer, 18 input and 15 output nodes. The back-propagation using a variable learning rate is selected as training method. The tests were conducted on a set of 300 hand gestures from five different people – 200 gestures for training set and 100 for testing set. With the use of these methods the authors reached gesture recognition performance of 98%.

In publication of Stergiopoulou and Papamarkos publication [37] can be found static gesture recognition through other type of neural network – Self-Growing and Self-Organized Neural Gas (SGONG). To quote the authors, SGONG is innovative neural network that grows according to the morphology of hands in a very robust way. The algorithms were tested for 31 hand gestures that derive from the combination of shown and hidden fingers. Data were collected from a camera, and the recordings of hand were created in a vertical position on uniform background. With these assumptions, gesture recognition rate of 90.45% have been reached but required average computation time was about 1.5 s, using a 3 GHz CPU.

Support-Vector Machine (SVM) is a classification method invented by Vapnik [7]. SVM is a supervised learning algorithm used for classification and regression analysis, based on the mapping of characteristics extracted from instances namely the feature vectors to points in space. SVM constructs in multi-dimensional space, set of hyperplanes, which non-linearly divide points in this space (input vectors) to different classes. Support Vector Machines can be called a maximum margin classifier, because the resulting hyperplanes maximize the distance between the 'nearest' vectors of different classes. These "nearest" vectors are called support vectors.

Chen and Tseng [5] presented system based on training SVM which allows effective recognition gesture in popular game, rock-paper-scissors. One of the challenges of their work was to teach the classifier to recognize multiple-angle hand gesture. The collection of training and testing data were images from video camera, which are preprocessed using conversion to grayscale and histogram equalization. Data were collected from 5 different people for the right hand only. For the learning set consisted of 420 images and testing set of 120 images, the recognition rate of 95% was achieved.

Rahman and Afrin [30] presented hand gesture recognition system which recognizes static hand gesture for alphabet of 10 letters using Biorthogonal Wavelet Transform and SVM. Input data in the form of images – in addition to filtering – are transformed by the Canny edge detection method and then processed sequentially through Radon and Biorthogonal Wavelet Transformations. Finally, the data in this form are transmitted to the SVM classifier. To achieve robustness of the method to varying conditions, authors used a large dataset – 800 positive samples and 1500 negative image samples. Average recognition rate was 87.4%.

Liu et al. [23] proposed recognition method based on SVM and Hu moments which applied to Chinese Driver Physical Examination System. For collection of 2416 positive samples and 3050 negative samples from 20 people recognition rate of 96.5% have been reached.

Ren and Zhang [31] proposed other recognition method named by them as MEB-SVM. This method combines the SVM with minimum enclosing ball (MEB) and – according to the authors – allows to reduce computation with effective separation of all kinds of vectors in hyperspace. The input data used to test this method are images which are initially binarized and then countour

line is retrieved. Finally, contour line is converted by means of Fourier transform, so that data are independent of translation, rotating and zooming. Their method achieved a recognition rate of 92.9%.

Dominio et al. [9] presented novel hand gesture recognition based on depth data using Kinect device. The proposed processing consists following, main steps: extraction hand region from the depth map and subdivided it into palm and finger samples, extraction set of features based on finger tips and center of the hand, classification by SVM. Based on 1,000 different depth maps with 10 gestures performed by 10 different people, they achieved mean recognition rate of 99.5%.

Other popular methods are *simple pattern recognition techniques*. This group includes methods based on a simple comparison the characteristics of new problem instance with instances seen in training, instead of performing explicit generalization. In the case of gesture recognition, output information of algorithm is evaluated on the basis of similarity of the gesture to other pre-defined or learned gestures, for which belonging to groups is known. Basis on this, it is concluded that the newly read gesture belongs to the group. These techniques are generally based on a efficient lazy learning methods such as instance-based learning methods. In the context of gesture recognition, the most widely used algorithm is the k-nearest neighbour.

Ziaie et al. [43] combine naïve Bayes classifier with k-nearest neighbors algorithm for static gesture recognition purposes. Based on the dataset from the camera consisting of 580 samples for the 3 gestures, authors have achieved 93% of recognition rate.

Chang et al. [4] proposed new approach for static gesture recognition based on Zernike moments (ZMS) and pseudo-Zernike moments (PZMs), which provide greater independence from the translation of gestures. For gesture classification, k-nearest neighbor has been used, which directly processes the data from the ZMS and PZMs. In addition, authors used a minimum bounding circle, which supports the decomposition of hand silhouette into the finger parts and palm part. For dataset of 600 samples with 6 gestures performed by 10 people, gesture recognition rate of 97.3% has been reached.

Dynamic gesture recognition

Referring to the previously proposed classification of gestures may be considered that dynamic gestures include a wide range of existing types of gestures and constitutes an important role in interpersonal communication as well as in HCl. Therefore, many approaches were proposed to gesture recognition taking into account the temporal aspect of gestures. Shen et al. [35] propose following division of these approaches:

- model-based methods,
- exemplar-based methods.

The first one includes methods that are based on the 3D model-based gesture representation, and which assume that the hand was detected and hand movements are tracked. Model-based approaches include Hidden Markov Models (HMM), Finite State Machines (FSM), Time Delay Neural Networks (TDNN) and self-organizing networks, which preserve topology.

Exemplar-based is group of methods, which perform recognition by exemplar or template matching. These methods use a visual representation of data such as spatio-temporal local features, motion trajectories, bag-of-features representation and other. However, this approach of dynamic gesture recognition will not be considered in thesis.

Chapter 4

Leap Motion controller

Chapter 5

Gesture recognition for Leap Motion

5.1 Classification of gestures

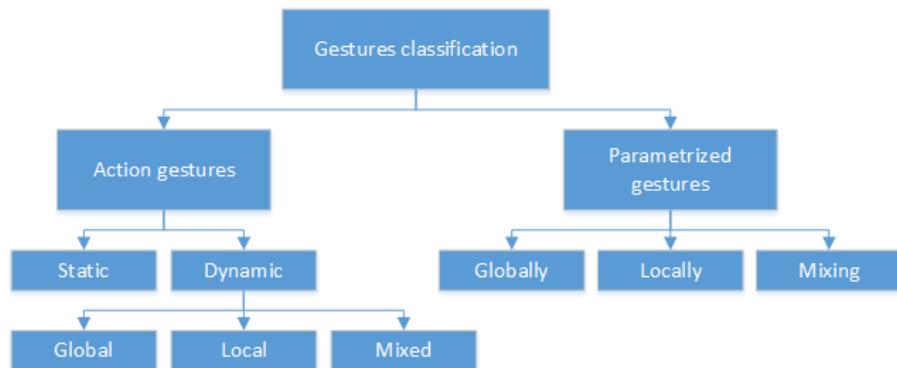


FIGURE 5.1: Classification of gestures proposed in the thesis

In section 3.1 have been presented the three main classifications of gestures based on the existing literature. Particularly important for further research and to properly define the problem of gesture recognition may be the last of the presented taxonomies, which defines a wide range of possible gestures occurring in human-computer interaction.

It should be noticed that the previously mentioned classifications are defined strictly due to the characteristics of individual gestures, which can make it difficult to analysis of choosing methods for gesture recognition. Therefore, it was decided to propose a new division that will be based on early taxonomies and in addition will be focused on the gesture recognition methods and will be defined strictly in the context of gesture recognition.

The proposed classification distinguishes basic division of gestures due to conveyed informations:

- action gestures,
- parametrized gestures.

The first group (1), represents gestures, which their recognition rely only on detection of gesture occurrence. This means that action gestures have assigned some meaning, and the occurrence of this type of gesture implies execution of pre-defined action. It should also be noted that gestures of this type do not convey any additional parameters which define or describe the action. An example of this type can be gesture showing open hand, which could mean stop playing music. Another

example might be a moving index finger in circles, which means rotating previously selected object of the specified number of degrees.

Gestures, which may be included to parameterized (2), are recognized on the basis of gesture detection and returning parameters associated with the context of the gesture. In contrast to previous group, parametrized gestures carry additional meaning or parameters needed to perform the gesture. An example of this kind of gestures might be a similar gesture to the previously mentioned instance of action gestures – making circles through moving finger – but in this case – in addition to returning information about recognized gesture – a value of angle, which should be used to rotate selected object, will also be conveyed.

Then action gestures (1) were divided into static (1.1) and dynamic (1.2). The first one relates to gestures, that are not variable in time, therefore hand posture, its position and rotation are fixed during hand gesture performing. It should also be mentioned that these gestures should be independent of the orientation relative to the environment and the occurrence of this gesture is detected only on the basis of hand posture.

Dynamic action gestures (1.2) refer to group of gestures, which are continuously and non-parametrized variable in time in terms of hand posture, position and rotation in space. This group is further divided into global, local and mixed. Global were detected based on a fixed hand posture and specified movement, which may relate both to changes in position and rotation. In the case of local, hand posture is only variable in time, while the mixed include gestures, for which both the hand posture, as well as the position and rotation vary with time.

For parameterized gestures (2) distinguished division of parameterized locally (1.1), globally (1.2) and mixed. Parameterized globally include gestures, whose parameters are determined by values of position and rotation or its changes of the hand. For instance, swipe gesture can be considered as globally parametrized gesture, where the parameter is the length of the swipe motion. Recognition of globally parametrized gestures is performed on the basis of unchanging hand posture or specific changes in the shape of the hand posture over time.

Locally parameterized group (1.2) include gestures whose the hand posture or its changes are parameterized, for example a distance between index finger and thumb of one hand may be a parameter for scaling gesture. Gestures of this kind should be independent of the position and rotation in space. Recognition should be based on the hand posture, which can be problematic when it will be time-varying.

Mixed parametrized group (1.3) represents gestures which both

- hand posture or hand posture changes over time,
- and hand posture values or changes of position and rotation

are parametrized. Gestures of this group are recognized based on hand posture. Classification presented above relates to the previously described taxonomy proposed by Aigner et al. at follows:

- pointing gestures – represented by mixed parametrized gestures,
- static semaphoric – represented by static action gestures,
- dynamic and stroke semaphoric – represented by dynamic action gestures,
- static iconic – for demonstrating the shapes of objects represented by static action gestures and for presenting the size of object by locally parametrized gestures,
- dynamic iconic – mainly represented by dynamic action gestures, but can also be globally parameterized gestures for indicating the size of the objects,

- pantomimic – represented mostly by mixed dynamic action gestures, assuming that they are always performed in the same way,
- manipulation – represented by all types of parametrized gestures.

5.2 Gesture data representation

5.3 Additional processing steps

Chapter 6

Static gestures recognition

As was already mentioned, the gestures that can be detected can be divided into two main groups: static gestures and dynamic gestures. The static gestures can be understood as a chosen position and orientation of fingers and a hand in a single moment, while dynamic gestures are defined as a movement of the hand and fingers in time. The problem of static gesture recognition is a subject of following chapter. Firstly, the proposed approach is presented, followed by the introduction to the evaluation scheme. In last section, the performed experiments are described, which were used to examine the effectiveness of proposed approach.

6.1 Proposed methods

The static gesture recognition problem can be stated as a problem invariant to time. That means that for each detected hand, the position and orientation can be treated as a new data, uncorrelated to previously classified hands. With this assumption, one can easily generate multiple samples from the sensors in short time, but it also gives an opportunity to look at the static gesture recognition problem as a problem of classification.

In literature, most of the approaches uses two-dimensional representation of gestures, which makes recognition problems relatively simple and therefore allow to use simple classification algorithms e.g. k-NNs [Gesture recognition toolbox]. The three-dimensional representation of data is more complicated to model by the set of features and finally successfully label. While dealing with 3D data, the position and orientation of hand can be easily affected by the height of the hand above the sensor or a small change in the orientation of the hand with respect to the sensor's coordinate system. All of those distortions should not prevent the system from recognizing those gestures as the one as they are similar. To meet those requirement, it is needed to define what is the "small" change in orientation resulting in treating the static gestures as the same one.

In application of the library presented in this thesis, it is assumed that the learning process can be done offline, while strict online requirements has to be met in the recognition part.

To meet listed requirements, the Support Vector Machines (SVMs)[7] were used as a robust classification algorithm. The SVMs were chosen as there exist a solid mathematical background supporting the simple idea of maximizing the margin between classes. Moreover, the SVMs were chosen due to the popularity in the current research trends and is commonly used technique in biology, robotics or IT for solving data classification problems. The another advantage is the existence of the open-source library libSVM [3], which contains the implementation of the SVM in different programming languages. In the original work, SVMs were designed to be only to classify between two classes, but the idea was expanded to utilize the one-vs-all classification allowing

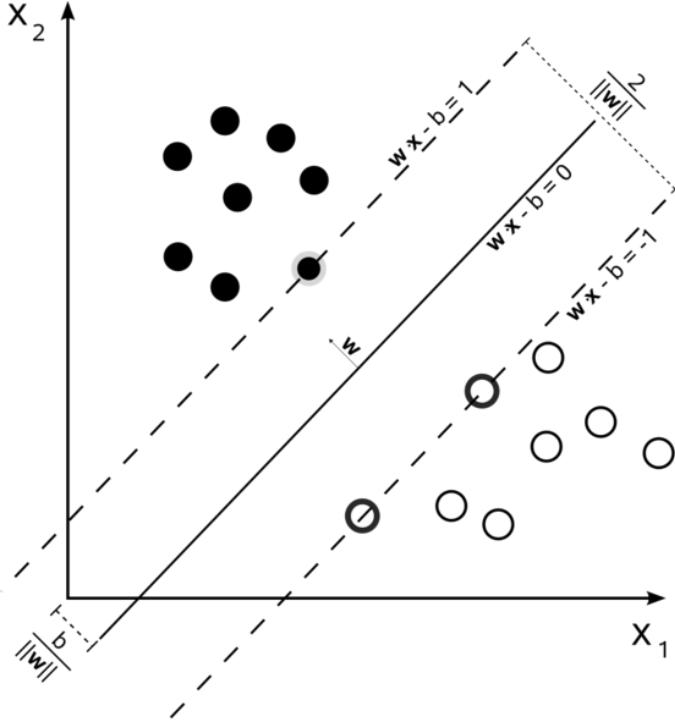


FIGURE 6.1: SVM is a classification algorithm looking for the hyperplane that maximizes the margin between classes¹

to differentiate between multiple classes. This might be the drawback, as the for n classes, it is necessary to train n SVMs each to classify one class against all samples from the rest of the classes. The efficiency of SVMs depends on correctly choosing the kernel function used to map the separation problem into higher dimension with expectation to achieve a classification problem that can be separated. There exists several kernel functions:

- linear: $K(x_i, x_j) = x_i^T x_j$.
- polynomial: $K(x_i, x_j) = (\gamma x_i^T x_j + r)^d, \gamma > 0$.
- radial basis function (RBF): $K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2), \gamma > 0$.
- sigmoid: $K(x_i, x_j) = \tanh(\gamma x_i^T x_j + r)$.

where γ , r , and d are kernel parameters. According to the authors of the library [3], linear kernels perform the best, when used for linearly separable problems, while RBF kernels are the most versatile ones and are recommended for most of the applications.

The problem of classification assumes that each sample consists of set of features, which describe a sample and can be used to distinguish it from the other samples. Additionally, each sample has a known or unknown label, which define the membership of the sample to the classification class. The samples with the known labels can be used to train the classification system. Trained classification system is used to predict the labels for the provided samples.

In application of gesture recognition the classification be divided into two flows: the training part (offline) and the recognition part (online). The static gesture processing flow is presented at fig. 6.2. In training part, the library will be provided with the recorded samples of static gestures with known correspondences to the static gesture classes. From those samples, the sets of features

¹http://en.wikipedia.org/wiki/File:Svm_max_sep_hyperplane_with_margin.png

are computed, which are used to train the classifier. The recognition part assumes to have trained classifier. The recognition part is provided with samples of static gestures without labels. For each sample the sets of features are computed and then given as input to the trained classifier. The classifier provides the information of the predicted gesture's class membership (label) for each input sample.

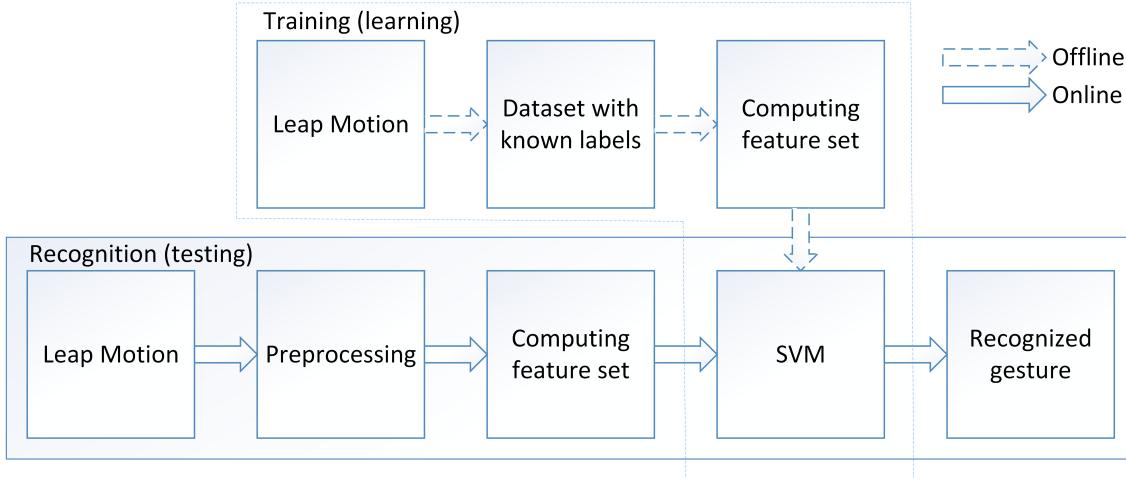


FIGURE 6.2: Proposed solution blocks for learning and recognition parts of static gestures recognition problem

While presented approach can be treated as state-of-the-art approach it still cannot be used without defining proper feature sets for gesture recognition. The naive solution would be to use the raw data from Leap motion sensor as the feature set. This solution was tested, but provided poor results as the proposed features were dependent on the position, orientation and scale of hand. Even small movement in any direction meant problems with stable recognition. The theoretical literature suggests to compute a set of features invariant to wanted transformations, which can allow to fully distinguish between different classes. The proposed feature sets are application dependant. Together the fact that the gesture recognition using the data similar to the data provided by the Leap Motion sensor has not been yet examined by the researches meant that seeking right feature set is a task of the thesis. The task is undertaken in experimental section 6.3.

6.2 Evaluation methodology

6.2.1 Assumptions

To provide user with library working in different conditions, it was assumed that the gesture is treated as the same one independently with respect to the translation, rotation and scale of the hand. This assumption means that the static gesture rotated by unknown angles, translated in sensor coordinate system and also with different hand sizes should still be recognized as the same gesture. Invariance to the rotation, translation and scale poses a great challenge to the recognition, but allows the future users of API to fully utilize the feasibility of the library. It is worth mentioning that, it does not reduce the possible applications of the library, as an assignment of static gesture to already defined class allows to find the transformation between the model of the class and observed gesture. And further, this transformation can be used to examine the parameters of the instance of the model e.g. rotation of the hand.

6.2.2 Recorded datasets

To propose and test the quality of the features twelve static gestures were chosen:

1. the peace sign,
2. a fist,
3. full hand with space between each finger,
4. from American Sign Language: “I love you” sign,
5. sign “gun” created by putting thumb and forefinger up, while holding the rest fingers in a fist,
6. all fingers in a fist with exception of thumb, which is up,
7. the sign X made with the forefingers of both hands,
8. the sign “Time” used e.g. by coaches in basketball games.
9. sign simulating rotating a knob by two fingers,
10. sign simulating rotating a knob by five fingers.

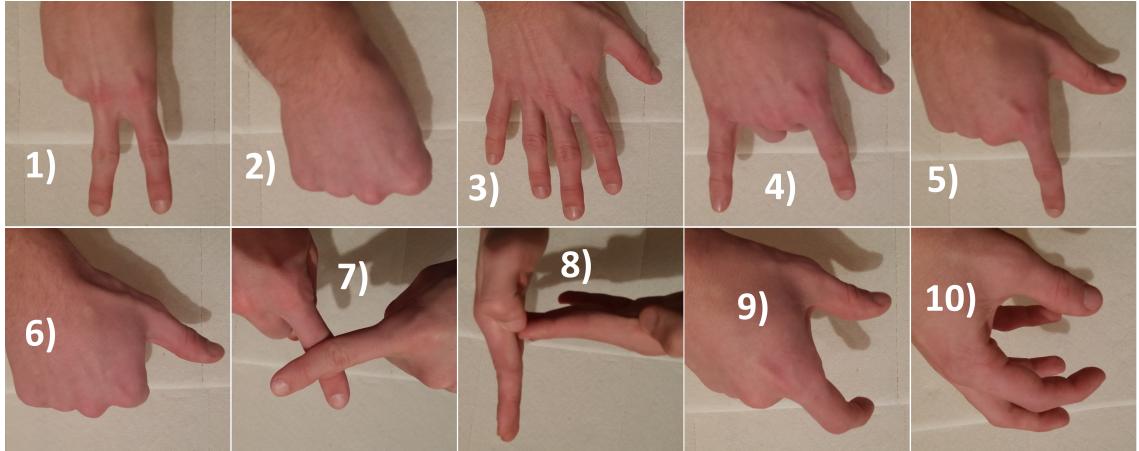


FIGURE 6.3: Recorded static gestures used for evaluation

The gestures are also presented at fig. 6.3.

The sample data of each gestures were recorded using the continuous mode of recording, while moving the hands in different directions and changing the orientation of the hands. The hands were moved and rotated as each repeated static gesture is always a little bit different from the previously recorded samples/ For each of the proposed gestures, each author of this thesis recorded approximately 1000 samples.

Having samples with known labels, the whole dataset was separated into training and testing sets in relation 2 : 1. For the training, the k-fold cross-validation (CV) scheme was used, which searches for optimal C and γ parameters of SVM trying to train the classification algorithm while preventing the algorithm from over-fitting the training data. This method is used to find the optimal parameters of the classification system, while estimating the performance is done on the data not used in the training part. In standard version of the CV, the gathered data is divided into two sets: one containing $k - 1$ parts of the data, the other 1 part of the data. The first is used

to train the classification system, while the rest of the gathered data is used to estimate how well the system is learnt. In the next iteration, the CV parts are shuffled by still having the ratio of training and validation sets. After the training, the final performance is estimated by calculating the percent of correctly recognized labels to the total size of the testing set known as recognition rate.

6.3 Experiments

Firstly, the experiments to find the proper set of features were conducted. The first proposed vector of features consisted of:

- number of fingers in frame,
- the euclidean distance between consecutive finger's tips,
- the absolute angles between consecutive fingers.

This feature set did not take into account the relative position of fingers to the hand. The second feature set is a first feature set extended by the distances between consecutive finger tips and the position of the hand's palm. The third feature set contains features from second feature set extended by the five angles between fingers and normal of hand's palm. Those, 3 proposed feature sets were firstly tested on all 10 recorded. The firstly tested set of all static gestures contained gestures, which were undistinguishable for the Leap Motion, because they did not take into account the way how the Leap Motion works. For gestures like fist or 'X' the recorded data contained almost no information how to classify those gestures. That's why the experiments were repeated on the five gestures, which could be easily distinguish using the data provided by Leap Motion. For this experiment the gestures peace, hand, "I love you", fist with thumb up and rotating knob by 5 fingers were chosen. The results achieved by those methods are presented in the table 6.3.

TABLE 6.1: Results obtained by experimental feature sets by the libSVM library

	5 gestures, CV	5 gestures, test set	10 gestures, CV	10 gestures, test set
feature set 1	87.072%	87.943%	69.987%	68.422%
feature set 2	87.072%	87.943%	69.987%	68.422%
feature set 3	87.072%	87.943%	69.987%	68.422%
feature set 4	81.150%	80.998%	68.438%	68.630%
feature set 5	86.544%	85.075%	77.101%	77.518%
feature set 6	92.762%	93.096%	80.543%	81.235%

For problem of recognition of 5gestures, three first feature sets resulted in over 87% recognition rate on testing sets, which was believed to be a good result. The same tests for problem of recognition of 10 gestures resulted in lower recognition rates. For feature sets 1 – 3 the recognition rate was below 70%, which could be unsatisfying from the perspective of the purpose of the application. The low recognition rate was analysed and revealed that the fingers are numbered accordingly to the position in Z axis of the tip of the finger. This means that when finger's tips are approximately on the same position in Z axis, the numbering can change rapidly and proposed features are compared between different fingers. To achieve features that would be invariant to the numbering of the fingers, the feature set was slightly modified. Instead of containing the absolute angles and distances between consecutive fingers, it was proposed to contain the five greatest values of angles and five greatest values of distances between all combinations of finger pairings. The

same sorting approach was used for the angles and distances between fingers and hand's palm. The feature sets 1, 2, 3 with sorting scheme were respectively called feature sets 4, 5, 6. Again, the same dataset with 5 and 10 gestures was used to evaluate those methods. The results are yet again presented in table 6.3. This approach was tested on the same training set and allowed to increase the recognition rate. The best results in both tasks were achieved in case of feature sets 6. The simple alleviation of finger numbering problem allowed to top the previous results with recognition rates 93.096% for 5 gesture problem and 81.235% for 10 gesture problem.

TABLE 6.2: Results obtained by experimental feature sets by the libLinear library

	5 gestures, CV	5 gestures, test set	10 gestures, CV	10 gestures, test set
feature set 1	78.282%	78.283%	50.596%	50.593%
feature set 2	78.282%	78.283%	50.596%	50.593%
feature set 3	78.282%	78.283%	50.596%	50.593%
feature set 4	78.205%	78.242%	50.658%	50.575%
feature set 5	79.527%	79.502%	55.284%	55.263%
feature set 6	88.076%	88.138%	64.747%	64.830%

While using more data and longer feature sets usually allows to achieve better results, it is worth to notice the growth of training time of classification technique. In case of 5000 training samples the typical training process of radial SVM took approximately 12 hours on standard desktop PC. This computing time can be unacceptable by the users of the library and the tests with another SVM library, libLinear [10] were performed. The libLinear's implementation of SVM utilizes the linear kernels and is dedicated for the large data training sets with multiple number of features. This library reduced the training time to about 5 seconds. Again, the same tests as for the libSVM were performed to compare both approaches. All achieved results are presented in tab. 6.3. For the 5 gesture case, the libLinear achieved 88.138% on testing set while using feature set 6 compared to the 93.096% achieved by libSVM in the same condition. In this case, the libLinear might be good choice as the recognition rate is lower by about 5%. For 10 gesture case, libLinear achieved 64.830% compared to 81.235% by LibSVM, which is a significantly lower recognition rate. It's up to user to decide which library to use, but in most recognition applications, the library can be learnt offline and only used online. That is why, the authors of the thesis recommend the SVM with RBF kernels for static gesture recognition task.

From the results obtained by the libSVM and libLinear on different feature sets, the feature set 6 was chosen as the one yielding the best results and used in further analysis.

Another factor, that might have an influence on the results is the preprocessing part, which should allow to partially remove noise from the data and thus increase the recognition rates. The preprocessing operates in the window of hands poses recorded over time, which size can be modified. The typical library usage allowed to gather data with $60Hz$, while it is assumed that the recognition can be performed with lower framerate. The difference can be efficiently used by defining the appropriate preprocessing window. For these reason, the experiments with no preprocessing and preprocessing with width size equal to 5, 10, 15, 20, 30 were performed and the influence on the recognition rate was examined. The results are presented in tab. 6.3. The achieved results confirm the need and importance of proper data preparation in task of data classification. For gesture recognition task containing 5 gesture poses, preprocessing allowed to increase the recognition rate from 93.096% to over 99% for window sizes equal or wider than 15. In task of correctly detecting 10 gestures, the preprocessing allowed to increase recognition rate from 81.235% to over 84% for windows sizes equal or wider than 10. From those results, the

preprocessing width of 10 was chosen as the one allowing to significantly improve recognition rate in almost any possible application of the library. The preprocessing of width 10 was used for all further experiments if not stated otherwise.

TABLE 6.3: The recognition rate achieved with feature set 6 and different parameters of preprocessing

preprocessing	5 gestures, CV	5 gestures, test set	10 gestures, CV	10 gestures, test set
off	92.762%	93.096%	80.543%	81.235%
width = 2	95.641%	96.527%	82.736%	83.172%
width = 5	95.861%	96.650%	83.078%	83.565%
width = 10	98.796%	98.965%	83.602%	84.071%
width = 15	98.981%	99.150%	84.112%	84.494%
width = 20	99.104%	99.242%	84.329%	84.834%
width = 30	99.232%	99.385%	84.848%	85.320%

All of already presented experiments, treated the classification results of consecutive hand poses as time-independent and not correlated with each other. In real applications, it can safely assumed that the consecutively detected hands are similar to each other and probably define the same gesture. The remaining question to be answered was the impact of combining the consecutive recognition results for the total recognition percentage. Firstly, it is important to use not only the class labels for tested dataset, but the whole information provided by SVM containing the measure of classification rate to all possible classes. This data can be combined to measure the membership rate for each class in a window of set width. Then the predicted class is the class with maximal measure. Formally, when there is a need to classify between k classes, the result for i -th data in dataset can be represented as:

$$l(i) = [l_{i1}, l_{i2}, \dots, l_{ik}] \quad (6.1)$$

where l_{ik} is the likelihood of belonging of i -th data to the k -th class. Then combining in a window of width w can be written as:

$$r(i, w) = \sum_{j=i-w}^i l(j) \quad (6.2)$$

where the sum of vectors l can be defined freely. The recognized label is the number of the vector component with the highest value:

$$\text{label} = \arg \max_k \{r(i, w)_k\} \quad (6.3)$$

The first approach to defining the sum operator is the operator of simple adding the elements of vectors l . The second proposition is to multiple the corresponding elements of vectors l . The third approach utilizes the idea of calculating the median elements of vectors l . The three approaches were compared for feature set 6, with preprocessing width equal to 10 for 5 and 10 gestures already used in previous experiments. Simultaneously, the test were performed for different widths of window and are presented at fig. 6.4. For almost all possible widths, the summation operator demonstrated the best recognition rate. For 5 static gesture, the summation with width equal to 10 allowed to achieve the recognition rate over 99.5% gaining over 0.5% when compared to solution without postprocessing. Interestingly, windows wider than 12 resulted in lower recognition rate than the best achieved with window size equal to 10. For 10 static gesture recognition problem, using window of size equal to 50 allowed to increase the recognition rate by over 5% to the value over 89%. In this case, wider window resulted in better results. Similarly to the preprocessing window size, also in postprocessing too wide window will result in delayed recognition rate of

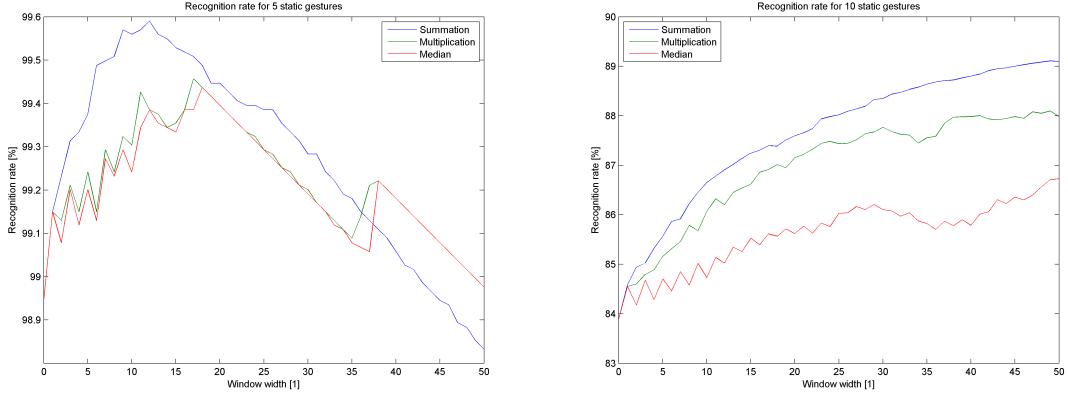


FIGURE 6.4: Evaluation of different operators for result combining in case of different window sizes

shown gesture. Also, too wide window may result in worse results. Therefore, the postprocessing window size of 10-15 is recommended, but it also depends on the application of the library.

With simple summation window, the currently achieved result and the results from the past are equally affecting the recognition result. Especially for wider windows, it can be assumed that the current measurement is more important than the measurement from the distant past. That is why, the idea of weighted sum was introduced. The weight distribution should have the highest weight for the current measurement and smaller values for results that were achieved earlier in time.

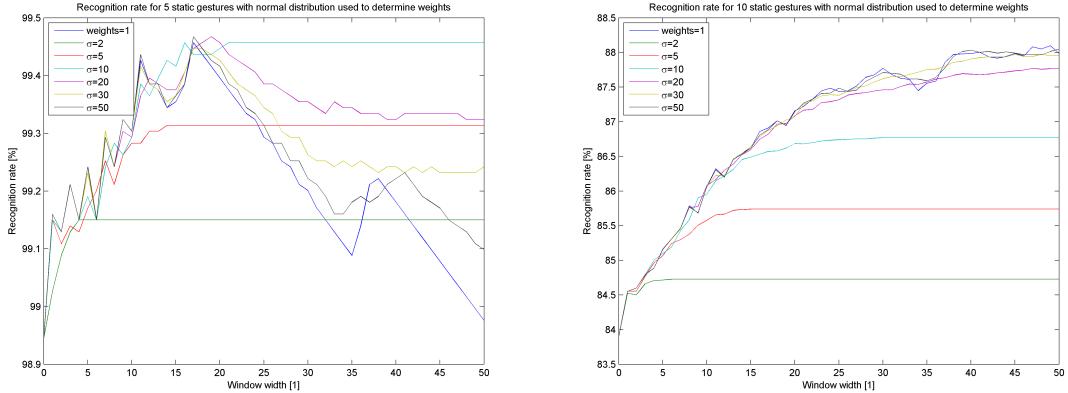


FIGURE 6.5: Evaluation of weights for postprocessing distributed according to normal distribution

For this task, the half of the Gaussian distribution of weights can be used with maximal peak reached for the currently achieved result with weights slowly decreasing for measurements further away in time. For Gaussian, the mean was assumed to be equal to 0. The standard deviation σ was the parameter, which different values were tested. The results were also performed for problems of recognition of 5 and 10 gestures. The achieved results are presented in fig. 6.5. For 5 gestures, too small σ prevented the increase of recognition rate due to the wider window size as the weights were equal to 0. For greater sigmas, the achieved recognition rate is similar. For this problem, the σ equal to 10 resulted in best recognition rate while the window was getting wider. For problem of 10 static gesture recognition problem, greater sigma values achieved results comparable with the the results achieved with weights equal to 1. As the usage of different weights also comes with

greater computing cost and still did not result in better results, this part of processing can omitted without effecting the recognition.

From the presented results, the simple summation with uniform weights provides the greatest gain when it comes to the recognition rate and is the method recommended by the authors.

6.4 Finger differentiation

Finger differentiation module is responsible to verify which particular fingers are detected during gesture recognition. The module can be used for additional processing on the user side, for which the information about the visible fingers can often be significant, It can also assist interpretation of data derived from parametrized gestures. The user can also use the module to better fit dataset features to the characteristics of selected gestures.

6.4.1 Methods

For this module were used the same classification algorithms (kernel functions) as those who had the best results in static gesture recognition. As in the case of static gesture recognition libSVM library has been used and RBF kernel has been selected.

6.4.2 Evaluation methodology

In the assumptions for static gesture recognitions was mentioned that the process must be independent of the position, rotation, sizes of hand and fingers. This presumption applies also to finger differentiation. However, this assumption is still does not contain one element, from which the differentiation process should be independent – arrangement of hand. Whether forefinger is straight or bent, it is the same class, where only this one finger is taken out.

Recorded dataset

Dataset was collected by two different people. Each of them has recorded 32 classes, which reflect all the possible permutations of hand arrangements. The next step was to select from 32 gestures, those which are most frequently used by people in everyday life. Those 15 selected gestures are natural for human and do not cause a pain. People uses them as specific signes like peace sign. From the dataset were created four data collections used in the experiments:

- 32 gestures recorded by two people,
- 15 gestures recorded by two people,
- 32 gestures recorded by one person,
- 15 gestures recorded by one person.

6.4.3 Features

As mentioned previously, it is important that the fingers differentiation process is independent of things like position of hand or size of finger.

6 features were selected, which are independent of those parameters:

1. finger count,

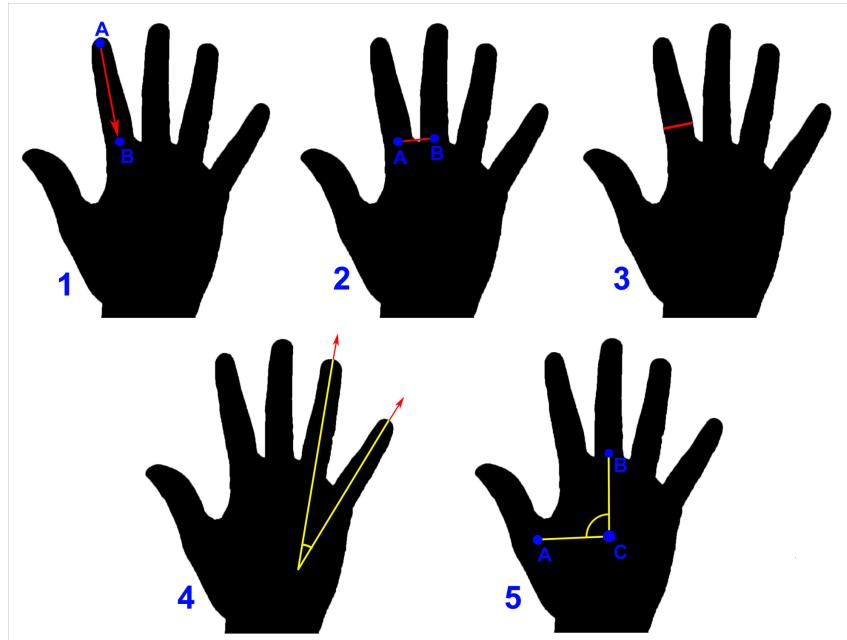


FIGURE 6.6: Features used for finger differentiation

2. distance between two nearest base points of a finger, To appoint a finger base point has to be reversed normalized direction vector multiplied by the length of a finger. Subsequent as the beginning of the vector thus formed is set finger tip position. The end of the vector indicates the finger base point (fig. 6.6.1). An example of distance between two nearest base points of a finger is showed on fig. 6.6.2.
3. ratio of the finger thickness to the maximal finger thickness, Leap Motion controller can return information about the thickness of the fingers (fig. 6.6.3). There are known relationships between the thicknesses of the fingers, for example that the thumb is the thickest finger. Just like in the previous feature ratio was used in order to become independent of size.
4. angles between two nearest fingers, In fig. 6.6.4 this feature has been shown. It is obtained by calculating the angle between the direction vectors of two adjacent fingers.
5. angles between finger and first finger relative to palm position. To calculate this angle three positions are needed : two finger tip positions and a palm position. The next step is to determine the line segments between palm position and finger tip positions. The searched angle is between two designated segments. This angle is marked in fig. 6.6.5.

During testing new feature was created based on the features 2:

- 2b. ratio of distance between two nearest base points of a finger to the minimal (non-zero) distance between two nearest base points This feature is very similar to the previous one. The difference is that in this case is used ratio of distance and not the actual values. Thanks to this feature is independent of the size of hand.

6.5 Experiments

For the tests created several combinations of previously described features:

- 1,2,4

- 1,3,5
- 1,3,4,5
- 1,2,3,4,5

All sets of features has been tested for all four datasets. The percentage results can be seen in table 6.5. From the data obtained it can be concluded that the best match is when the all features are used. However, feature 2 is not independent of the size of the hand. Therefore, a new feature has been created using the ratio of two values obtained from the original feature. Then the test was performed for a set of: 1, 2b, 3, 4, 5 for all the dataset. The result for this test is in the last line of the table 6.5. Interestingly, despite the fact that the modified feature is independent of the size of hand, attribute 2b had an inferior result than the original feature, which operates on the actual value of the distance apart fingers.

TABLE 6.4: Results obtained by experimental feature sets for different data collections

features	32 classes, 2 people	15 classes, 2 people	32 classes, 1 person	15 classes, 1 person
1,2,4	75.393%	87.288%	78.411%	90.608%
1,3,5	76.291%	87.504%	79.638%	91.548%
1,3,4,5	78.082%	87.504%	79.638%	91.548%
1,2,3,4,5	83.180%	90.278%	85.825%	93.036%
1,2b,3,4,5	76.553%	87.619%	79.798%	91.629%

It is worth to mention that it seemed that feature 4 – angles between two nearest fingers – will have a substantial contribution to finger differentiation. However, the results showed that this feature has impact only on the largest amount of data, while in other cases brought nothing.

In the second round of tests was verified how preprocessing will affect the finger differentiation process. For the best set of features from the previous test was performed a new test with preprocessing for all datasets.

TABLE 6.5: The recognition rate achieved with feature set 1,2,3,4,5

features	32 classes, 2 people	15 classes, 2 people	32 classes, 1 person	15 classes, 1 person
1,2,3,4,5	%	%	%	%

Chapter 7

Detection of dynamic gestures

7.1 Proposed methods

The dynamic gesture recognition problem is a problem, where the input data consist of several consecutive positions and orientations of hand and fingers. Moreover, the important factor for recognition is the time dependencies between data frames. The slower and faster gestures should be recognized as the same dynamic gesture.

The proposed solution utilizes parts of the solution used for recognition of the static gestures. Each frame of the captured data is described by the same features as in the static recognition part. The set of features for each frame is then processed by the Hidden Markov Model scheme.

7.1.1 Hidden Markov Model

Hidden Markov Model is model of a system with Markov property. The first introduction of HMM comes from the L. E. Baum et. al.^[1], who proposed the mathematical background for HMMs. A HMM can be considered a finite, N -element set of states, which are associated with the probability distribution. The transitions between states are represented by the transition probabilities usually stored in $N \times N$ matrix T . In every state, one of the observation from the finite, K -element observation set can be generated with observation probability usually represented by the $N \times K$ emission matrix E . The finite set of all possible observation is called the alphabet. To fully define the model, the definition of N -element vector of initial state probabilities Π . Each HMM can be fully defined by the (T, E, Π) .

The example HMM can be seen at fig. ???. The HMM at the figure consists of 3 states $X1, X2, X3$ and 4 possible observations $y1, y2, y3, y4$. The probabilities a_{ij} define the transition probability from state j -th to state i -th. The probabilities b_{ij} define the observation probability of generating j observation while being in state i .

The HMM can be understood as a directed graph with two types of weighted vertices. This way, each state is represented by one type of vertices while observations can be shown as second type of vertices. The edges between states contain and are an equivalent to the transition matrix. There are no edges between vertices representing observations. The edges between states and observations are equal to the observation matrix.

xx There are three main algorithms used with the HMMs:

- Forward-Backward algorithm,
- Viterbi algorithm,

¹<http://en.wikipedia.org/wiki/File:HiddenMarkovModel.svg>

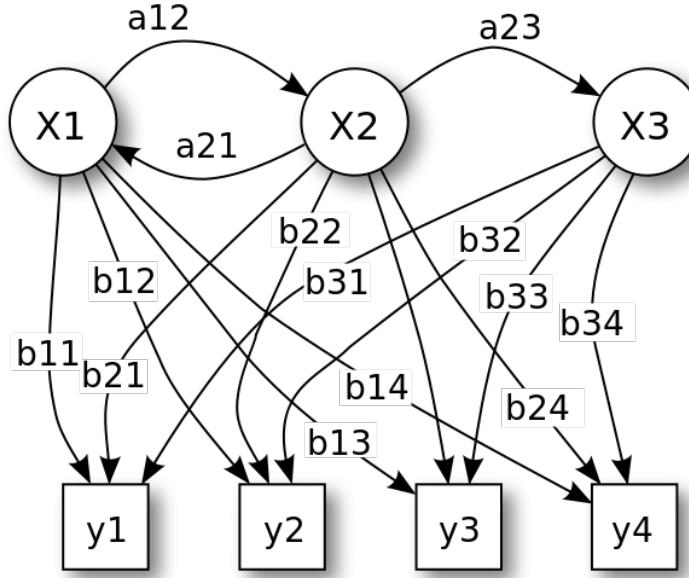


FIGURE 7.1: Solution blocks of learning and testing parts in task of dynamic gesture recognition¹

- Baum-Welch algorithm.

The Forward-Backward algorithm is used to find the posterior probability of given states given the set of observations. For the whole set of state variables $X = X_{i=1}^N$ and a set of observations $o_{1:t} = o_1, o_2, \dots, o_t$, the algorithm computes the $P(X_i|o_{1:t})$. The algorithm utilizes a dynamic programming approach, by performing 3 steps in a loop:

1. computing forward probabilities,
2. computing backward probabilities,
3. computing smoothed values.

Firstly, the forward pass phase for $i = 1, \dots, N$ computes the $P(X_i, o_{1:k})$, where k is smaller than t , which represent the probability of ending up in state X_i after first k observations. The backward pass computes the $(P(o_{k+1:t})|X_i)$, which are the probabilities of observing the rest of the observations in state X_i . The smoothing part, uses Bayes rule to compute the probability of state X_i given the whole observation sequence:

$$P(X_k|o_{1:t}) = P(X_k|o_{1:k}, o_{k+1:t}) \propto P(o_{k+1:t}|X_k)P(X_k|o_{1:k}) \quad (7.1)$$

The time complexity of this algorithm is $O(N^2T)$, where T is the length of observation sequence and N is the number of possible states.

The Viterbi algorithm is used to find the most likely sequence of hidden states that best explain the set of observations. The set of those states is usually called the Viterbi path. Introducing the sequence:

$$V_{1,k} = P(o_1|k) * \Pi_k \quad (7.2)$$

$$V_{t,k} = P(o_t|k) * \text{argmax}_{x \in X}(a_{x,k} * V_{t-1,x}) \quad (7.3)$$

Finally, the problem can be understood as finding the $\text{argmax}_{x \in X}(V_{T,x})$. The algorithm time complexity is the same as forward-backward $O(N^2T)$.

The Baum-Welch algorithm is the algorithm used to find the unknown parameters of the HMM. For each training example, the algorithms updates the state transition probabilities, observation probabilities in each state and initial probabilities to maximize the likelihood of observation. Having the set of sets of observations, the algorithm can be used to train HMM to detect the sequences similar to the ones used in learning process.

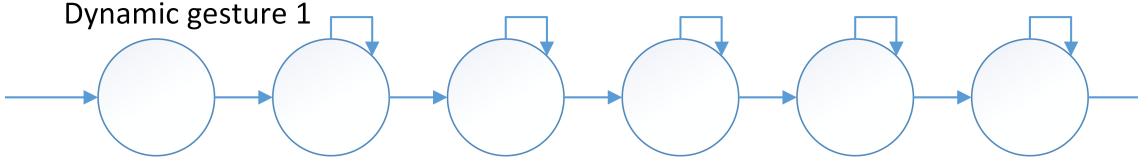


FIGURE 7.2: Non-zero state transitions and states of the HMM used to detect single dynamic gesture

In the dynamic gesture recognition task, we adopted a structure of HMM with state having non-zero transition probabilities to self and to the next state in the sequence. The proposed structure is presented at fig. 7.2 and was firstly proposed in paper [42]. The states can be understood as the phases of hand movement and position that happen when the wanted gesture is performed. The self-transitions are used to model the different speeds of the gestures. This structure after training process can be used to measure the probability of that the dynamic gesture occurred given the set of observations.

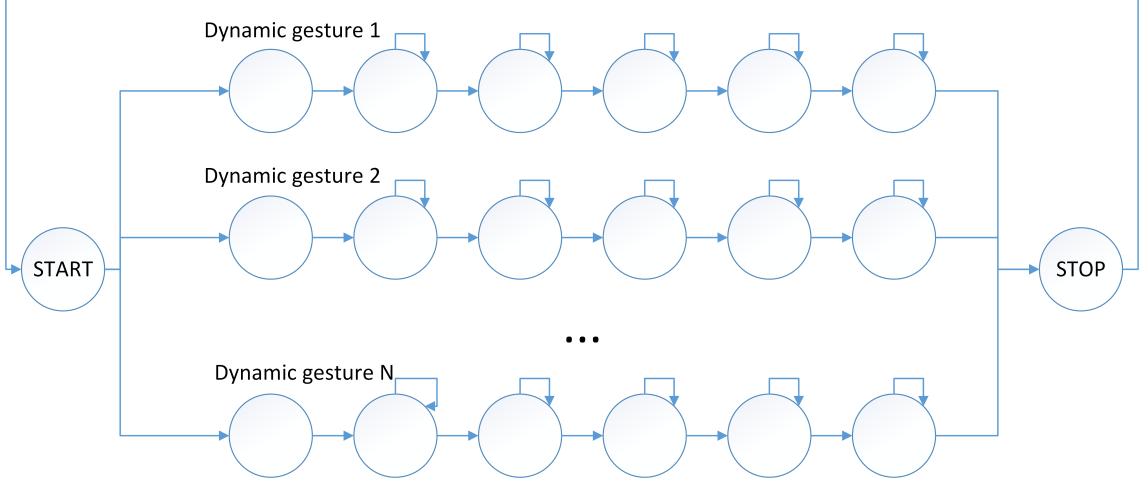


FIGURE 7.3: Non-zero state transitions and states for the HMM enabling simultaneous detection of m dynamic gestures

Having single dynamic gesture recognition problem modelled as the sequence of n states in which k th state is connected by the edges to the k and $(k+1)$ state and to all observations. Having the problem of distinguishing m gestures translates to the m sequential graphs. The problem of finding if and what dynamic gesture occurred is the problem of finding the probabilities from each m HMM and checking them against preset threshold. Alternatively, the m HMMs can be combined into one HMM where those single HMMs are treated as parallel paths. The structure

of proposed single HMM is presented at fig. 7.3. Then the recognition process is a process of finding, which of the parallel paths is the most probable. The only remaining problem is the problem of preparing the set of cured, 1-dimensional observations from the Leap Motion data.

7.1.2 HMM observation from Leap Motion data

Most of the proposed solutions in the literature assume that the observations are 1-dimensional. There were attempts to expand the HMM to the 2-dimensional observations [], but the raw data from Leap Motion is defined in higher-dimensional than any found solution. That is the reason, why the need to reduce the dimensionality of the observation of the single hand in the predefined moment in time.

To do this, unsupervised clustering algorithms were introduced. For this task we examined 3 methods:

- Chinese whispers [2, 6],
- Girvan-Newman clustering [12],
- k-means with additional methods to determine the number of clusters. [36, 24].

Chinese whispers algorithm is an efficient randomized graph-clustering algorithm, which is time-linear in the number of edges. The algorithm was proposed in [2] and thoroughly examined in Natural Language Processing problems. The main advantage of the algorithm is the ability to independently determine the number of classes. In tests, we used the implementation available in dLib-ml library[21], which provides multiple machine learning algorithms.

Girvan-Newman clustering is a hierarchical algorithm used to detect clusters in a data represented as a graph. The algorithm progressively removes edges from the original network until it reaches the maximal number of iterations or the error condition is met. The solution iteratively calculates the eigenvalues of a created matrix using the power iteration method. In case of a complete graph, this algorithm is relatively slow.

The another proposed approach is a k-means clustering algorithm. The idea for the algorithm comes from polish mathematician Hugo Steinhaus [36], but the first practical application has been presented in [24]. The algorithms initially randomly selects k starting points and performs iteratively 2 steps utilizing the idea of expected-maximization [8]:

1. For each data sample the algorithm calculates the distances to k centroids and labels the samples with the label of the centroid with the smallest distance.
2. For each class, the algorithm recalculates centroids position to the averaged position of all samples belonging to class.

The algorithm stops after maximum number of iterations or when the change between consecutive iterations is smaller then defined epsilon. Unfortunately, the algorithm requires the knowledge of the number of expected classes. Although there exists a heuristics methods aiding in correctly choosing the number of classes. One of those methods is plotting the error sum of squares (SSE) within classes against the number of chosen classes. In this plot the drastic drop of SSE is sought and the number of classes it happens is considered a good choice. Another, more formal approach introduces the measure of dissimilarity and is called Silhouette width [33]. Declaring $a(i)$ as the average dissimilarity with all the other data within the same cluster and $b(i)$ as the lowest average

dissimilarity to any other cluster which i is not a part of, we can write a measure:

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \quad (7.4)$$

The big value of $s(i)$ implies good clustering, while small value of $s(i)$ suggests that the number of classes was not properly chosen.

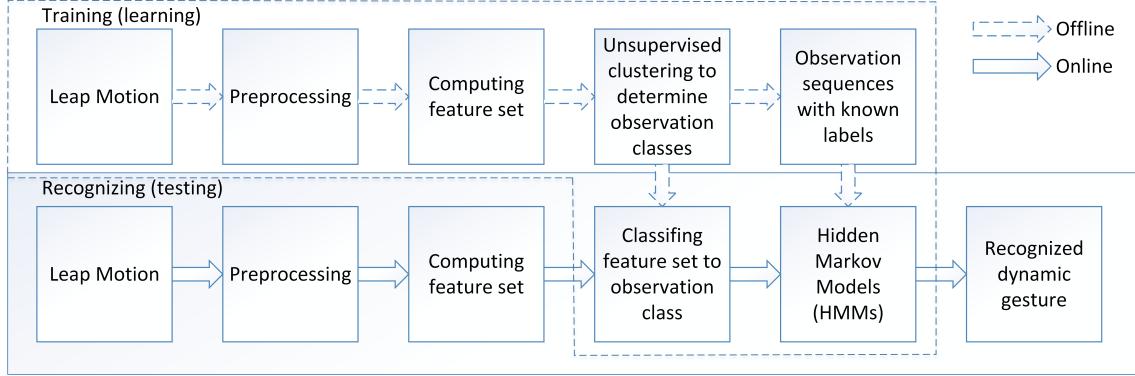


FIGURE 7.4: Solution blocks of learning and testing parts in task of dynamic gesture recognition

Finally, the whole processing flow has been designed and is presented at fig. 7.4. The whole solution consists of two parts: offline learning and online recognition. For learning, the raw data is preprocessed and then the feature are extracted. Similarly to the static gestures, the features set were computed. In training part, those features are extracted from all recorded positions for all dynamic gestures. Then one of the unsupervised clustering algorithms is used to find the classes of observations. Knowledge of the classes is used to represent each dynamic gesture as a series of one-dimensional observations. After that for each dynamic gesture, the corresponding HMM is learned by running the Baum-Welch algorithm on the sequence of observations. The Baum-Welch algorithm finds the matrices that maximizes the likelihood for one sequence of observations, therefore the problem of training on multiple training data was encountered. Following by the idea presented in [27], each matrix trained by one training example is incorporated into trained matrix by using simple addition with learning rate μ . Denoting the transition matrix by T and learnt transition matrix T_{new} , it can be written:

$$T = (1 - \mu) \times T + \mu \times T_{new} \quad (7.5)$$

The same idea applies also to the emission matrix, denoted by E and learnt emission matrix E_{new} :

$$E = (1 - \mu) \times E + \mu \times E_{new} \quad (7.6)$$

The problem with the matrices by computing equations ?? is the fact, that they do not fulfill the probability constraints. Therefore, each row of new matrices is normalized to 1.0. When dealing with multiple training examples, the solution may be overfitted to the training samples. To deal with this problem, the K-fold cross validation was utilized in a similar approach as in the static gestures. To complete the training process, each of the HMMs is put into one recognition structure. During the recognition part, observed sequence is treated as input to the Viterbi algorithm performed on each of the HMMs and the output is analysed. The HMM with the highest likelihood is chosen, and the observed sequence is labelled with the label of this HMM.

In the online working, the raw data from leap motion is preprocessed, then each frame is labelled accordingly to the classes learnt by the unsupervised clustering algorithm. The next step

include providing the set of observations to HMM and running the Viterbi algorithm. The Viterbi algorithm finds the most likely sequence in one the parallel paths, which informs us about the number of the found gesture. If the likelihood is above preset threshold, the gesture is assumed to be correctly recognized.

7.2 Evaluation methodology

To evaluate the quality of recognition achieved by proposed processing pipeline, 6 dynamic gesture were chosen for which the data was recorded. For each of those gestures, we recorded 120 samples (30 samples per each author). The recorded data were in different positions, with different speed to measure the wanted invariance. The chosen dynamic gestures are:

1. “123” – it’s a gesture when performing counting to 3 using hands,
2. “door” – it’s a gesture performed while trying to grasp a handle of door and then open the door,
3. “circle” – making a circle in the air with the forefinger,
4. “scissors” – simulating the cutting with the scissors made by a forefinger and a middlefinger,
5. “gun” – only the thumb and the index finger are not in the fist. The hand moves up simulating the movement during the firing from the gun.
6. “moving the object” – performing the task of grasping an invisible object, moving it and letting it go in different place.

To test the proposed approach, similarly to the static gesture recognition problem, the $120 \times 6 = 720$ samples were divided into training set containing 66,7% of recordings and testing group containing the remaining 240 recordings. The part of the gestures recorded by Katarzyna were used to find the proper number of clusters than could be used by form the observations for the HMMs. The training set was used to train each HMM separately on the training data of each gesture. To prevent overfitting to the data and achieving the measure when the training process should be stopped, the k-fold cross validation was utilized. Preliminary results revealed that initialization of HMM matrices has an impact on the achieved results. Due to no prior knowledge, the random initialization was chosen. As this was susceptible to not yield reliable solution, each training process is performed 10 times and the best model from cross-validation is returned. Each training cycle of all HMMs takes between 1 – 10 minutes, with total training part taking around 25 – 45 minutes. If not stated otherwise, the learning rate μ was set to 0.1, k-cross validation was performed for $k = 5$ and the number of states in one HMM was set to 10.

7.3 Experiments

The performed experiments started with testing the unsupervised clustering methods to determine the correct number of observations. Based on the successful static gesture recognition, each of the hand poses was represented by the vector of values containing:

- number of detected fingers,
- 5 greatest angles between the finger tip vector and palm normal,
- 5 greatest angles between the fingers tip vectors,

- 5 greatest distances between the tip positions of fingers.

In order to compare different hand poses, the distance function was introduced as the L2 norm between feature vectors:

$$d(x, y) = \sqrt{\sum_{i=1}^{16} (x_i - y_i)^2} \quad (7.7)$$

The Chinese whispers and Girven-Newman uses the similarity function, which was defined as:

$$s(x, y) = \frac{1.0}{\max \{d(x, y), \text{eps}\}} \quad (7.8)$$

where eps was the numerical epsilon.

For the Chinese Whispers and Girven-Newman clustering the typical parameters from dlib-ml library were used — the Chinese Whispers maximal iterations were set to 100, while the Girven-Newman was run with maximal iterations equal to 2000 and precision eps set to $1e-4$. For the SSE analysis for different cluster number, the publicly available script was used[26]. The script automatically calculates the k-means clustering algorithm for chosen range of k values and presents the figures, which can be a help in determining the correct number of classes. For the Silhouette width analysis own, R script was written.

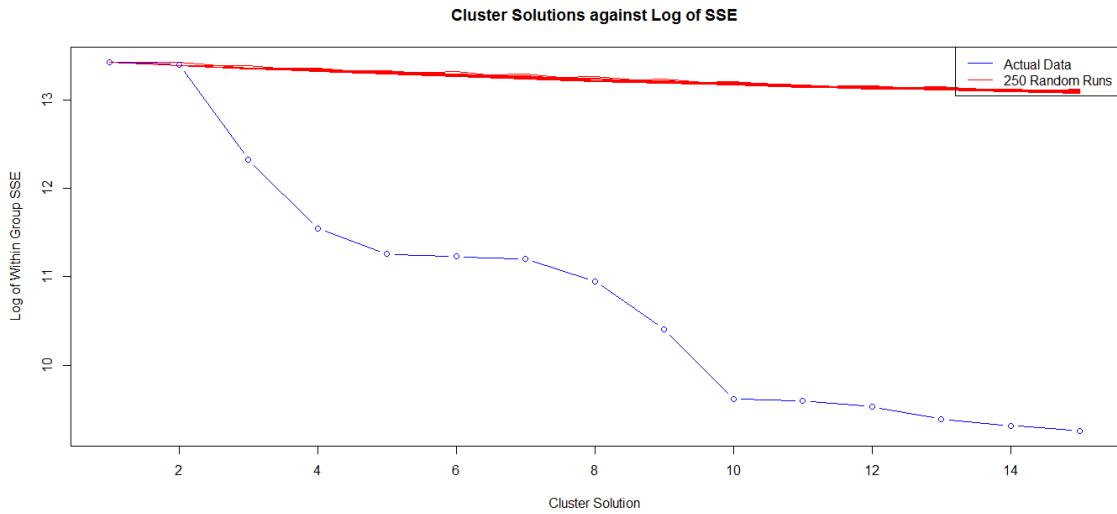


FIGURE 7.5: Different representations of SSE as a function of number of clusters for the k-means algorithm

The SSE analysis usually does not yield the direct answer. The toll provides several plots, but the most interesting one is presented at fig. 7.5. Looking at this plot, there is a significant drop in SSE for the cluster number equal to 5 and 10. It is expected that either of those numbers is the correct choice when it comes to the number of classes for k-means algorithm.

TABLE 7.1: Comparison of suggested number of clusters for the dataset containing all positions of hand in all dynamic gesture recordings made by Katarzyna

Girven-Newman recordings by Katarzyna	Chinese whispers	SSE vs clusters	Silhouette width
x	xx	5 or 10	xx

The proposed values of clusters for different methods are presented in the table 7.1. Due to the inconclusive results, the correct number of classes was determined by performing experiments

with different class number parameter. The analysis of the correct class number was postponed until the correct feature set is chosen. Until stated otherwise, the class number equal to 10 was chosen.

The next challenge was to find the feature set describing one gesture in time, that would contain the information about the dynamic nature of the gesture to be recognized. The first feature set consisted of features encoding the information about the speed of the hand. For each i -th hand position in recorded sequence, the $(i - 10)$ -th position was also considered. The feature set consisted of the recorded displacement of hand in X, Y, Z directions. The number of fingers in both positions were also added, resulting in feature set containing:

- finger count in i -th hand position,
- finger count in $i - 10$ -th hand position,
- displacement of hand position in X ,
- displacement of hand position in Y ,
- displacement of hand position in Z .

The presented approach allowed to achieve 64.583% on cross-validation set. Due to the small testing set, the total recognition rate was calculated on training and testing sets. This approach resulted in 61.806% on all samples, which is too small to be used in real applications.

The previous approach suffers from the fact, that the encoded speeds are represented in the coordinate system of the Leap Motion. Executing the gesture with different angle to the Leap Motion's coordinate system makes the recognition tasks harder and sometimes results in mislabelling. Therefore, the feature set was supposed to encode the information with respect to the local coordinate system of hand. In new approach, the magnitude of the displacement is calculated. It is independent of the chosen orientation of the coordinate system. To represent the direction of movement, the normalized displacement vector is computed as a vector connecting the center of $i - 10$ -th hand position with i -th hand position. Then, the dot product between normalized displacement vector and palm normal in $i - 10$ position is calculated. As this not encode the whole information, the dot product between the normalized displacement vector and the plan direction is also computed. The resulting feature is contains also 5 elements:

- finger count in i -th hand position,
- finger count in $i - 10$ -th hand position,
- magnitude of the hand's displacement,
- dot product of the normalized displacement vector and the palm's normal in $i - 10$ -th position,
- dot product of the normalized displacement vector and the palm's direction in $i - 10$ -th position.

This approach allowed to improve the recognition rate in cross-validation () and on the whole set (), although was still not satisfying from the practical point of view.

The simple observation on the feature set and recorded gestures reveals that sometimes only the fingers are changing positions, while position of hand is steady. The previous approaches, did not encode this information in feature vectors, which made the those changes invisible to the further processing pipeline. Therefore, the new feature vector contained also the information about the

displacement of each corresponding finger. In order to keep the feature size small, it was decided to only incorporate the magnitude of those displacements. Due to the nature of Leap Motion's erratic numbering of fingers, those magnitudes were also sorted. The total of 4 top values were added to feature vector resulting in:

- finger count in i -th hand position,
- finger count in $i - 10$ -th hand position,
- magnitude of the hand's displacement,
- dot product of the normalized displacement vector and the palm's normal in $i - 10$ -th position,
- dot product of the normalized displacement vector and the palm's direction in $i - 10$ -th position,
- 4 greatest magnitudes of finger displacements.

The results were not satisfying as were even worse than the ones achieved for the feature set 2 – 66.111% compared to the previously received 67.500%. The further inspection revealed that, while adding displacement of fingers helps in most situations, it can also pose a great risk. In case of different order of finger numbering, the resulting displacement may be calculated between different fingers, thus resulting in big artificial displacement of fingers. Therefore, the idea of using displacement of fingers was abandoned.

Looking at the recorded gestures, it seemed that speed information is not sufficient to correctly detect the sequences of gestures combining in one dynamic gesture. The new feature set, contained also the information about the static hand positions allowing to distinguish dynamic gestures that are comprised of sequential static gestures. The new feature set contained the already proposed speed part with additional static encoding that was successfully used in static gesture recognition:

- finger count in i -th hand position,
- finger count in $i - 10$ -th hand position,
- magnitude of the hand's displacement,
- dot product of the normalized displacement vector and the palm's normal in $i - 10$ -th position,
- dot product of the normalized displacement vector and the palm's direction in $i - 10$ -th position,
- 4 greatest euclidean distances between all combination of finger's tips in i -th position,
- 4 greatest absolute angles between all combination of finger's vectors in i -th position,
- 4 greatest euclidean distances between finger's tips and palm's position in i -th position,
- 4 greatest absolute angles between finger's vectors and palm's normal in i -th position.

The experiments showed that the proposed approach is not as successful as it was believed. It was believed that the sophisticated static part of feature set is dominating dynamic part of feature set, which prevents the method from achieving better results. Therefore, the 5-th proposed feature set consist of static part represented by the 4 greatest euclidean distances between all combination of finger's tips in i -th position and 4 greatest absolute angles between all combination of finger's vectors in i -th position. Similarly, feature sets 6, 7, 8 were also the propositions containing the

TABLE 7.2: Results obtained by the HMM approach using different feature sets, 10 classes

	5 gestures, CV	5 gestures, whole dataset
feature set 1	64.583%	61.806%
feature set 2	69.583%	67.500%
feature set 3	68.542%	66.111%
feature set 4, full info	77.917%	76.111%
feature set 5	66.667%	63.889%
feature set 6, only dist	72.912%	70.417%
feature set 7, only angles	71.875%	69.167%
feature set 8, dist fing, ang finpalm	78.125%	77.222%

reduced static parts, but did not improve significantly over already received results. The all results are presented in tab. 7.3.

The best results were obtained for feature set 8, but the results obtained by the full static representation (feature set 4) are not significantly worse. For further assessment, the feature set 4 was chosen as the one that describes each gesture with more information and therefore is believed to work better for gestures that were not included in proposed experiments. The mentioned feature set will be used in further evaluation.

The next experiments were performed to test what is the best number of observation classes for k-means algorithm. The tests were conducted using feature set 4 and learning rate $\mu = 0.05$.

TABLE 7.3: Results obtained by the HMM approach using feature set 4 and different class clusters

	5 gestures, CV	5 gestures, whole dataset
$k = 4$	77.292%	74.722%
$k = 5$	79.583%	77.639%
$k = 6$	79.583%	77.639%
$k = 8$	77.708%	75.977%
$k = 10$	77.917%	76.111%
$k = 12$	78.542%	77.361%

The best obtained results are for the $k = 5$ and $k = 6$ allowing to increase the recognition rate from cross-validation to 79.583%, which allowed to achieve recognition rate equal to 77.639% on the whole dataset.

The next experiments involved finding the best learning rate. In our application, the stable learning rate was chosen to minimize the number of parameters to tune. The further developments may involve using learning rate that is changing depending on the already received learning rate. Using predefined learning rate may be dangerous as small value may mean small convergence, while too big step may result in not converging to the locally best solution. The achieved results are presented in 7.3. In proposed application, the learning rate $\mu = 0.05$ allowed to achieve the best results.

TABLE 7.4: Results obtained by the HMM approach using feature set 4, 5 observation classes and different learning rates

	5 gestures, CV	5 gestures, whole dataset
$\mu = 0.01$	77.500%	76.944%
$\mu = 0.05$	80.417%	79.028%
$\mu = 0.1$	79.583%	77.639%
$\mu = 0.2$	69.583%	69.7222%

The only parameter not tested in previous experiments is the number of state each HMM is composed of. The previously used value was equal to 10 states. The performed experiments confirmed the choice as smaller number of states (5) resulted in lower recognition rate, while greater number of states did not improve the result. The drawback of taking HMMs with big state numbers is the learning time. The performed experiments and theoretical computational complexity confirm that the computational complexity grows linear with the chosen number of states, which in practical applications means linear grow of taken processing time.

TABLE 7.5: Results obtained by the HMM approach using feature set 4, 6 observation classes and different state elements

	5 gestures, CV	5 gestures, whole dataset
K = 5	73.958%	73.611%
K = 10	79.583%	77.639%
K = 20	77.500%	76.528%
K = 30	77.917%	76.25%

Chapter 8

LMGesture library dedicated for Leap Motion controller

8.1 Architecture

8.1.1 Model

While data from the LMR files are being processed, they are stored using a specially created class representing the data. GestureFrame represents a single frame obtained from Leap Motion Controller. All gathered data is stored in a vector containing elements of GestureFrame type. GestureFrame holds information such as:

- timestamp,
- list of hands occurring in the frame, stored in a vector containing elements GestureHand type.

GestureHand stores parameters of hand performing gesture. In one instance of GestureFrame can be stored many instances of GestureHand. GestureHand holds following information:

- hand id,
- plam position,
- stabilized palm position,
- palm normal vector,
- palm direction,
- list of fingers of particular hand, stored in a vector containing elements GestureFinger type,
- ordered value, obtained during hand sorting.

GestureFinger stores parameters of one finger. In one instance of GestureHand can be stored many instances of GestureFinger. GestureFinger contains:

- finger id,
- tip position,
- stabilized tip position,

- finger direction,
- finger length,
- finger width,
- ordered value, obtained during finger sorting.

8.2 Processes

8.2.1 The learning process

The learning process is a process, which shows how user can teach LMGesture library a new gesture. First gestures must be recorded gestures in LMR format. To obtain recordings in this format user may use Gesture Recorder – a module included in the library, which records data from Leap Motion Controller and saves it as LMR file. When all desired gestures are prepared, the user starts the process of learning by [tutaj dodać informacje jak uruchomić ten proces - prawdopodobnie jakas komenda + jakieś parametry typu -d uczenie gestów dynamicznych, -s uczenie gestów statycznych + pliki LMR jako argumenty].

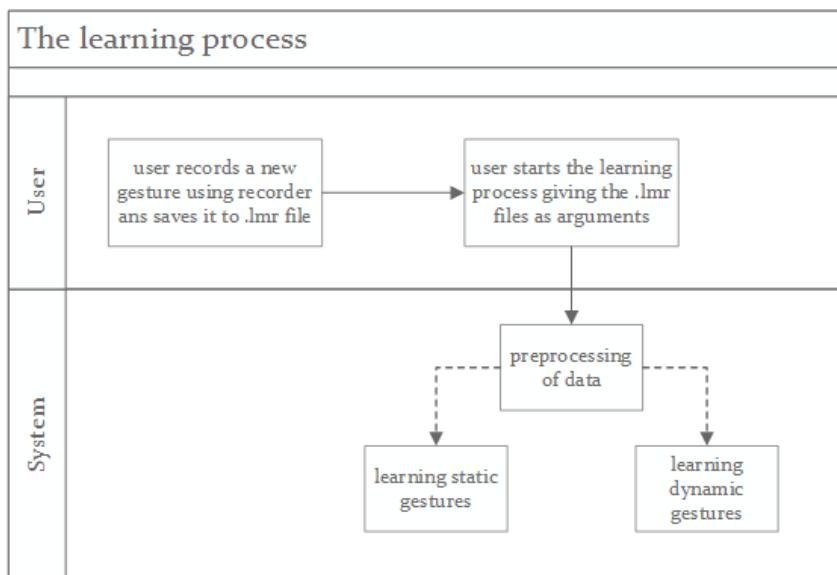


FIGURE 8.1: Diagram showing the learning process

The first step in the learning process, which is performed by the system is the preprocessing of data. Training data should be deprived of noise or should not have lost fingers for several frames. According to user-specified parameters, the system will learn the gesture as a static or dynamic gesture. Appropriate way of learning corresponds to the appropriate gesture. For the learning static gesture system uses a support vector machine (SVM), and for dynamic gestures uses hidden Markov model (HMM). The learning process for static gestures:

1. In the first stage, data which will be used during the learning is preprocessed.
2. Data are scaled using internal scaling module. The result of this step are scale and range files. In the scale file scaled data is stored. The range file contains information that enables to scale any data in the same way.

3. Scale file with scaled data is transmitted to training module. Data is processed using SVM algorithm. Additionally during learning process k-validation is performed and results are returned to the user. In process model file is created, which is used during static gesture recognition process.

[opisać learning process dla dynamicznych]

8.2.2 The recognition process

The recognition process is a process in which the system attempts to match gesture given by a user, to the set of gestures obtained by the learning process. [Opis procesu recognition process, przedstawienie scenariusza uzycia w stosunku do architektury]

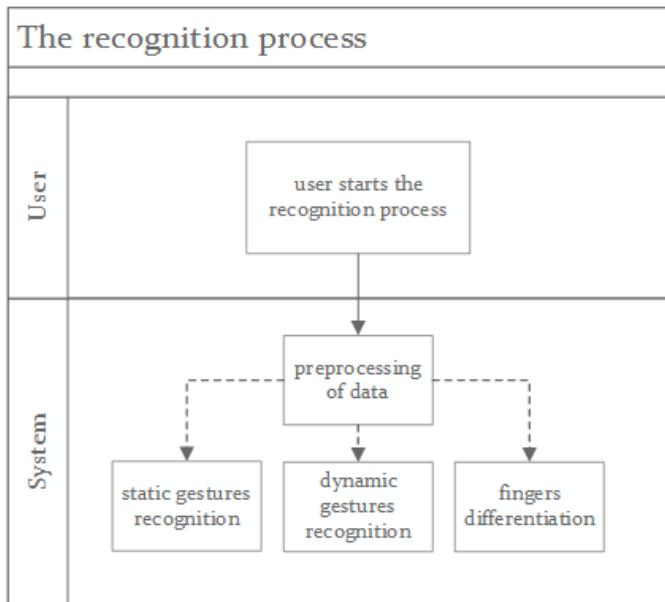


FIGURE 8.2: Diagram showing the recognition process

Modules handlers are implemented using the Observer pattern. For each of the modules are specified relevant events that are reported in the key moments of recognition such as: moment in which particular gesture began to be recognized, moment in which the particular gesture stopped being recognized, moment in which frame processing is finished. The user can handle events received in specific module by implementing appropriate listener interface and adding reference to listener list in specific library module. During gesture recognition process user can use from 1 to 3 listeners. Each listener is on a separate thread.

Methods for a static gesture observer:

- onStart – signals the start of the gesture,
- onFrame – returns a list of recognized gestures with matching probabilities,
- onGesture – signals the end of the gesture and returns a list of recognized gestures with matching probabilities.

Methods for a dynamic gesture observer:

- onStart, onFrame – have the same purposes as in the case of static observer,

- onGesture – signals the end of the gesture and returns a list of recognized gestures with matching probabilities and parameter values.

Methods for a finger differentiation observer:

- onFrame – returns a list of matched classes with probabilities,
- onChange – signals a change of fingers arrangement.

As in the learning process for recognizing static gestures system uses SVM and for dynamic gestures – HMM. For finger differentiation system uses support vector machine. However, there are other classes than those used for static gestures recognition.

8.3 Gesture recorder and visualizer

Recorder and visualizer module is an additional part of library that allows users for easy management of gestures recordings. The recorder collects data from Leap Motion Controller, converts it into data representation described in section regards architecture of library [weryfikacja po sklecienniu pracy] and writes it to the LMR file, which is supported by the LMGesture library. Visualizer enables users to see the recorded gestures stored in LMR format. Data gathered from Leap Motion, used to recognize gestures are very large and the data processed in the library have a specific format. Therefore, it was necessary to create an auxiliary program, that it would facilitate the work with data in the fastest and most user-friendly way.

8.3.1 LMR files

This is a file format specially developed for the LMGesture library, supported by various modules for example by the visualizer. The file structure is as follows:

- Line represents one frame.
- One frame contains: timestamp and hand parameters.
- Hand parameters include: hand id, palm position, stabilized palm position, palm normal vector, palm direction vector and detected fingers parameters.
- Finger parameters include: finger id, finger tip position, stabilized tip position, finger direction vector, finger length and finger width.

Technical information regards LMR files:

- Timestamp and hands are separated by “#”.
- In specified hand occurs hand parameters and fingers.
- Hand parameters are separated by space, and finger are separated by “f” (hands parameters and fingers are separated by “f” too).
- Specified finger has parameters, which are split by “ ”.
- Values in the trivalent parameters are separated by a semicolon.

8.3.2 Visualizer

Visualizer presents the contents of the lmr files. As mentioned earlier, each line of the file is a separately converted frame obtained from Leap Motion Controller. In one moment only one frame is displayed. Almost all fields of the model are visualized:

- palm position,
- palm normal vector,
- palm direction,
- tip position,
- finger direction,
- finger length.

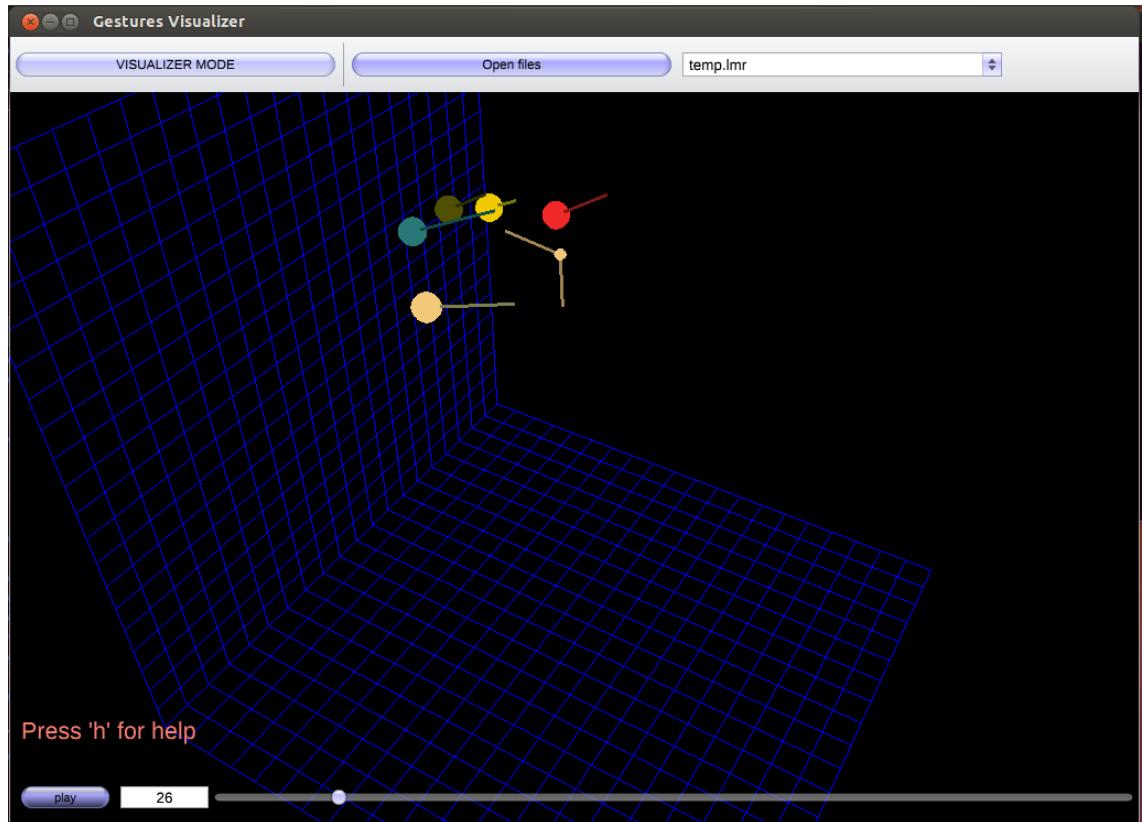


FIGURE 8.3: Screenshot of the visualizer

An example of parameter, which is not visualized, is finger width, in order not to obscure the image.

- Program can read the lmr files chosen by the user.
- The user can load many files at once and select one of the recordings from the drop-down list.
- To facilitate the user work with visualizer, program has implemented windowing interface.
- The user has ability to rotate and zoom camera.

- Slider is available in order to move between frames of the recording.
- Visualizer has option to play recorded gesture, which the user can turn on and off by pressing the play/stop button.
- There is also a button to enter the recorder mode.

8.3.3 Recorder

Recorder is part of the described module, which is responsible for collecting information from Leap Motion Controller and saving it to lmr file. Each frame read from the Leap Motion is captured, and then converted to the previously described model. Then the model is saved by the appropriate sub-module to LMR file. The conversion process contains also sorting hands and fingers. Hands are sorted by X coordinate of palmPosition. In the case of sorting fingers, the usual sort by X coordinate is not enough. The order of the fingers must be independent of hand rotation, therefore a different method for fingers sorting had to be proposed. Fingers are sorted by distance between finger tip position and plane, which is perpendicular to the surface of the hand and contains a direction vector of a hand. This plane can be determined using palm position, direction vector and normalized normal vector of the hand. Below is the formula for the distance (d) between finger tip position and designated plane.

$$d = -(f - pp) \cdot (\hat{hd} \times \hat{hn})$$

Where:

f : is the finger tip position

pp : is the palm position

hd : is the hand direction vector

hn : is the hand normal vector

List of Recorder features:

- Turning on and off recording is done by using space key.
- After recording window automatically appears, in which can be chosen where to save the file.
- The recorder cooperates with the visualizer. During recording performed gesture is visualized.
- There is a button to enter the visualizer mode.

8.4 Used libraries

8.5 Samples of code using the library dedicated to Leap Motion Controller

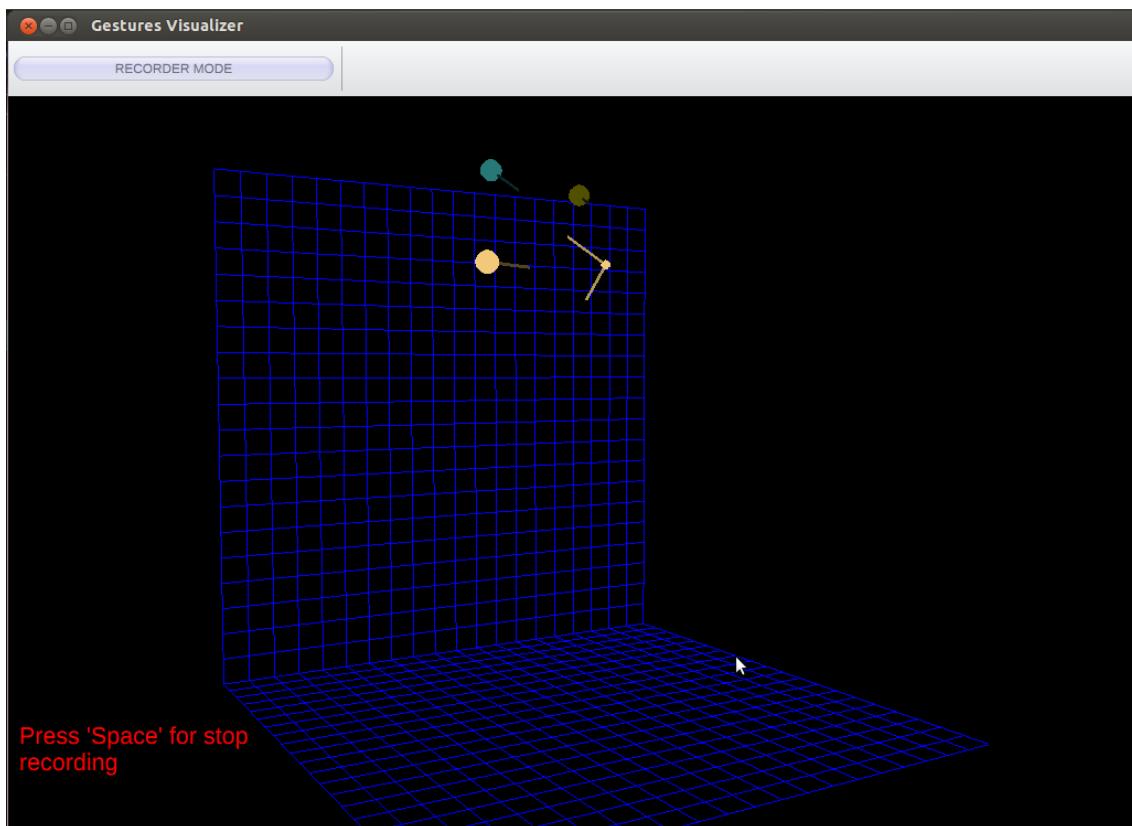


FIGURE 8.4: Screenshot of the recorder

Chapter 9

Conclusions

Bibliography

- [1] 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, 12 1966.
- [2] C. Biemann. Chinese whispers - an efficient graph clustering algorithm and its application to natural language processing problems. *Proceedings of the HLT-NAACL-06 Workshop on Textgraphs-06*, 2006.
- [3] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [4] Chin-Chen Chang, Jiann-Jone Chen, Wen-Kai Tai, and Chin-Chuan Han. New approach for static gesture recognition. *J. Inf. Sci. Eng.*, 22(5):1047–1057, 2006.
- [5] Yen-Ting Chen and Kuo-Tsung Tseng. Developing a multiple-angle hand gesture recognition system for human machine interactions. In *Industrial Electronics Society, 2007. IECON 2007. 33rd Annual Conference of the IEEE*, pages 489–492, 2007.
- [6] Sven Teresiak Chris Biemann. *Disentangling from Babylonian Confusion – Unsupervised Language Identification*. Springer Berlin Heidelberg, 2005. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [7] Vladimir Vapnik Corinna Cortes. *Support-vector networks*.
- [8] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 1–38, 1977.
- [9] Fabio Dominio, Mauro Donadeo, Giulio Marin, Pietro Zanuttigh, and Guido Maria Cortelazzo. Hand gesture recognition with depth data. In *Proceedings of the 4th ACM/IEEE International Workshop on Analysis and Retrieval of Tracked Events and Motion in Imagery Stream*, ARTEMIS '13, pages 9–16, New York, NY, USA, 2013. ACM.
- [10] Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. Liblinear: A library for large linear classification. *J. Mach. Learn. Res.*, 9:1871–1874, June 2008.
- [11] Laurene Fausett, editor. *Fundamentals of Neural Networks: Architectures, Algorithms, and Applications*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1994.
- [12] M. Girvan and M. E. J. Newman. Community structure in social and biological networks. *Proceedings of the National Academy of Sciences*, 99(12):7821–7826, 2002.
- [13] Haitham Hasan and S. Abdul-Kareem. Static hand gesture recognition using neural networks. *Artificial Intelligence Review*, pages 1–35, 2012.
- [14] Andreas Höfer, Aristotelis Hadjidakos, and Max Mühlhäuser. Gyroscope-based conducting gesture recognition. In *NIME 2009 Proceedings: International Conference on New Interfaces for Musical Expression*, pages 175–176, 2009.

- [15] Thomas S. Huang and Vladimir I. Pavlovic. Hand gesture modeling, analysis, and synthesis. In *In Proc. of IEEE International Workshop on Automatic Face and Gesture Recognition*, pages 73–79, 1995.
- [16] M.B. Kaâniche. *Human Gesture Recognition*. 2009.
- [17] Dietrich Kammer, Georg Freitag, Mandy Keck, and Markus Wacker. Taxonomy and overview of multi-touch frameworks: Architecture, scope and features. In *Workshop on Engineering Patterns for Multitouch Interfaces*, Berlin, 2010.
- [18] Maria Karam and M. C. Schraefel. A taxonomy of gestures in human computer interactions. Technical report, 2005.
- [19] N.Y.Y. Kevin, S. Ranganath, and D. Ghosh. Trajectory modeling in gesture recognition using cybergloves reg; and magnetic trackers. In *TENCON 2004. 2004 IEEE Region 10 Conference*, volume A, pages 571–574 Vol. 1, 2004.
- [20] Jonghwa Kim, Stephan Mastnik, and Elisabeth André. Emg-based hand gesture recognition for realtime biosignal interfacing. In *Proceedings of the 13th International Conference on Intelligent User Interfaces*, IUI '08, pages 30–39, New York, NY, USA, 2008. ACM.
- [21] Davis E. King. Dlib-ml: A machine learning toolkit. *J. Mach. Learn. Res.*, 10:1755–1758, December 2009.
- [22] Jiayang Liu, Zhen Wang, Lin Zhong, J. Wickramasuriya, and V. Vasudevan. uwave: Accelerometer-based personalized gesture recognition and its applications. In *Pervasive Computing and Communications, 2009. PerCom 2009. IEEE International Conference on*, pages 1–9, 2009.
- [23] Yun Liu, Zhijie Gan, and Yu Sun. Static hand gesture recognition and its application based on support vector machines. In *Software Engineering, Artificial Intelligence, Networking, and Parallel/Distributed Computing, 2008. SNPD '08. Ninth ACIS International Conference on*, pages 517–521, 2008.
- [24] J. MacQueen. Some methods for classification and analysis of multivariate observations. Proc. 5th Berkeley Symp. Math. Stat. Probab., Univ. Calif. 1965/66, 1, 281–297 (1967)., 1967.
- [25] Vladimir I. Pavlovic, Rajeev Sharma, and Thomas S. Huang. Visual interpretation of hand gestures for human-computer interaction: A review. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19:677–695, 1997.
- [26] Matthew Peeples. R Script for K-Means Cluster Analysis. [on-line] <http://www.mattpeeples.net/kmeans.html>, 2011.
- [27] Valery A Petrushin. Hidden markov models: Fundamentals and applications. In *Online Symposium for Electronics Engineer*, 2000.
- [28] Qifan Pu, Sidhant Gupta, Shyamnath Gollakota, and Shwetak Patel. Whole-home gesture recognition using wireless signals. In *Proceedings of the 19th Annual International Conference on Mobile Computing & Networking*, MobiCom '13, pages 27–38, New York, NY, USA, 2013. ACM.
- [29] Francis Quek, David McNeill, Robert Bryll, Susan Duncan, Xin-Feng Ma, Cemil Kirbas, Karl E. McCullough, and Rashid Ansari. Multimodal human discourse: Gesture and speech. *ACM Trans. Comput.-Hum. Interact.*, 9(3):171–193, September 2002.
- [30] Md. Hafizur Rahman and Jinia Afrin. Hand gesture recognition using multiclass support vector machine. *International Journal of Computer Applications*, 74(1):39–43, July 2013. Published by Foundation of Computer Science, New York, USA.
- [31] Yu Ren and Fengming Zhang. Hand gesture recognition based on meb-svm. In *Embedded Software and Systems, 2009. ICESS '09. International Conference on*, pages 344–349, 2009.

- [32] Hrvoje Benko Michael Haller David Lindbauer Alexandra Ion Shengdong Zhao Roland Aigner, Daniel Wigdor and Jeffrey Tzu Kwan Valino Koh. Understanding mid-air hand gestures: A study of human preferences in usage of gesture types for hci. Technical report, November 2012.
- [33] Peter J. Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20(0):53 – 65, 1987.
- [34] Alexandre Savaris and Aldo von Wangenheim. Comparative evaluation of static gesture recognition techniques based on nearest neighbor, neural networks and support vector machines. *J. Braz. Comp. Soc.*, 16(2):147–162, 2010.
- [35] Xiaohui Shen, Gang Hua, Lance Williams, and Ying Wu. Dynamic hand gesture recognition: An exemplar-based approach from motion divergence fields. *Image Vision Comput.*, 30(3):227–235, March 2012.
- [36] Hugo Steinhaus. Sur la division des corps matériels en parties. *Bull. Acad. Pol. Sci., Cl. III*, 4:801–804, 1957.
- [37] E. Stergiopoulou and N. Papamarkos. Hand gesture recognition using a neural network shape fitting technique. *Eng. Appl. Artif. Intell.*, 22(8):1141–1158, December 2009.
- [38] Robert Y. Wang and Jovan Popović. Real-time hand-tracking with a color glove. *ACM Trans. Graph.*, 28(3):63:1–63:8, July 2009.
- [39] Sabine Weibel, Jens Keil, and Michael Zoellner. Multi-touch gestural interaction in x3d using hidden markov models. In Steven Feiner, Daniel Thalmann, Pascal Guittton, Bernd Fröhlich, Ernst Kruijff, and Martin Hatchet, editors, *VRST*, pages 263–264. ACM, 2008.
- [40] Ying Wu and Thomas S. Huang. Vision-based gesture recognition: A review. In *Proceedings of the International Gesture Workshop on Gesture-Based Communication in Human-Computer Interaction*, GW '99, pages 103–115, London, UK, 1999. Springer-Verlag.
- [41] Deyou Xu, Wuyun Yao, and Yongliang Zhang. Hand gesture interaction for virtual training of spg. In *ICAT Workshops*, pages 672–676. IEEE Computer Society, 2006.
- [42] Jie Yang and Yangsheng Xu. Hidden markov model for gesture recognition. Technical report, DTIC Document, 1994.
- [43] Pujan Ziaie, Thomas Müller, Mary Ellen Foster, and Alois Knoll. A naïve Bayes classifier with distance weighting for hand-gesture recognition. In *Advances in Computer Science and Engineering*, volume 6 of *Communications in Computer and Information Science*, pages 308–315. Springer Berlin Heidelberg, 2009.



© 2014 Michał Nowicki, Olgierd Pilarczyk, Jakub Wąsikowski, Katarzyna Zjawin

Poznań University of Technology
Faculty of Computing
Institute of Computing Science

Typeset using L^AT_EX in Computer Modern.

BibT_EX:

```
@mastersthesis{ key,  
    author = "Michał Nowicki \and Olgierd Pilarczyk \and Jakub Wąsikowski \and Katarzyna Zjawin",  
    title = "{Gesture recognition library for Leap Motion Controller}",  
    school = "Poznań University of Technology",  
    address = "Poznań, Poland",  
    year = "2014",  
}
```