

Online Handwritten Gurmukhi Character Recognition using Support Vector Machine

*Thesis submitted in partial fulfillment of the requirements for the award of
degree of*

**Master of Technology
in
Computer Science and Application**

Submitted By
Dushyant Khurana
(Roll No. 601103005)

Under the supervision of
Dr. R.K. Sharma
Professor, SMCA



**SCHOOL OF MATHEMATICS AND COMPUTER APPLICATIONS
THAPAR UNIVERSITY
PATIALA - 147004**

July 2013

CERTIFICATE

I hereby certify that the work which is being presented in the thesis entitled, "**Online Handwritten Gurmukhi Character Recognition using Support Vector Machine**", in partial fulfillment of the requirements for the award of degree of Master of Technology in Computer Science and Applications submitted in School of Mathematics and Computer Applications, Thapar University, Patiala, is an authentic record of my own work carried out under the supervision of **Dr. R. K. Sharma, Professor, School of Mathematics and Computer Applications.**

The matter presented in this thesis has not been submitted for award of any other degree of this or any other University.



(Dushyant Khurana)

This is to certify that the above statement made by the candidate is correct and true to the best of my knowledge.



(Dr. R. K. Sharma) 15.7.13

Professor

SMCA

Countersigned by:



(Dr. Rajesh Kumar)

Head,

School of Mathematics and Computer Applications

Thapar University,

Patiala.



Dean (Academic Affairs)

Thapar University,

Patiala

ACKNOWLEDGEMENTS

First of all I would like to thank the Almighty, who has always guided me to work on the right path of the life.

This work would not have been possible without the encouragement and able guidance of my supervisor **Dr. R. K. Sharma**. I thank my supervisor for their time, patience, discussions and valuable comments. Their enthusiasm and optimism made this experience both rewarding and enjoyable.

I am equally grateful to **Dr. Rajesh Kumar**, Associate Professor and Head, School of Mathematics and Computer Applications, for motivation and inspiration that triggered me for the thesis work.

I will be failing in my duty if I don't express my gratitude to **Dr. S. K. Mohapatra**, Senior Professor and Dean of Academic Affairs, Thapar University, for making provisions of infrastructure such as library facilities, computer labs equipped with net facilities, immensely useful for the learners to equip themselves with the latest in the field.

I am also thankful to the entire faculty and staff members of School of Mathematics and Computer Applications Department for their direct-indirect help, cooperation, love and affection, which made my stay at Thapar University memorable.

Last but not least, I would like to thank my parents for their wonderful love and encouragement, without their blessings none of this would have been possible. I would also like to thank my brothers, since they insisted that I should do so. I would also like to thank my close friends for their constant support.

Date: July 14th, 2013

Place: Thapar University, Patiala


(Dushyant Khurana)

ABSTRACT

Communication is always an important part of our life, either in the form of speech or writing. Natural Handwriting is one of the easiest ways to exchange information. In this world of technology exchanging information between user and computer is of immense importance and input devices such as keyboard and mouse have limitations as they cannot provide natural handwriting as input. Online handwriting recognition system can be used as an easiest and natural way of communication between user and human computers. Therefore Pen-Based interfaces are becoming more and more popular and hence a lot of research is being done for recognition. Research work presented in this thesis aims to recognize character with higher accuracy written in Gurumukhi script using Support Vector Machine (SVM) by improving processes of pre-processing phase used for recognition. Gurumukhi is a script of Punjabi Language which is widely spoken across the globe. This thesis is divided into five chapters. A brief outline of each chapter is given in the following paragraphs.

First chapter of this report consists of introduction to online handwritten recognition system, issues in online handwritten recognition system overview of Gurmukhi script and literature review. Issues in online handwriting recognition system includes: handwriting style variations; constrained and unconstrained handwriting; personal, situational and material factors; writer dependent vs. writer independent recognition system. In literature review, a detailed literature survey on each phase of established procedure of online handwriting recognition has been presented.

Second chapter gives the detailed work carried out in three phases. They are data collection, pre-processing and feature extraction. In data collection phase, input handwritten strokes are collected is shown. Phases of pre-processing are discussed and algorithms are presented. In the end feature extraction is explained.

Third chapter describes the recognition techniques that can be used for online handwritten recognition system. In this work Support Vector Machine (SVM) is used as a classifier for the recognition. This chapter also illustrates the use of post-processing phase.

Fourth Chapter contains the results of the algorithm proposed in this thesis work and a comparison is made with the proposed algorithms in Agrawal, (2012). The cross validation testing has been done for calculating the accuracy and 3, 4 and 5 fold cross validation testing is applied on a sample of 30, 50 and 70 of each zone.

Finally, the result of the thesis is concluded in this chapter. Future scope of the work is also discussed.

TABLE OF CONTENTS

CERTIFICATE	i
ACKNOWLEDGEMENTS	ii
ABSTRACT	iii
TABLE OF CONTENTS	v
LIST OF FIGURES	vii
LIST OF TABLES	xi
ABBREVIATIONS	xii
Chapter 1 INTRODUCTION	1
1.1 Classification of Handwriting Recognition	2
1.1.1 Offline Handwriting Recognition	2
1.1.2 Online Handwriting Recognition	3
1.1.3 Online Recognition vs. Offline Recognition	3
1.2 Introduction to Gurmukhi Script	4
1.2.1 Vowels in Gurmukhi Script	7
1.2.2 Zone in Gurmukhi Script	8
1.3 Issues in Recognition	9
1.3.1 Handwriting styles variations	9
1.3.2 Constrained and Unconstrained Handwriting	10
1.3.3 Personal and Situational Aspects	11
1.3.4 Writer dependent vs. writer independent recognition system	11
1.4 Online Handwritten Character Recognition System	12
1.4.1 Data Acquisition	13

1.4.2 Pre-processing	13
1.4.3 Feature Extraction	14
1.4.4 Recognition	15
1.5 Literature Survey	16
1.5.1 Work related to pre-processing	17
1.5.2 Work Related to Character Recognition	19
1.5.3 Work on Gurmukhi Character Recognition	20
Chapter 2 DATA COLLECTION, PRE-PROCESSING AND FEATURE EXTRACTION	22
2.1 Data Capturing	22
2.2 Pre-processing Phase	23
2.2.1 Removal of Duplicate Points	23
2.2.2 Size Normalization and Centering	24
2.2.3 Missing Point Interpolation	26
2.2.4 Resampling	28
2.3 Feature Extraction	30
Chapter 3 GURMUKHI CHARACTER RECOGNITION AND POST- PROCESSING	31
3.1 Introduction to SVM	31
3.2 Recognition of Gurmukhi Character	34
3.3 Post-Processing	45
Chapter 4 RESULTS AND DISCUSSIONS	47
4.1 Scheme 1: k -fold cross validation on date set of 30 samples for each zone	48
4.2 Scheme 2: k -fold cross validation on date set of 30 samples for each zone	54
4.3 Scheme 3: k -fold cross validation on date set of 30 samples for each zone	60
Chapter 5 CONCLUSION AND FUTURE SCOPE	67
REFERENCES	69

LIST OF FIGURES

Figure 1.1	Areas of handwriting recognition	2
Figure 1.2	A tablet digitizer, input sampling and communication to the computer	3
Figure 1.3	Three zones and headline in Gurmukhi word	8
Figure 1.4	Variation in some of the Gurmukhi Characters by five writers	10
Figure 1.5	Different styles of writing	10
Figure 1.6	Types of writing styles	11
Figure 1.7	Phases of online handwritten character recognition	12
Figure 1.8	Commonly used hardware devices for data acquisition	13
Figure 1.9	Pre-processing phase of recognition system	14
Figure 2.1	A sample stroke of Gurmukhi character	22
Figure 2.2	A sample Gurmukhi character of 3 strokes	23
Figure 2.3	A sample stroke showing duplicate or overlapped points	24
Figure 2.4	Handwriting character after normalization of a sample Gurmukhi character (ka)	25
Figure 2.5	A sample stroke of Gurmukhi showing (a)-Inputted stroke (b) normalized stroke (c) Interpolated stroke	28
Figure 2.6	Shows the entire algorithm (a) input stroke by the writer (b) normalized and centering (c) interpolation missing points (d) resampling of points	29
Figure 3.1	Linear classification of objects by SVM into two classes	32

Figure 3.2	Mapping of SVM classifications from complex to linear: (a) classification by a complex kernel function (b) mapping of input space of figure (a) into feature space	32
Figure 3.3	Concept of support vectors	33
Figure 3.4	Pre-processed points in SVM format	43
Figure 3.5	Scaled file of pre-processed points	44
Figure 3.6	Model file of pre-processed points	45
Figure 4.1	Surface graph showing variation in two methods for 3 fold cross validation on a data set of 30 samples for lower zone	49
Figure 4.2	Surface graph showing variation in two methods for 4 fold cross validation on a data set of 30 samples for lower zone	49
Figure 4.3	Surface graph showing variation in two methods for 5 fold cross validation on a data set of 30 samples for lower zone.	50
Figure 4.4	Surface graph showing variation in two methods for 3 fold cross validation on a data set of 30 samples for middle zone.	51
Figure 4.5	Surface graph showing variation in two methods for 4 fold cross validation on a data set of 30 samples for middle zone.	51
Figure 4.6	Surface graph showing variation in two methods for 5 fold cross validation on a data set of 30 samples for middle zone.	52
Figure 4.7	Surface graph showing variation in two methods for 3 fold cross validation on a data set of 30 samples for upper zone.	53
Figure 4.8	Surface graph showing variation in two methods for 4 fold cross validation on a data set of 30 samples for upper zone.	53

Figure 4.9	Surface graph showing variation in two methods for 5 fold cross validation on a data set of 30 samples for upper zone.	54
Figure 4.10	Surface graph showing variation in two methods for 3 fold cross validation on a data set of 50 samples for lower zone.	55
Figure 4.11	Surface graph showing variation in two methods for 4 fold cross validation on a data set of 50 samples for lower zone.	55
Figure 4.12	Surface graph showing variation in two methods for 5 fold cross validation on a data set of 50 samples for lower zone.	56
Figure 4.13	Surface graph showing variation in two methods for 3 fold cross validation on a data set of 50 samples for middle zone.	57
Figure 4.14	Surface graph showing variation in two methods for 4 fold cross validation on a data set of 50 samples for middle zone.	57
Figure 4.15	Surface graph showing variation in two methods for 5 fold cross validation on a data set of 50 samples for middle zone.	58
Figure 4.16	Surface graph showing variation in two methods for 3 fold cross validation on a data set of 50 samples for upper zone.	59
Figure 4.17	Surface graph showing variation in two methods for 4 fold cross validation on a data set of 50 samples for upper zone.	59
Figure 4.18	Surface graph showing variation in two methods for 5 fold cross validation on a data set of 50 samples for upper zone.	60
Figure 4.19	Surface graph showing variation in two methods for 3 fold cross validation on a data set of 70 samples for lower zone.	61
Figure 4.20	Surface graph showing variation in two methods for 4 fold cross	61

	validation on a data set of 70 samples for lower zone.	
Figure 4.21	Surface graph showing variation in two methods for 5 fold cross validation on a data set of 70 samples for lower zone.	62
Figure 4.22	Surface graph showing variation in two methods for 3 fold cross validation on a data set of 70 samples for middle zone.	63
Figure 4.23	Surface graph showing variation in two methods for 4 fold cross validation on a data set of 70 samples for middle zone.	63
Figure 4.24	Surface graph showing variation in two methods for 5 fold cross validation on a data set of 70 samples for middle zone.	64
Figure 4.25	Surface graph showing variation in two methods for 3 fold cross validation on a data set of 70 samples for upper zone.	65
Figure 4.26	Surface graph showing variation in two methods for 4 fold cross validation on a data set of 70 samples for middle zone.	65
Figure 4.27	Surface graph showing variation in two methods for 5 fold cross validation on a data set of 70 samples for middle zone.	66

LIST OF TABLES

Table No.	Title of Table	Page No.
Table 1.1	Character Set of Gurmukhi Script	4
Table 1.2	Unique Vowel Character	6
Table 1.3	Six Special consonants in Gurmukhi	7
Table 1.4	Dependent vowel lists	7
Table 3.1	Lower zone stroke ids	34
Table 3.2	Upper zone stroke ids	35
Table 3.3	Middle zone stroke ids	35
Table 4.1	Cross validation accuracy of lower zone for 3, 4, and 5 cross validation for 30 samples	48
Table 4.2	Cross validation accuracy of middle zone for 3, 4, and 5 cross validation for 30 samples	50
Table 4.3	Cross validation accuracy of upper zone for 3, 4, and 5 cross validation for 30 samples	52
Table 4.4	Cross validation accuracy of lower zone for 3, 4, and 5 cross validation for 50 samples	54
Table 4.5	Cross validation accuracy of middle zone for 3, 4, and 5 cross validation for 50 samples	56
Table 4.6	Cross validation accuracy of upper zone for 3, 4, and 5 cross validation for 50 samples	58
Table 4.7	Cross validation accuracy of lower zone for 3, 4, and 5 cross validation for 70 samples	60
Table 4.8	Cross validation accuracy of middle zone for 3, 4, and 5 cross validation for 70 samples	62
Table 4.9	Cross validation accuracy of upper zone for 3, 4, and 5 cross validation for 70 samples	64

ABBREVIATIONS

Abbreviations

PC

PDA

HMM

SVM

SVC

SVR

OCR

OHCR

Expanded Forms

Personal Computer

Personal Digital Assistant

Hidden Markov Model

Support Vector Machine

Support Vector Classifier

Support Vector Regressor

Optical Character Recognition

Online Handwritten Character Recognition

CHAPTER 1

INTRODUCTION

In this era, everybody wants to perform one's tasks as easily as possible. Interchanging information between user and the computer can be done by natural handwriting and this seems to be the easiest way. Hence, the research on the topic of online handwriting recognition system is on boom. Most of the users now use touch screen mobile phones, tablets for their needs and ease but still they cannot use their own handwriting for communication with their phones even with this high processing capability of the computers and mobile phone specifically in Indian languages. In order to fully utilize the high processing capability of the computer or tablets, a user interface is needed which should not only be efficient but also natural to the user. This naturality can be there in terms of acceptance of input by computers in user's own handwriting in his own language.

In general it takes good amount of time to input data into a computer through keyboard and the one who is not literate cannot use it. Same is for the output devices that were of huge size. But now both input and output process can be done on a single compact device, where a user gives input through an electronic pen and finds the output simultaneously and hence the data can be recorded. Sometimes we need textual form of data and online handwritten recognition systems got placed which make it very much easier for the user to input the data that can be converted into its digital or text form without using the keyboard and the additional benefit is the data can be entered at a faster rate. Thus, the handwriting recognition field has great potential to improve the communication between the user and the computer. Handwriting recognition is the additional ability of a computer to interpret the data which is intelligibly handwriting input from various resources such as paper documents, electronic pen tablets, touch screen and other devices and then generate a description of the data in a desired format.

Handwriting recognition is in research for many decades and many researches are going on all over the world. Great advance have been made in this field and due to this reason, the usage and reliability of online handwriting based devices such as Tablet PCs and Personal Digital Assistants (PDAs) have increased a lot. There are

many applications of the handwriting recognition systems. Some examples are: this can act as a communication tool for entering information, tool for reading bank cheques fields, converting document into textual form, *etc.*

1.1 Classification of Handwriting Recognition

The classification of handwriting recognition can be done into two major categories shown in Figure 1.1.

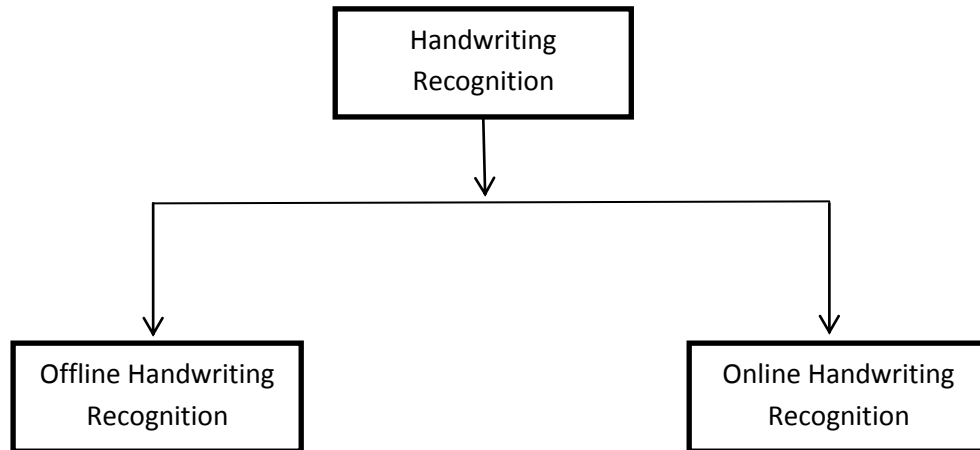


Figure 1.1: Classification of handwriting recognition

1.1.1 Offline Handwriting Recognition

Under this type recognition systems, the writing is available in the form of image which is captured with the help of a scanner or a camera which captures the writing optically. This means that the text is not recognised at the same time when it is produced but after the writer finishes writing, *i.e.*, in this case, the text is written on a surface such as paper and from there it is then passed to computer in the image form. It is then recognised by the handwriting recognition system. This is done in a sequential process: firstly the image is stored digitally in grey scale format, *e.g.*, bitmap image, and then further processing is done on it to get good recognition accuracy. Features for recognition can also be enhanced and then extracted from the stored bitmap image by applying desired functions. This type of recognition is called as Optical Character Recognition. Recognition of machine printed text or characters is also a part of OCR. Offline methods are less suitable for communication between man and machine because it has no real time interactivity. They are more suitable for automatic conversion of paper documents to electric documents which may be intercepted by computers.

1.1.2 Online Handwriting Recognition

In comparison to the offline methods of recognition, online handwriting recognition is done in real time, *i.e.*, at the same time as the handwriting is produced. The surface used for handwriting is usually electronic device called as digitized tablet and is used along with a pen called as “Stylus” with the help of which writer writes on it. As the pen moves on the surface of the digitized pad, the two dimensional coordinates will be captured and stored as a function of time as shown in Figure 1.2.

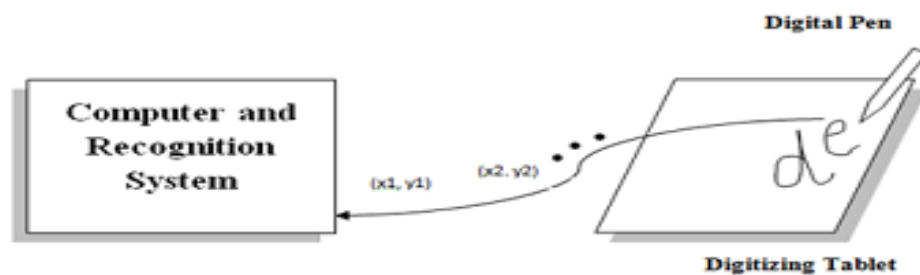


Figure 1.2: A tablet digitizer, input sampling and communication to the computer

1.1.3 Online Recognition vs. Offline Recognition

- **Real Time Interactivity:** Online Handwriting recognition is real time process, *i.e.*, recognition is done when data is written but in case of offline recognition is done after the writing is completed or at a later time.
- **User Adaptation:** In case of online methods, recognition can be corrected immediately. If a particular character is not being recognized properly, user can change his writing style then and there, and hence adapts to online method.
- **Pre-processing:** Pre-processing operations such as smoothing, normalization, detection of loops *etc.*, is faster in case of online methods. Thus, lesser time is required for pre-processing operations in case of online handwriting as compared to offline handwriting.
- **Easy Segmentation:** As dynamic information is collected for each stroke in case of online recognition methods, therefore segmentation operations are much easier in online handwriting recognition as compared to offline.

- Special electronic devices are needed in case of online handwriting recognition *e.g.*, Digital tablets, PDAs, *etc.*, but in case of offline we need just a pen and paper. Offline methods are well suited for recognising printed documents.

In this work, we have focussed on improving of recognition accuracy by improving the processes of pre-processing. The stages of a handwriting recognition system, including the pre-processing stage are described in section 1.4. The present work has also focussed on the improvement of recognition accuracy of strokes that are used to form a Punjabi akshara. Next section introduces, in brief, the Gurmukhi script, which is used to write Punjabi language.

1.2 Introduction to Gurmukhi Script

Gurmukhi is the script used for writing Punjabi language in India and across the world which is the world's 14th most widely spoken language. The name Gurmukhi derived from the old Punjabi term "Gurmukhi", meaning "from the mouth of the Guru". Gurmukhi script is written in left-to-right direction and in top-down approach. Most of the characters have a horizontal line on the upper part by which characters of the word are connected. This is called as the headline. Gurmukhi has 41 consonants, 9 vowel symbols, and two symbols for nasal sound, 1 symbol which duplicates the sound of any consonant (addak), 3 subjoined forms of the consonant Rara, Haha and Vava. Table 1.1 shows these Gurmukhi characters.

The Gurmukhi script contains 35 letters. Out of which the three letters are distinct in Gurmukhi script as they form the basis of vowel and are distinct because they form the basis for vowels and are not consonants.

Table 1.1: Character Set of Gurmukhi Script

Character	Character Name
ੳ	URHA
ਊ	ERHA
ਏ	EERHI

ਸ	SUSSA (Sa)
ਹ	HAHA (Ha)
ਕ	KUKKA (Ka)
ਖ	KHUKHA (Kha)
ਗ	GUGGA (Ga)
ਘ	GHUGGA (Gha)
ਙ	UNGGA (Nga)
ਚ	CHUCHA (Ca)
ਛ	CHHUCHHA (Cha)
ਜ	JUJJA (Ja)
ਝ	JHUIHA (Jha)
ਞ	YANZA (Nya)
ਟ	TAINKA (Tta)
ਠ	THUTHA (Ttha)
ਡ	DUDDA (Dda)
ਢ	DHUDDA (Ddha)
ਣ	NAHNHA (Nna)
ਤ	TUTTA (Ta)
ਥ	THUTHA (Tha)
ਦ	DUDA (Da)
ਧ	DHUDDA (Dha)

ਨ	NUNNA (Na)
ਪ	PUPPA (Pa)
ਫ	PHUPHA (Pha)
ਬ	BUBBA (Ba)
ਭ	BHUBBA (Bha)
ਮ	MUMMA (Ma)
ਯ	YAIYYA (Ya)
ਰ	RARA (Ra)
ਲ	LULLA (La)
ਵ	VAVA (Va)
ੜ	RARA (Rha)

Out of these character first three are unique vowels, these characters are never used on their own. Remaining are the consonants of the Gurmukhi script.

Table 1.2: Unique Vowel Characters

ੳ	URHA
ਐ	ERHA
ੲ	EERHI

In addition to these, there are six consonants created by placing a dot (bindi) at the foot (pair) of the consonant these are shown in the Table 1.3.

Table 1.3: Six Special consonants in Gurmukhi

ਸ਼	SUSSA PAIR BINDI (Sha)
ਖ਼	KHUKHA PAIR BINDI (Khha)
ਗ਼	GUGGA PAIR BINDI (Ghha)
ਜ਼	JAJJA PAIR BINDI (Za)
ਫ਼	PHUPHA PAIR BINDI (Faa)
ਲ਼	LALLA PAIR BINDI (Lla)

1.2.1 Vowels in Gurmukhi Script

Gurmukhi follows similar concepts to other Brahmi scripts. All consonants are followed by an inherent ‘a’ sound. This inherent vowel sound can be changed by using dependent vowel signs which attach to bear consonants. In some cases, dependent vowel signs cannot be used at the beginning of a word or syllable for instance and so an independent vowel character is used instead. Table 1.4 shows a list of dependent vowels.

Table 1.4: Dependent vowel lists

◌੍	MUKTA (a)
◌ਾ	KANNA (aa)
◌ਿ	SIHARI (i)
◌ੀ	BIHARI (ii)
◌ੇ	LAVAN (ee)
◌ੈ	DULAVAN (ai)
◌ੁ	ONKAR (u)

ੳ	DULANKAR (uu)
ੳ	HORA (oo)
ੳ	KANAURA (au)

1.2.2 Zones in Gurmukhi Script

A word in Gurmukhi can be portioned into three zones these are:

Upper Zone: It is the region above the headline, where vowels reside

Middle Zone: It represents the area where the consonants and some other parts of vowels reside, *i.e.*, the area below the headline but above the lower zone.

Lower Zone: It represents the area below the middle zone where some vowels, halant or certain half character lie at the foot of the consonants. These zones are shown in Figure. 1.3 for the Gurmukhi word (ਦੁਲੰਕਾਰ).



Figure 1.3: Three zones and headline in Gurmukhi word

Characteristics of Gurmukhi Scripts

- Gurmukhi alphabet is syllabic, *i.e.*, all consonants have an inherent vowel.
- Unlike Roman characters, Gurmukhi is written below the line and it has no concept of lower and upper case.
- The Gurmukhi script, unlike the Greek and Roman alphabets, is arranged in a logical fashion: vowel first, then consonants and semivowels
- Upper zone and middle zone are separated by a line called head line or *siro rekha* as present in other Indian scripts also such as Devnagari.

1.3 Issues in Handwritten Character Recognition

The online handwriting recognition is used for identification of characters and it is used with the devices such as PDAs, touch screen phones, cross pad and tablet PCs, but due to variability in handwriting styles and distortions caused by digitizing process, even the best handwriting character recognizer is unreliable. Distortion can be because of the electronic device used which takes the input and send to the recognizer and if the input is distorted good result cannot be gained. In order to use these input devices, the accuracy achieved by the handwriting recognizer must be sufficiently high so that it is acceptable by the user. Some of the other issues are explained, in brief, in next subsections.

1.3.1 Handwriting Styles Variations

This is the major issue in recognition of a character as different writer has a different way of writing a character, means different writer different handwriting, different styles. Many times a person finds himself/herself unable to recognize his/her own handwriting. Hence, practically it is much difficult for a recognizer recognize handwriting efficiently. Variations like Deformed Geometry, slants, skews, overlapping, distortions are inserted by different writers in different ways. Geometric properties like aspect ratio, position and size vary. Some writer writes very small some writes very big characters. Figure 1.4 shows a sample of variations in five writers writing the same character of the Gurmukhi script although such samples share a high degree of similarities. The character may look similar but the no. of strokes, direction and the drawing may vary considerably.

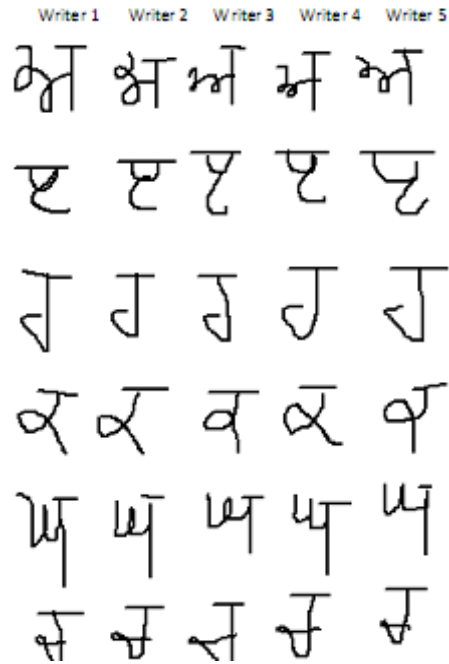
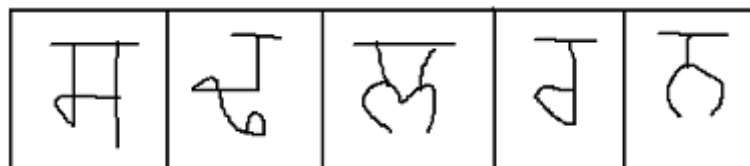


Figure 1.4: Variation in some of the Gurmukhi Characters by five writers

1.3.2 Constrained and Unconstrained Handwriting

Handwriting styles could be constrained or unconstrained. Constrained handwriting is boxed discrete or spaced discrete in nature. In boxed discrete handwriting, unconstrained handwriting is cursive or mixed cursive Figure 1.4 shows the constrained and unconstrained handwriting as the writing changes from constrained to unconstrained the difficulty level of recognizing also increases. As in constrained writing each character is written separately with spaces and no character touches the other character but in other case character in one word is connected and strokes are used more than once in individual character, it is referred to as cursive writing. As there is no clear boundaries are specified between characters to distinguish them hence recognition is difficult. Also, it is observed that many users write in mixed cursive writing illustrated in Figure 1.5 (Agrawal, 2011).



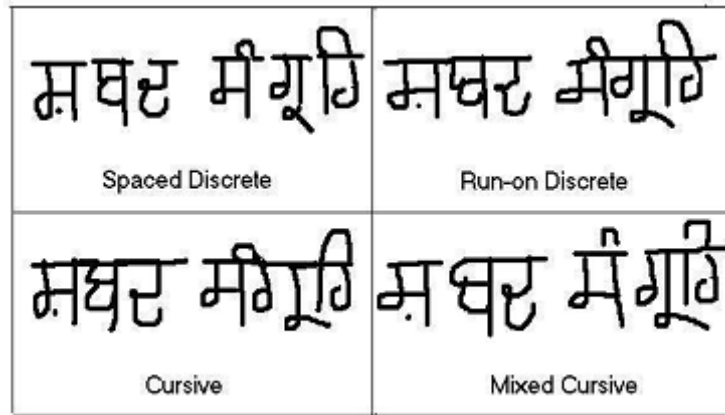


Figure 1.5: Different styles of writing

1.3.3 Personal and Situational Aspects

Personal factors include writer's writing style which might be affected by handedness, either left handed or right handed and many writers are habitual to write random or specific inclined text lines, Some users write very fast some write very slow. A good recognition machine requires a neat and clean handwriting and this writing style also depends on profession of writers to some extent as shown in Figure 1.6 (Tappert, 1984).

The situational aspects depend on the facts either writer is interested or not to write, how much attention or concentration writer is paying, whether there was any interruption while writing. The material used in writing may provide comfort or distortion to writer that result in variations in handwriting. This also includes the position and size of the writing pad.

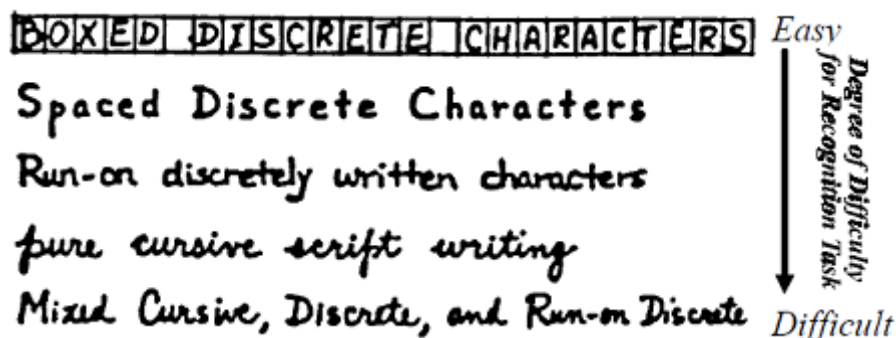


Figure 1.6: Types of writing styles

1.3.4 Writer dependent vs. writer independent recognition system

The writer dependent recognition system is used to recognize the samples of only those writers whose samples are taken to train the recognition system. It is specific to

a group of writers or a single writer, if it is train with the sample taken from a single writer. In writer dependent system, all possible style variations can be trained to the system, hence a higher recognition rate can be achieved.

On the other hand writer independent system is a generalized recognition system which works for the entire writer means, it recognizes the handwriting samples of unknown writers. Hence, it needs to train the system with all possible and commonly used style variations. Therefore it needs to train the system with large number of samples from a large number of different writers, to make the recognition system generalized, therefore the writer independent system has a comparatively lower recognition rate. In practice, the writer independent system is in more demand because of generalized application.

1.4 Online Handwritten Character Recognition System

It involves many steps to completely recognize and produce machine encoded text. These phases are termed as Data acquisition, Pre-processing, Feature extraction, recognition and post-processing. The architecture of these phases is shown in Figure 1.7 and these phases are listed below with a brief description. These phases are elaborated in the next sections.

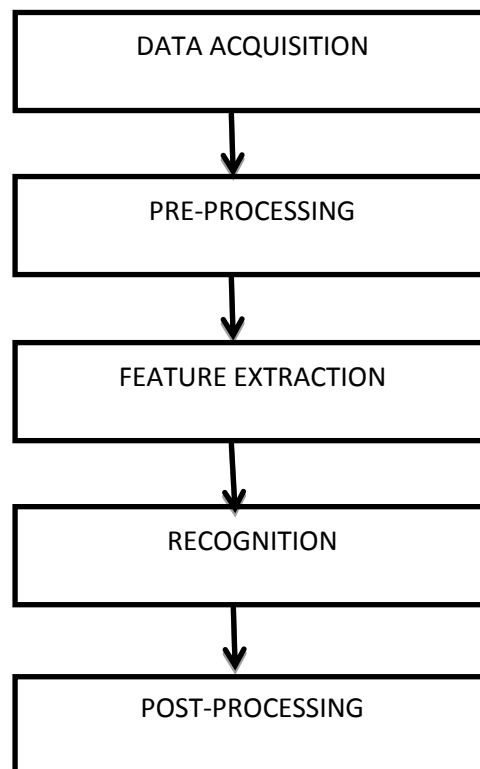


Figure 1.7: Phases of online handwritten character recognition

1.4.1 Data Acquisition

Online handwriting recognition requires a device that captures the writing as it is written by the writer. The most common devices are electronic tablet or digitizer. These devices use a pen that is digital in nature. Data collection is the first phase in online handwriting recognition that collects the sequence of coordinate points a moving pen. A typical pen includes two actions, namely, PenDown and PenUp. The pen traces, between PenDown and PenUp is called as stroke. These pen traces are sampled at a constant rate, therefore these pen traces are evenly distributed in time not in space. The appearances of Personal Digital and Assistants, cross pad and tablet are shown in Figure 1.8 used for data acquisition.



Figure 1.8: Commonly used hardware devices for data acquisition

1.4.2 Pre-processing

This phase in handwriting recognition system is applied to remove noise or distortion present in input stroke due to hardware and software limitations like: irregular size of the stroke, missing points of coordinates while capturing pen movement, jitter present in the text. Left or right bend in handwriting and uneven distances of points from neighbouring positions with many more problems like,

- Accidental pen lifts which occur when the writing pressure falls below the tablet activation threshold
- Repetitive or redundant samples when the pen is stationary or moving slowly
- Isolated samples which are distant from the general outline due to the irregular motion of the writer's hand or imperfections in the tablet

To remove these types of error pre-processing is required in online handwriting recognition system, it includes five common steps, namely, removal of duplicate points, size normalization, centering, interpolation missing points, and resampling of points. These steps are shown in Fig 1.9.

1.4.3 Feature Extraction

Feature extraction is extracting information from raw data which is most relevant for classification purpose and that minimizes the variation in class and maximizes the variation between classes.

Selection of a feature extraction method is probably the single most important factor in achieving high recognition performance in recognition system. The computational complexity of a classification problem can also be reduced if suitable features are selected.

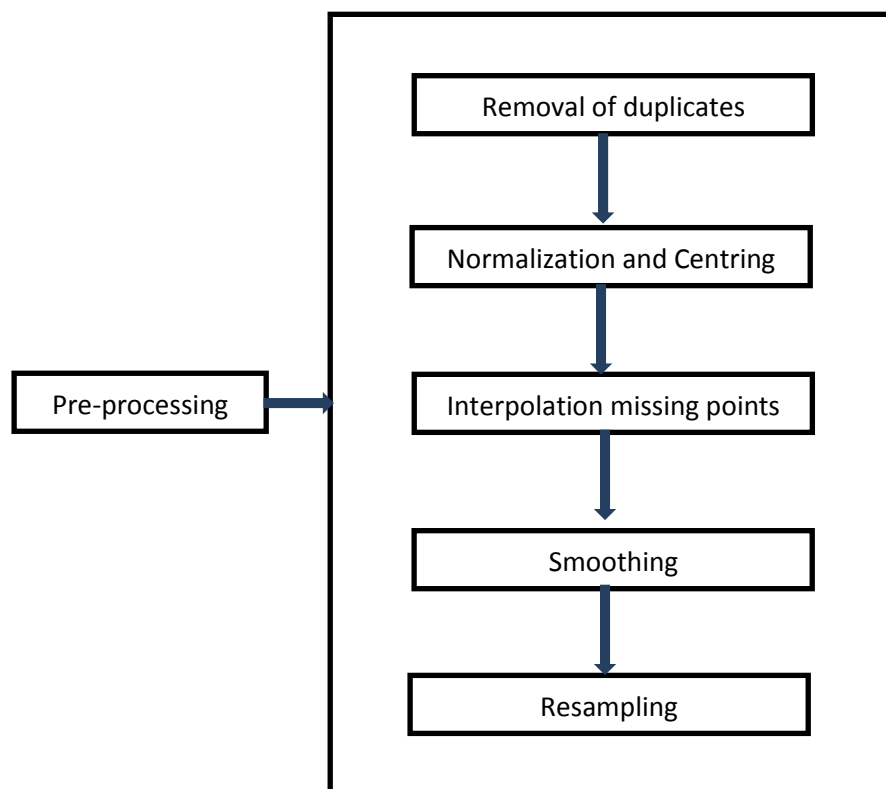


Figure 1.9: Pre-processing phase of recognition system.

1.4.4 Recognition

This is the most important phase of the online handwriting recognition system and uses the features, extracted in the previous stage, to identify the input stroke. In order to obtain the class labels of character or stroke, variety of classification methods can be used like statistical methods, structural and syntactical method, elastic matching method, neural network methods *etc.*, a brief introduction of these methods are described below.

- **Statistical Classification Methods**

Statistical classification methods are based on the Bayes decision theory, which aims to minimize the loss of classification with given loss matrix and estimated probabilities. According to the class-conditional probability density estimation approach, statistical classification methods are divided into parametric and nonparametric ones. Parametric methods assume a functional form for the density of each class *e.g.*, Gaussian classifier. Nonparametric classifier does not assume any functional form for the conditional distributions. The two approaches related to non-parametric are k-nearest neighbour and Parse widow.

- **Artificial Neural Networks**

Artificial neural networks are similar to neuron system. A neural network is composed of a number of interconnected neurons, and the manner of interconnection differentiates the network models into feed forward networks, recurrent networks, self-organizing networks, and so on. Single layer neural network and multilayer perception are the example of ANN classifier.

- **Support Vector Machines (SVM)**

Support Vector Machine is relatively new classification technique producing efficient recognition results. SVM is a new type of hyperplane classifier developed and based on statistical learning theory given by Vapnik. Basically SVM is a binary (two class) linear classifier in kernel induced feature space and is formulated as a weighted combination of kernel functions on training examples. The kernel function represents the inner product of two vectors in linear/nonlinear feature space. Multiclass classification is accomplished by combining multiple binary SVM classifiers.

- **Structural Classification Methods**

Structural methods recognize patterns via elastic matching of strings, graphs, or other structural descriptions. Structural pattern recognition methods are used more often in online character recognition than in offline character recognition. Instead of representing the character pattern as a feature vector of fixed dimensionality, structural methods represent a pattern as a structure (string, tree or graph) of flexible size. The structural representation records the stroke sequence or topological shape of the character pattern, and hence resembles well to the mechanism of human perception. In recognition, each class is represented as one or more structural template, the structure of the input pattern is matched with the templates and is classified in the class of the template of minimum distance or maximum similarity.

- **Multiple Classifier Methods**

By multiple classifier methods, the classifications of multiple classifiers are combined to reorder the class. Different classifier show variable performance, varying classification accuracy and speed, and different errors on concrete patterns. Hence, it is natural to combine the strengths of different classifier to achieve higher performance. This is usually accomplished by combining the decision of multiple classifiers. In research, many combination methods have been proposed, and the applications to practical problems have proven the advantage of the ensemble over individual classifier. The combination of different classifier can be categorized as parallel (horizontal) and sequential (vertical, cascaded). A parallel combination is most often adopted for improving the classification accuracy, whereas sequential combination is mainly used for acceleration the classification of a large category set.

1.5 Literature survey

The basic idea of doing this survey is to find that how many different types of techniques have been applied for the recognition and to see the recognition accuracy in percentage of their respective methods. The survey is done on both online and offline handwriting recognition systems with a special emphasis on the Indian scripts.

In our literature survey we have studied many research approaches practiced on many languages and scripts particularly in the Indian context. Our emphasis is to study and analyze the best pre-processing and feature extraction approaches. Observed or reported results and many other relevant issues considerably to any new research work on character recognition. After our detailed survey we were able to discover the theme of our proposed work including the techniques we have incorporated in various phases, to implement it and to evaluate our results and conclusions.

1.5.1 Work Related to Pre-processing

A survey carried out by Tappert *et al.* (1988) and by Tappert *et al.* (1990) identifies different pre-processing techniques like External and internal segmentation; and noise reduction including smoothing, filtering, wild point correction, dehooking, dot reduction and stroke connection; normalization. The normalization technique consists of deskewing, base line drift correction, size normalization and stroke length normalization.

Pre-processing technique for online recognition given by (Nair and Leedham, 1991). They present solutions to some of the major problems that occur in recognition system. They describe seven steps for pre-processing a stroke before they can move on to recognition phase. These steps are:

- Minimum distance filtering.
- Modified gear backlash smoothing.
- Minimum distance filtering.
- Straight line average smoothing.
- Minimum distance filtering.
- Straight line average smoothing with angle constraint.
- Serif Removal.

Homayoon *et al.* (1994) have mainly discussed two pre-processing steps for an online handwriting recognition system, which are size normalization and slope correction. According to them while performing size normalization, the base-line and the mid-line need to be estimated. The area surrounded by the base-line and the mid-line is the only part of any word which is always non-empty. Once accurate estimates of the base-line and the mid-line are obtained, a magnification factor can be computed from

the ratio of the nominal mid-portion size and that of the input. The entire input data may then be scaled using the obtained magnification factor. Then slope correction is performed which is based on the evaluation of the mean velocities in the x and y directions. The angle between these velocity vectors is used to estimate the angle by which the words should be rotated. For special cases such as in words with all block capital letters, special provisions was taken to avoid wrong estimation of the rotation angle.

Aparna *et al.* (2004) proposed a pre-processing technique that consists of interpolation, smoothing and normalization of strokes. For normalizing a stroke they first determine the bounding box of the entire stroke. Then they divide both dimensions of the stroke by the 'height' of the horizontal block. This preserves the relative size of strokes in a horizontal block. The strokes are then converted onto curve length base (sampled uniformly along curve length) and then smoothed independently along t-axis using a Gaussian filter.

Nair *et al.* (2008) discussed the problem arises on writing like accident pen lifts which occur when the writing pressure falls below the tablet activation threshold, redundant samples when the pen is stationary, isolated samples due to irregular motion, serifs at the start and end of pen strokes. They have also given the solution for these problems in their work.

Zheng and Zhu (2006) demonstrated that hardware system collects handwriting signature signals, including the position signals of and in real time. The space resolution is 4900x4900 and the sampling interval is 10ms. First, they removed the gaps between strokes, and these gaps reflect the time that the handwriting pen does not contact the handwriting pad and the positions that strokes begin and end. Simultaneously, they normalized the size, the length of the signature

Hosny *et al.* (2011) has given pre-processing technique for Arabic characters using four steps:

- **Duplicate Points Removal:** By checking whether the coordinates of any two points in a stroke are the same. If so, one of them is removed.
- **Interpolating Points:** To add any missing points by linear interpolation.

- **Smoothing:** To eliminate hardware imperfections and trembles in writing each point is substituted with the weighted average of its neighbouring points.
- **Re-sampling:** Due to the variation in writing speed, the acquired points are not distributed evenly along the stroke trajectory. This operation is used to get a sequence of points which is equidistant.

1.5.2 Work Related to Character Recognition

For Indian languages most of the research work is performed on firstly on Devnagri and secondly on Bangla script. Pal and Chaudhary (2004) presented a survey on Indian Scripts Character Recognition. This paper introduces the properties of Indian scripts and work and methodologies approached to different Indian scripts. They have presented the study of the work for character recognition on many Indian languages scripts including Devnagari, Bangla, Tamil, Gurmukhi, Gujarati and Kannada.

Wakabayashi *et al.* (2009) also presented a comparative study of Devnagari and handwritten character recognition using different features and classifier. They used four sets of features based on curvature and gradient information obtained from binary as well as grayscale images and compared result using 12 different classifiers as concluded the best results 94.94% and 95.19% for features extracted from binary and gray image respectively obtained with Mirror Image Learning (MIL) classifier. They also concluded curvature features to use for better results than gradient features for most of classifiers.

A later review of research on Devnagari character recognition is also presented by Dungre *et al.* (2010). They have reviewed the technique the techniques available for character recognition. They have introduced image pre-processing techniques for thresholding, skew detection and correction, size normalization and thinning which are used in character recognition. They have also reviewed the feature extraction using global transformation and series expansion like Fourier transform, Gabor transform, Wavelet moments, Statistical features like zoning, projections, crossing and distance and some geometrical and topological features commonly practiced. They also reviewed the classification using template matching, statically techniques, neural network, SVM and the combination of a classifier for better accuracy is practiced recognition.

Mukherji and Rege, (2009) used shape features and fuzzy logic to recognize offline Devnagari Character recognition. They segmented the thinned character into strokes using structural features like endpoint, cross-point, junction points, and thinning. They classified the segmented shapes or strokes as left curve, right curve, horizontal stroke, vertical stroke, slanted lines *etc.*, They used tree and fuzzy classifier and obtained average 86.4% accuracy.

Vamvakas *et al.*, (2007) described the statistical and structural features they have used in their approach of Greek handwritten character recognition. The statistical features they have used are zoning, projections and profiling, and crossing and distances features. In direction features they used directional histograms of contour and skeleton images. In addition to normal profile features they described in and out the profiles of contour of images. The structural features they have depicted are end point, crossing point, loop, horizontal and vertical projections histogram, radial histogram, radial out-in and in-out histogram.

Sarabjit *et al.* (2007) described projections based statistical approach for handwritten character recognition. They proposed four sided projections of characters and projections were smoothed by polygon approximation.

Araki *et al.* (2008) proposed a statistical approach for character recognition using Bayesian Filter. They reported good recognition performance in spite of the simplicity of Bayesian algorithm.

Arora *et al.* (2008) used intersection shadow features. Chain code histograms and straight line fitting features and weighted majority voting techniques for combining the classification decision obtained from different Multilayer Perception (MLP) based classifier. They also used chain code histogram and moment based features to recognize handwritten Devnagari characters. Chain code was generated by detecting the direction of the next in-line pixel in the scaled contour image. Moment features were extracted from scaled and thinned character image.

1.5.3 Work on Gurmukhi Character Recognition

In particular to Gurmukhi script, earliest major contributors founded are Chandan Singh and G. S. Lehal. They proposed Gurmukhi script recognition system (Lehal and Singh, 2000). Later they develop a complete machine printed Gurmukhi OCR system

(Lehal and Singh, 2006). Some of their other research works related to Gurmukhi Script recognition is done (Lehal and Singh, 2002) in which they proposed feature extraction, classification and post-processing approaches for Gurmukhi scripts.

Sukhpreet Singh *et al.* (2012) In this manuscript handwritten Gurmukhi character recognition for isolated characters is proposed. We have used Gabor Filter based method for feature extraction. Our database consists of 200 samples of each of basic 35 characters of the Gurmukhi script collected from different writers. These samples are pre-processed and normalized to 32*32 sizes. The highest accuracy obtained is 94.29% as 5-fold cross validation of whole database with SVM classifier used with RBF kernel.

Sharma *et al.* (2008) and Sharma *et al.* (2009) have presented the implementation of three approaches: elastic matching technique, small line segments and HMM based technique, to recognize online handwritten Gurmukhi character and reported 90.8%, 94.59% and 91.59% recognition accuracies respectively. In elastic machine technique, first they recognized the strokes and then evaluated the character based on the strokes.

DATA COLLECTION, PRE-PROCESSING, AND FEATURE EXTRACTION

The phases of the recognition process have been discussed in Chapter 1. In this chapter we will elaborate each phase and discuss about the algorithm used in this thesis work.

2.1 Data Capturing

The characters are inputted with the help of mouse or digital pen (stylus) for writing on a writing pad. The time sequential signal from the movement of the pen on the writing pad captures the dynamic position of the pen with the help of the sensors attached to the pad, means it captures the sequence of consecutive points on the x-y plane in the form of coordinates.

The example of a sample stroke is given in Figure 2.1. A stroke $S = \{(X_0, Y_0), (X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)\}$ is expressed by the time sequence of co-ordinates as shown in Figure 2.1 represents the generation of data for a sample of a Gurmukhi character in which (X_0, Y_0) and (X_n, Y_n) denote the initial and the final point of the generated time sequence of co-ordinates for the particular sample.

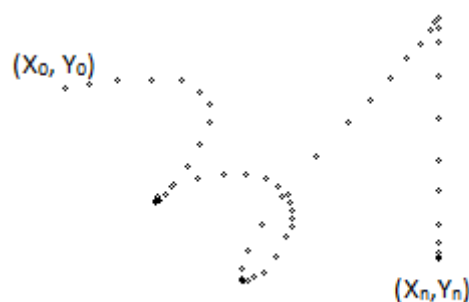


Figure 2.1: A sample stroke of Gurmukhi character

A stroke is defined as the sequence of sample points occurring between consecutive PenDown and PenUp transitions. A character could possibly be a uni-stroke character or formed from two or more strokes. In Gurmukhi Script also single stroke as well as multi-stroke both types of characters are present. In Figure 2.2 we have Gurmukhi Character ‘ੴ’. We can observe that ‘ੴ’ is a 3-stroke character.

Speed of writing directly affects the recognition process. When the speed of writing is slow, the sample points are located densely, whereas quick writing produces sparsely located points. One can note that the speed of writing typically slows down on sharp corners, in the beginning of the stroke and at the end of the stroke. The selection of a suitable level of sampling rate and resolution depends on the writing speed and the scale of the meaningful pen trace features. If the sampling rate is too slow, odd corners will be introduced to the sampled pen trace and some of the real corners and miniscule trace features can be missed.

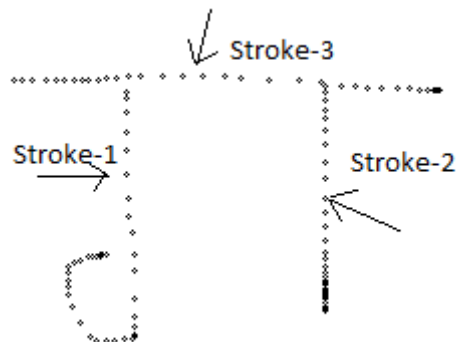


Figure 2.2: A sample Gurmukhi character of 3 strokes.

2.2 Pre-processing Phase

The input stroke contains noise and distortion due to hardware and software limitations. As discussed in chapter 1, the issues in online handwriting in stroke has to be removed before so as to extract efficient feature which can be proven as the higher accuracy of recognition. Pre-processing phase for Gurmukhi character recognition consists of following stages are:

2.2.1 Removal of Duplicate Points

Sometimes while writing, the stroke is overlapped due to the shape of the stroke or the character, and hence while writing the pixels captured contains duplicate points, as shown in Figure 2.3, which has to be removed for getting equal spaced resampling points.

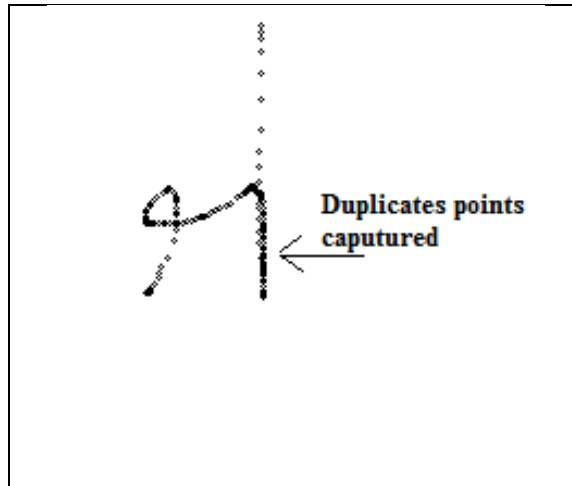


Figure 2.3: A sample stroke showing duplicate or overlapped points

Algorithm 2.1 for removing duplicate points:

1. Store the captured X and Y co-ordinate value captured in different array.
2. Set i = total number of pixels captured
3. Set $j = 0$ and $k = 0$.
4. Repeat step 4 until $j < i$ and increment j by 1.
5. Set $k = j + 1$ and repeat step 6-8 until $k < i - 1$ and increment k by 1.
6. If $(X[j] == X[k] \text{ and } Y[j] == Y[k])$.
7. Delete $X[k]$ and $Y[k]$ from the array.
8. End if.

2.2.2 Size Normalization and Centring

The size of the input stroke depends on the movement of a pen, the writing pad and also on the writer. Strokes are not of the same size, when we move the pen again and again. Size normalization normalizes every stroke to the same size and centered to a constant frame and places the text at a fixed distance from the origin. This can be achieved by comparing the input character border frame or taking the

maximum value of the co-ordinates of x-axis and y-axis with an already assumed fixed frame and the input character can be moved along with an assumed center location for centring.

The algorithm implemented in this work fixes the size of an input character to a fixed size of 300x300 pixels and places in the center of the fixed frame. Thus, if a stroke is entered, less than 300x300 size. It is transformed into a fixed size of 300x300 pixels. Similarly if the stroke entered, is greater than the size of 300x300 then also it will be transformed to a size of fixed size 300x300. This can be seen in the Figure 2.4 where in the first case, the smaller size stroke is entered and in the second window after transformation normalized and centered stroke is shown.

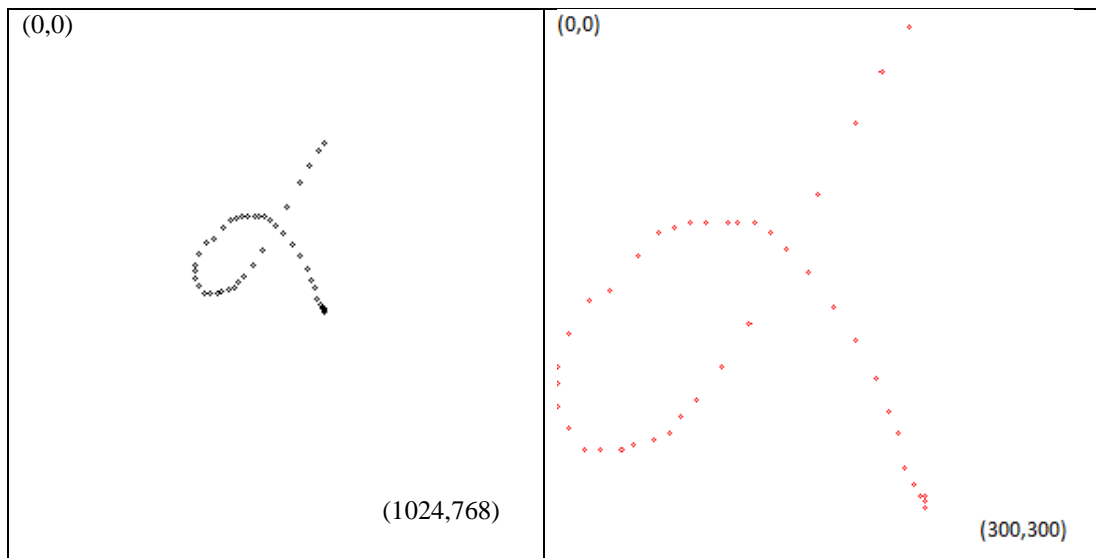


Figure 2.4 Handwriting character after normalization of a sample Gurmukhi character (ka)

Algorithm 2.2 Normalization and centring

1. Set i = total no of pixel or points captured.
2. Set $frame-x = 300$ and $frame-y = 300$ (for scaling to a 300x300 area).
3. Find the maximum of X and Y co-ordinates values, as $x-max$ and $y-max$ and minimum of co-ordinates value, as $x-min$ and $y-min$
4. Calculate $S_x = frame-x / (xmax - xmin)$
5. Calculate $S_y = frame-y / (ymax - ymin)$
6. Take $\min(S_x, S_y)$ as scaling factor, $fact$.
7. Repeat step 8 to 10 for all the points captured.
8. $P_x[i] = P_x[i] * fact$, where P_x is the X co-ordinate point.

9. $P_y[i] = P_y[i] * fact$, where P_y is the y co-ordinate point.
10. $i = i - 1$.

This algorithm normalizes the stroke in a fixed size of 300x300 pixels.

Centring

1. Set $i =$ total no of pixel or points captured.
2. Find the min of co-ordinates value of x-axis and y-axis and set as $xmin$ and $ymin$ respectively.
3. Repeat step 4 to 6 untill $i = 0$,
4. $P_x[i] = P_x[i] - xmin$, where P_x is the x co-ordinate point.
5. $P_y[i] = P_y[i] - ymin$, where P_y is the y co-ordinate point.
6. $i = i - 1$.

2.2.3 Missing Point Interpolation

When a character is drawn with high speed, the writing pad sensors may miss some point to capture as hardware issue or it can be a software problem the interface algorithm cannot capture a high speed writing input stroke. These missing points can be interpolated using various techniques such as Bezier and B-Spline curve interpolation. In this thesis work, cubic Bezier interpolation has been used. In this piecewise interpolation technique, a set of consecutive four points is considered for obtaining the curve. The next four points give the next Bezier curve and so on.

Bezier interpolation is applied between the points where the distance is greater than one. First, the algorithm finds the difference between the pair of consecutive points and calls Bezier where the distance between the points is greater than one. Interpolation is done between those points whose distance is greater than only by making a Bezier curve. Figure 2.5 shows how interpolation of missing points is done. Figure consists of three columns where first column shows the inputted stroke second column shows the normalized and third column shows the stroke after interpolation.

Algorithm 2.3 Interpolating Missing Points.

1. Create an empty list L for storing the points generated.
2. Set $t =$ number of strokes in the list and set $k = 1$.

3. Repeat step 4 for each stroke k , until $k \leq t$.
4.
 - i) Calculate I as the total number of points in the current stroke.
 - ii) If $(I \geq 4)$ then GOTO step 5
5. Set $P = 1$
6. Repeat step 7 until $P = I-4$
7. Repeat step 8 until distance $(P_i, P_{i+1}) > 1$
8. Call Bezier $(P_i, P_{i+1}, P_{i+2}, P_{i+3})$
 - i) Update List L by incorporating the new points as the consecutive points obtained through Bezier Function.
 - ii) Set $P = P+1$
9. Exit

Function Bezier $(P_i, P_{i+1}, P_{i+2}, P_{i+3})$

1. u is a variable such that $0 \leq u \leq 1$.
2. Set $u = 0.1$.
3. Repeat steps 4 and 5 until $u \leq 1$.
4. Calculate x co-ordinate of the new point as

$$(1-u)^3 P_{ix} + 3u(1-u)^2 P_{(i+1)x} + 3(1-u)u^2 P_{(i+2)x} + u^3 P_{(i+3)x}$$
 And similarly calculate Y co-ordinate

$$(1-u)^3 P_{iy} + 3u(1-u)^2 P_{(i+1)y} + 3(1-u)u^2 P_{(i+2)y} + u^3 P_{(i+3)y}$$
5. $u = u + 0.1$.
6. Return.

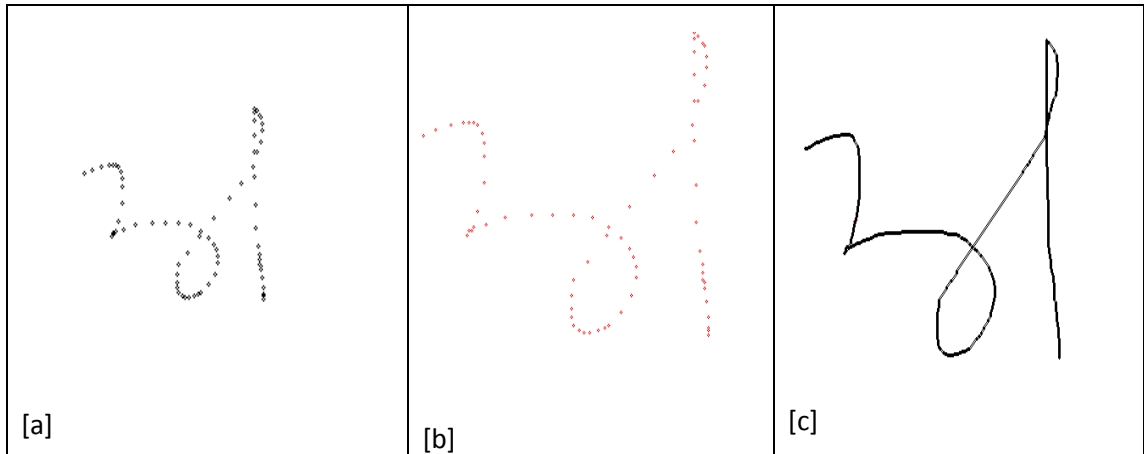


Figure 2.5: A sample stroke of Gurmukhi showing (a) -Inputted stroke (b) normalized stroke (c) Interpolated stroke

2.2.4 Resampling

Resampling of points is done to fix the number of points in the input handwritten character and also retain the original shape of the character. A filter is applied such that only a fixed number of points are selected. According to the literature survey, the results we obtained said that the best results can be obtained when the numbers of resampled points are 64. The points are selected in such a way that the shape of the character will not get affected. Since, the distance between the two points after interpolation is less than one and a half, removal of points shall include two options: remove all points between the pair of points having distance less than one, or, remove points at constant distances, *i.e.*, two or three and so on. Figure 3.4 illustrates the resampling done on the input handwritten character.

Algorithm 2.4 Resampling

1. For all the points in the list,
 - If (distance between the two points) > 1.5 then
 - Call algorithm (interpolation missing points)
 - End if
2. For all the points calculate P as a total no of points in a stroke.
3. For getting the 64 points calculate $X = P/64$
4. If X is absolute then
 - For all the points of a stroke in a list, remove points at X distance with respect to total number of points in the stroke

End if

5. Else take $X = (\text{int}) X$ and $Y = X + 1$

Two linear equations will be made

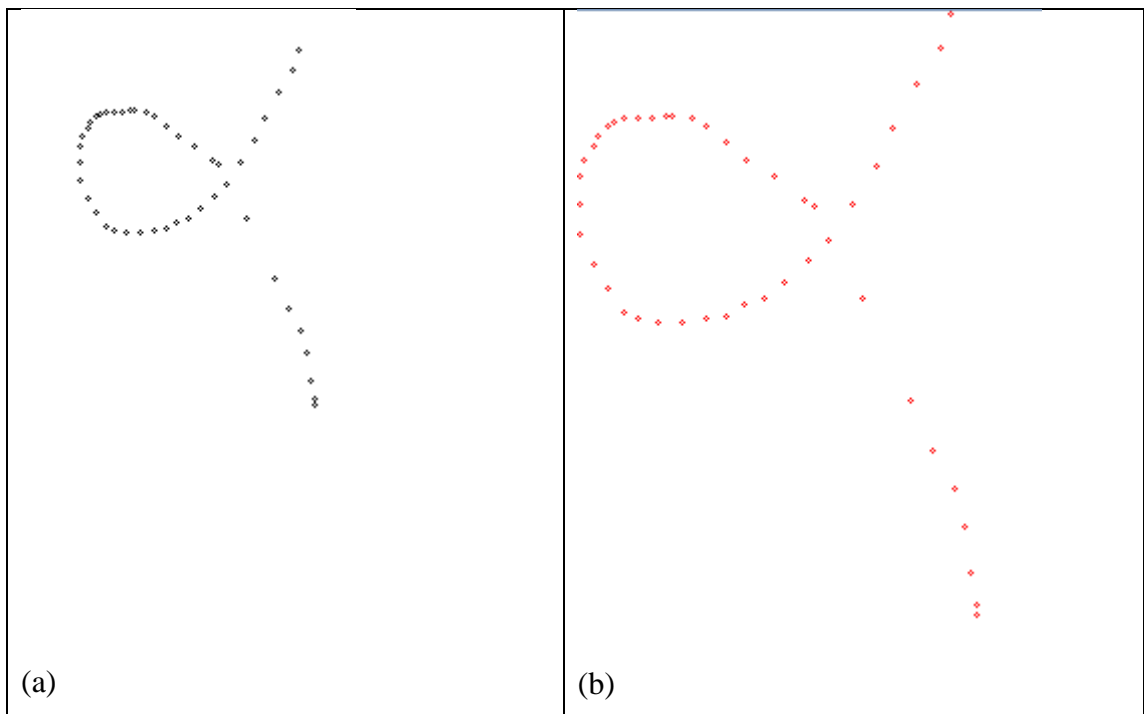
$$X * a + Y * b = P \text{ (Total no of points in a stroke)}$$

$$a + b = 64$$

Find the value of a and b

For all the points of a stroke, remove the points at X distance ' a ' number of times and Y distance ' b ' number of times.

End else.



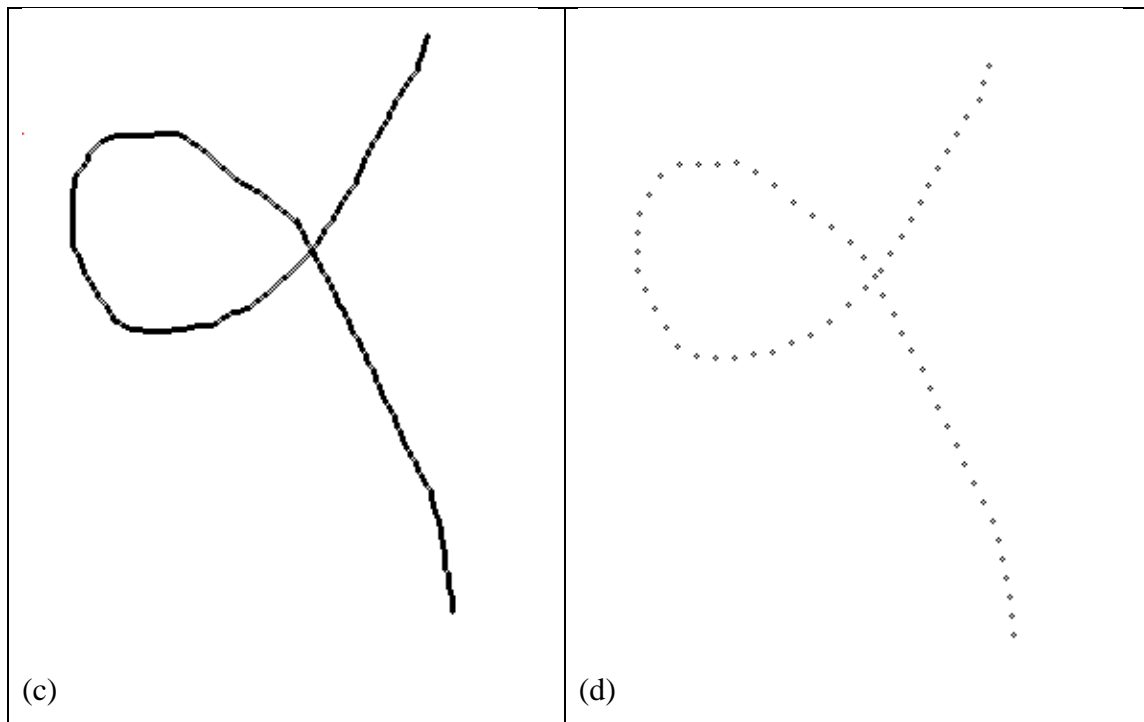


Figure 2.6: Shows the entire algorithm (a) input stroke by the writer (b) normalized and centering (c) interpolation missing points (d) resampling of points.

2.3 Feature Extraction

Feature extraction is a very important step and recognition accuracy largely depends upon the features extracted. Points generated after pre-processing phase is used as a feature for recognition. The number of points for each stroke will be fixed for each stroke and are equidistant as far as possible. In this thesis work the number of points taken is 64. The extracted feature, *i.e.*, the x and y co-ordinates of each stroke are storing it in a format of a used classifier for recognition as it will be passed to a classifier which tells about the class of these points they belongs to in recognition phase.

GURMUKHI CHARACTER RECOGNITION AND POST-PROCESSING

In this chapter the recognition phase of the online handwritten recognition system is discussed. Before recognition of the stroke, it is first, pre-processed as discussed in chapter 2. The pre-processed handwritten stroke points are then recognized using Support Vector Machine (SVM) which returns stroke-id. Introduction to SVM is discussed in section 3.1, the Recognition phase is discussed in section 3.2 and post-processing part is discussed in section 3.3.

3.1 Introduction to SVM

Support Vector Machine (SVM) is a popular classification tool used for pattern recognition and other classification purposes. Support Vector Machines are a group of supervised learning methods that can be applied to classification or regression. The standard SVM classifier takes the set of input data and predicts to classify them into one of the two distinct classes. SVM classifier is trained by a given set of training data and a model is prepared to classify test data based upon this model. For multiclass classification problem, we decompose multiclass based into multiple binary classification problems, and we design suitable combined multiple binary SVM classifier.

Support Vector Machines are based on the concept of decision planes that defines decision boundaries. A decision plane is one that separates between a set of objects having different class membership. The Figure 3.1 shows the classification of objects having class one of the two: either triangle or diamond. The separating line defines a boundary on the right side of which all objects are diamond and to the left of which all objects are triangulated. Any new object falling to the right is labelled, *i.e.*, classified, as diamond (or classified as triangle if its fall to the left of the separating line).

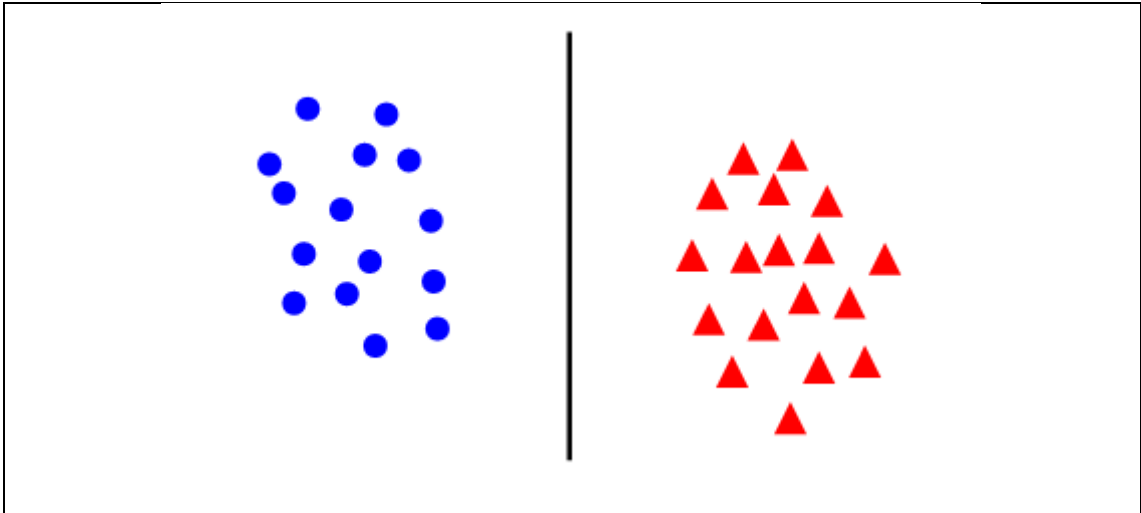


Figure 3.1: Linear classification of objects by SVM into two classes

The above Figure is a classic example of a linear classifier, *i.e.*, a classifier that separates a set of objects into their respective classes (diamond and triangle in this case) with a line. Most classification tasks, however, are not that simple, and often more complex structures are needed in order to make an optimal separation, *i.e.*, correctly classify new objects (test case) on the basis of the examples that are available (train case). This complex separation is represented by a complex function rather than linear function and such classification of objects is shown in Figure 3.2 (a). Compared to the previous schematic, it is clear that a full separation of the diamond and triangle objects would require a curve (which is more complex than a line).

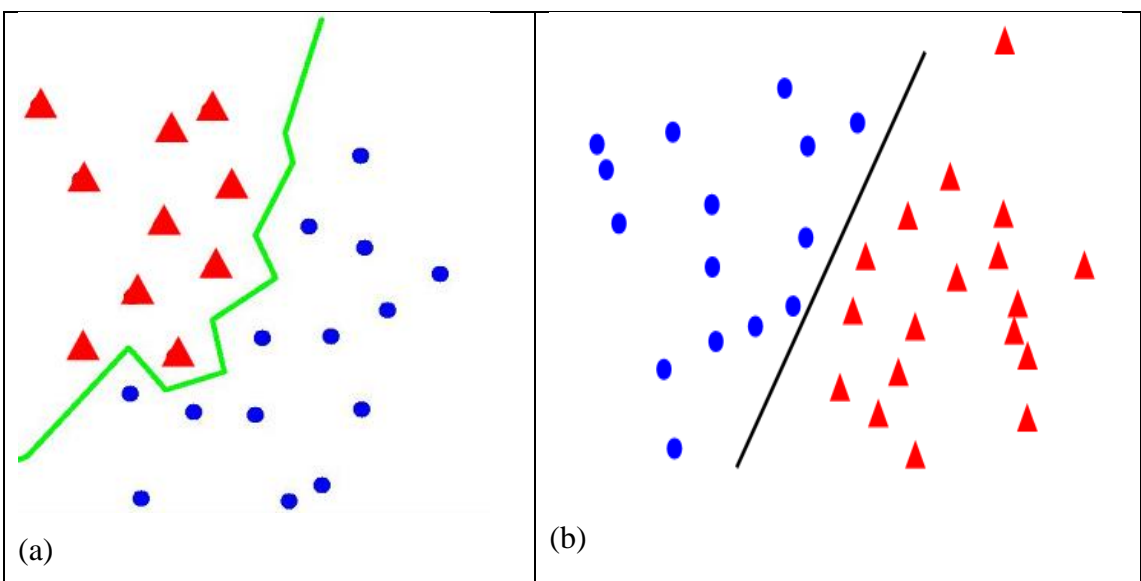


Figure 3.2 Mapping of SVM classifications from complex to linear: (a) classification by a complex kernel function (b) mapping of input space of Figure (a) into feature space

Classification tasks based on drawing separating to distinguish between objects of different class membership are known as hyperplane classifier. It linearly separates any two classes by finding a maximum margin between the two classes. The margin means the minimal distance from the separating hyperplane to the closest data points. SVM learning machines search for an optimal separating hyperplane where the margin is maximal. The outcome of the SVM is based only on the data points that are at the margin and are called ‘Support Vectors’. This concept is illustrated in Figure 3.3.

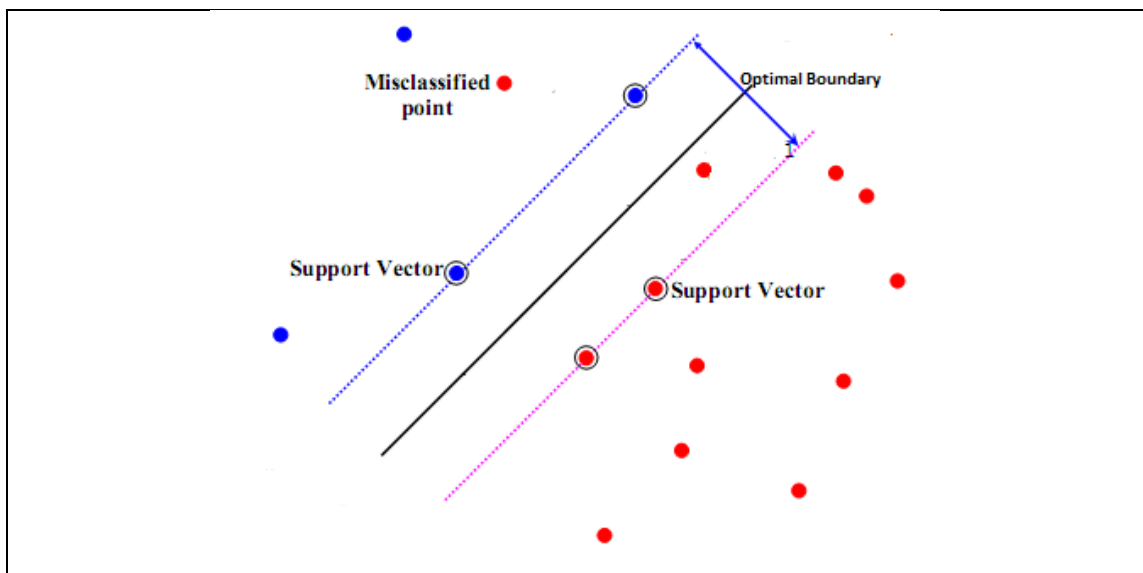


Figure 3.3: Concept of support vectors

The black line in Figure 3.3 determines the mapping done from the input space to the feature space in SVM.

Support Vector Machine classifier has gained immense popularity in recent years providing excellent recognition results in various applications. Some of the applications of SVM are:

- Text (and hypertext) categorization.
- Image classification
- Bioinformatics (protein classification, cancer classification)
- Topic drift in page ranking algorithms

- Handwriting character recognition

SVM in the present dissertation has been used for handwriting character recognition.

3.2 Recognition of Gurmukhi Characters

As discussed earlier, a Gurmukhi character can be written in one stroke or more than one stroke, it depends on the user. A character of Gurmukhi is divided into three zones, *i.e.*, lower zone, middle zone and upper zone, for recognition, stroke id is given to every possible stroke that can make a Gurmukhi character. Ids given to a particular stroke is shown in Table 3.1 for lower zone, Table 3.2 for upper zone and Table 3.3 for middle zone.

With the help of electronic devices and tool developed in C#. Net, samples, which are given in the below tables of each zone of all the strokes are captured. In total 11 users, at an average of 10 has written each stroke of every class. Stroke ids given in the tables has 77 strokes in the middle zone, 12 strokes in upper zone and 7 strokes in the lower zone. 100 samples of each stroke for all the zones is captured and the co-ordinates of every sample stroke, *i.e.*, X, Y co-ordinates are stored in a text file, called as raw data, in a specified format which is discussed in the next section.

Table 3.1: Lower zone stroke ids



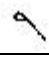
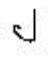



S.NO	STROKE ID	SYMBOL	DESCRIPTION
1)	101		Onkar
2)	102		Bindi
3)	103		Ra
4)	104		Ha
5)	105		Va
6)	106		Dulenkar
7)	107		Halant

Table 3.2: Upper zone stroke ids

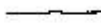


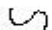






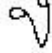




S.NO	STROKE ID	SYMBOL	DESCRIPTION
1)	121		Upper bar
2)	122		Lawan
3)	123		Dalawan
4)	124		Hora
5)	125		Ghanoura bar
6)	126		Adhak
7)	127		Bindi
8)	128		Tippi
9)	129		Associate bar
10)	130		Dalawan
11)	131		Dalawan
12)	132		Lawan

Table 3.3: Middle zone stroke ids

S.NO	STROKE ID	SYMBOL	DESCRIPTION
1)	141		Urah
2)	142		Urah
3)	143		Urah

4)	144	ᄃᄃ	Aarah
5)	145	ᄃᄃ	Aarah
6)	146	ᄃ	Eedi
7)	147	ᄃ	Eedi
8)	148	ᄃ	Eedi
9)	149	ᄃ	Eedi
10)	150	ᄃ	Eedi
11)	151	ᄃ	Sassa
12)	152	ᄃ	Sassa
13)	153	ᄃ	Sassa
14)	154	ᄃ	Sassa
15)	155	ᄃ	Haaha

ᄃ

16)	156	३	Haaha
17)	157	२	Kakka
18)	158	२	Kakka
19)	159	५	Khakha, Pappa
20)	160	५	Khakha,
21)	161	२	Khakha, Pappa
22)	162		Vertical Bar
23)	163	—	Khakha, Dhadha
24)	164	४	Gagga, Rara
25)	165	४	Gagga, Rara
26)	166	३	Ghagga

27)	167	୧	Ghagga
28)	168	୧	Aeyaa
29)	169	୧	Aeyaa
30)	170	୧	Chachaa
31)	171	୧	Chachaa
32)	172	୧	Chhachha
33)	173	୧	Chhachha
34)	174	୧	Jajja
35)	175	୧	Jajja
36)	176	୧	Jhajja
37)	177	୧	Jhajja
38)	178	୧	Jhajja

39)	179	᳚	Neyaa, Vawa
40)	180	᳛	Neyaa, Vawa
41)	181	᳜	Vawa
42)	182	᳝	Tenqa
43)	183	᳞	Thathaa
44)	184	᳟	Thathaa
45)	185	᳠	Dadda
46)	186	᳡	Dadda
47)	187	᳢	Dhadhaa
48)	188	᳣	Dhadhaa
49)	189	᳤	Nahnaa

50)	190	ሐ	Nahnaa
51)	191	ሠ	Tatta
52)	192	ሡ	Tatta
53)	193	ሢ	Dadhaa
54)	194	ሣ	Dadhaa
55)	195	ሤ	Dadhaa
56)	196	ሥ	Nanhaa
57)	197	ሦ	Nanhaa
58)	198	ሧ	Nanhaa
59)	199	ረ	Half Vertical bar
60)	200	ሪ	Faffa
61)	201	ራ	Faffa

62)	202	ब	Babba
63)	203	भ	Babba
64)	204	ब	Bhabhaa
65)	205	ब	Bhabhaa
66)	206	म	Mamma
67)	207	य	Yaeya
68)	208	य	Yaeya
69)	209	य	Yaeya
70)	210	ल	Lalla
71)	211	ल	Lalla
72)	212	र	Rara

73)	213	।	Rara
74)	214	੨	Rara
75)	215	੫	Cehari
76)	216	੯	Behari
77)	217	੪	Mamma
78)	218	/	Jhajja
79)	219	੨	Rara
80)	220	੫	Cehari
81)	222	੫	Jajja
82)	223	੫	Jajja

For recognition of a Gurmukhi character, which is in the form of strokes, first pre-processed and 64 resemble points are recorded, *i.e.*, 64 *X* co-ordinates and 64 *Y* co-ordinates. These points are stored in a specified format of the classifier. Here we have used SVM classifier, therefore, we have stored the resampled points in SVM format. The format of the SVM is given as below:

**Stroke ID: 1:X co-ordinate 2:Y co-ordinate 3:X co-ordinate 4:Y co-ordinate,...,
127:X co-ordinate 128:Y co-ordinate.**

Hence, all the 100 samples of each stroke are firstly pre-processed and the result is stored in the SVM format as shown in Figure 3.4. As discussed in earlier section the SVM needs to be trained with a set of training data, in the recognition, which gives model file as an output. Before training the training data is first scaled in a lower and upper bound. As discussed earlier, after pre-processing of raw data, co-ordinates values range from (0, 0) to (300, 300), these co-ordinates was scaled with a lower limit of 1 and upper limit of 9, using the tool of LIBSVM.

```

lower_new_30 - Notepad
101 1:0 2:0 3:4 4:0 5:8 6:0 7:12 8:0 9:19 10:0 11:22 12:0 13:27 14:0 15:32 16:0 17:35 18:0 19:41 20:0 21:44 22:0 23:49 24:0 25:52 26:0 27:55 28:0 29:59 30:0 31:64 32:0
33:69 34:0 35:73 36:0 37:78 38:0 39:82 40:0 41:87 42:0 43:91 44:0 45:96 46:0 47:100 48:0 49:105 50:0 51:108 52:0 53:114 54:0 55:117 56:0 57:123 58:0 59:126 60:0 61:132
62:0 63:135 64:0 65:141 66:0 67:144 68:0 69:150 70:0 71:153 72:0 73:159 74:0 75:162 76:0 77:168 78:0 79:171 80:0 81:177 82:0 83:182 84:0 85:188 86:0 87:195 88:0 89:200
90:0 91:206 92:0 93:212 94:0 95:217 96:0 97:224 98:0 99:230 100:0 101:235 102:0 103:242 104:0 105:248 106:0 107:253 108:0 109:260 110:0 111:266 112:0 113:271 114:0
115:276 116:0 117:280 118:0 119:284 120:0 121:288 122:0 123:292 124:0 125:296 126:0 127:300 128:0
101 1:0 2:0 3:6 4:35 5:11 6:34 7:19 8:33 9:24 10:33 11:29 12:31 13:34 14:32 15:39 16:30 17:44 18:30 19:50 20:28 21:54 22:27 23:59 24:27 25:66 26:27 27:71 28:27 29:78
30:27 31:84 32:27 33:89 34:27 35:96 36:27 37:100 38:27 39:106 40:26 41:113 42:26 43:117 44:26 45:121 46:24 47:124 48:20 49:130 50:18 51:134 52:18 53:137 54:22 55:137
56:26 57:141 58:22 59:144 60:18 61:151 62:18 63:157 64:20 65:162 66:24 67:169 68:27 69:173 70:27 71:179 72:26 73:186 74:26 75:191 76:27 77:193 78:23 79:196 80:19
81:200 82:17 83:200 84:13 85:200 86:9 87:206 88:9 89:211 90:8 91:216 92:8 93:222 94:8 95:228 96:9 97:228 98:5 99:228 100:1 101:233 102:0 103:239 104:0 105:244 106:0
107:251 108:0 109:257 110:0 111:262 112:0 113:269 114:0 115:274 116:0 117:278 118:0 119:282 120:0 121:286 122:0 123:290 124:0 125:295 126:0 127:300 128:0
101 1:0 2:0 3:5 4:1 5:8 6:1 7:12 8:3 9:15 10:3 11:18 12:5 13:22 14:6 15:27 16:8 17:31 18:9 19:35 20:12 21:41 22:12 23:44 24:12 25:49 26:12 27:53 28:12 29:57 30:12
31:62 32:12 33:66 34:12 35:70 36:12 37:75 38:12 39:78 40:12 41:83 42:12 43:88 44:12 45:92 46:12 47:97 48:12 49:100 50:12 51:104 52:12 53:110 54:12 55:113 56:12 57:118
58:12 59:123 60:12 61:129 62:12 63:135 64:12 65:141 66:12 67:147 68:12 69:153 70:12 71:158 72:12 73:164 74:12 75:170 76:12 77:176 78:12 79:181 80:12 81:187 82:12
83:193 84:12 85:199 86:12 87:204 88:12 89:210 90:12 91:216 92:12 93:222 94:12 95:227 96:12 97:233 98:12 99:239 100:12 101:245 102:12 103:250 104:12 105:256 106:12
107:262 108:12 109:266 110:11 111:272 112:11 113:277 114:11 115:281 116:11 117:285 118:11 119:286 120:8 121:287 122:3 123:291 124:0 125:295 126:0 127:300 128:0
101 1:0 2:0 3:5 4:0 5:9 6:0 7:13 8:0 9:17 10:0 11:20 12:0 13:23 14:0 15:27 16:0 17:33 18:0 19:38 20:0 21:41 22:0 23:45 24:0 25:49 26:0 27:55 28:0 29:61 30:0 31:65 32:0
33:69 34:0 35:74 36:0 37:78 38:0 39:82 40:0 41:88 42:0 43:91 44:0 45:96 46:0 47:99 48:0 49:103 50:0 51:106 52:0 53:109 54:0 55:113 56:0 57:118 58:0 59:122 60:0 61:126
62:0 63:131 64:0 65:135 66:0 67:141 68:0 69:147 70:0 71:153 72:0 73:159 74:0 75:164 76:0 77:171 78:0 79:176 80:0 81:183 82:0 83:188 84:0 85:195 86:0 87:200 88:0 89:206
90:0 91:211 92:0 93:217 94:0 95:223 96:0 97:228 98:0 99:232 100:0 101:237 102:0 103:243 104:0 105:248 106:0 107:254 108:0 109:260 110:0 111:265 112:0 113:269 114:0
115:273 116:0 117:278 118:0 119:282 120:0 121:286 122:0 123:291 124:0 125:295 126:0 127:300 128:0
101 1:0 2:0 3:5 4:0 5:12 6:0 7:17 8:1 9:22 10:1 11:27 12:2 13:32 14:2 15:37 16:2 17:46 18:2 19:48 20:2 21:44 22:2 23:42 24:4 25:48 26:4 27:53 28:4 29:58 30:4 31:63
32:4 33:72 34:6 35:76 36:6 37:80 38:6 39:86 40:6 41:93 42:8 43:98 44:8 45:106 46:8 47:111 48:9 49:118 50:9 51:124 52:10 53:129 54:11 55:133 56:12 57:138 58:12 59:145
60:12 61:150 62:15 63:155 64:15 65:162 66:15 67:167 68:15 69:174 70:15 71:180 72:15 73:185 74:15 75:192 76:15 77:197 78:15 79:204 80:15 81:210 82:15 83:215 84:15
85:222 86:15 87:226 88:14 89:232 90:14 91:238 92:14 93:242 94:13 95:249 96:13 97:255 98:15 99:255 100:11 101:255 102:7 103:255 104:3 105:256 106:0 107:260 108:0
109:264 110:0 111:268 112:0 113:271 114:3 115:275 116:6 117:280 118:10 119:285 120:15 121:285 122:9 123:289 124:5 125:293 126:3 127:300 128:0
101 1:0 2:0 3:3 4:0 5:7 6:0 7:13 8:0 9:16 10:0 11:20 12:0 13:25 14:0 15:28 16:0 17:33 18:0 19:37 20:0 21:40 22:0 23:47 24:0 25:52 26:0 27:58 28:0 29:64 30:0 31:68 32:0
33:72 34:0 35:77 36:0 37:85 38:0 39:90 40:0 41:95 42:0 43:100 44:0 45:109 46:0 47:115 48:0 49:120 50:0 51:128 52:0 53:133 54:0 55:140 56:0 57:145 58:0 59:151 60:0
61:156 62:0 63:160 64:0 65:165 66:0 67:173 68:0 69:177 70:0 71:184 72:0 73:189 74:0 75:193 76:0 77:198 78:0 79:203 80:0 81:209 82:0 83:213 84:0 85:218 86:0 87:225 88:0
89:229 90:0 91:234 92:0 93:240 94:0 95:245 96:0 97:249 98:0 99:254 100:0 101:261 102:0 103:265 104:0 105:269 106:0 107:273 108:0 109:278 110:0 111:283 112:0 113:287
114:0 115:291 116:0 117:296 118:0 119:295 120:0 121:292 122:4 123:289 124:8 125:286 126:13 127:282 128:10
101 1:0 2:15 3:4 4:15 5:10 6:15 7:17 8:15 9:22 10:16 11:28 12:16 13:33 14:17 15:38 16:17 17:43 18:17 19:48 20:17 21:56 22:17 23:63 24:17 25:63 26:17 27:59 28:17 29:55
30:17 31:55 32:19 33:62 34:19 35:66 36:19 37:72 38:19 39:77 40:19 41:82 42:19 43:87 44:19 45:94 46:21 47:99 48:21 49:104 50:21 51:109 52:21 53:116 54:21 55:122 56:21
57:122 58:21 59:121 60:23 61:126 62:23 63:131 64:23 65:139 66:23 67:143 68:24 69:151 70:24 71:156 72:25 73:163 74:25 75:168 76:26 77:173 78:26 79:178 80:27 81:184
82:27 83:191 84:27 85:196 86:30 87:202 88:30 89:208 90:30 91:212 92:29 93:219 94:29 95:225 96:30 97:229 98:29 99:236 100:28 101:239 102:27 103:243 104:24 105:247
106:21 107:251 108:16 109:256 110:15 111:260 112:15 113:265 114:15 115:270 116:15 117:271 118:11 119:275 120:9 121:279 122:8 123:284 124:6 125:291 126:4 127:300 128:0
101 1:0 2:48 3:4 4:46 5:9 6:45 7:14 8:45 9:19 10:43 11:23 12:42 13:28 14:42 15:34 16:41 17:40 18:41 19:38 20:41 21:37 22:42 23:42 24:40 25:47 26:39 27:52 28:39 29:61
30:40 31:67 32:39 33:72 34:37 35:78 36:36 37:83 38:36 39:88 40:34 41:93 42:35 43:98 44:33 45:103 46:34 47:108 48:32 49:112 50:32 51:117 52:32 53:125 54:32 55:130 56:32

```

Figure 3.4: Pre-processed points in SVM format

The scaled file was used for training the SVM machine. The scaling can be done with any lower and upper limit depends on any data set values. In this thesis, to scale the Data, lower limit of 1 and upper limit of 9 is taken. The scaled file is shown in Figure 3.5. The training data for training the SVM should be accurate to get correct recognition results, thus respective stroke-id of stroke is given, captured at the time of data collection.

```

lower_pre_Testing.txt - Notepad
File Edit Format View Help
101 1:1 2:3.27586 3:1.17897 4:2.8125 5:1.43636 6:2.44 7:1.64249 8:2.20231 9:1.78553 10:2.03993 11:1.91733 12:1.93458 13:2.29609 14:1.73059 15:2.89333 16:1.65753 17:3.2
109:2.25333 110:1.58667 111:1.98667 112:1.58667 113:1.72 114:1.58667 115:1.45333 116:1.58667 117:1.37333 118:1.69333 119:1.74667 120:2.04 121:2.04 122:2.33333 123:2.30
101 1:1 2:1.82759 3:1.21477 4:1.175 5:1.45714 6:1.64 7:1.70466 8:1.55491 9:1.88889 10:1.49485 11:2.17333 12:1.4486 13:2.49721 14:1.43836 15:3.02667 16:1.43836 17:3.3466
41333 102:1.66667 103:4.17333 104:1.66667 105:3.93333 106:1.66667 107:3.69333 108:1.66667 109:3.45333 110:1.66667 111:3.21333 112:1.66667 113:2.97333 114:1.66667 115:2
101 1:1.46532 2:1.89655 3:1.73378 4:1.625 5:2.1013 6:1.48 7:2.36788 8:1.36994 9:2.65375 10:1.3299 11:3.00533 12:1.2243 13:3.36872 14:1.12665 15:4.25333 16:1.18265 17:4
7 114:1.42667 115:3.45333 116:1.34667 117:3.74667 118:1.34667 119:4.17333 120:1.34667 121:4.49333 122:1.34667 123:4.92 124:1.24 125:5.26667 126:1.13333 127:5.88 128:1
101 1:1 2:1 3:1.26846 4:1.0625 5:1.66494 6:1.16 7:1.97409 8:1.23121 9:1.323 10:1.24742 11:2.68533 12:1.33645 13:3.12291 14:1.32877 15:3.93333 16:1.43836 17:4.30667 18:
06:1 107:5.37333 108:1 109:5.77333 110:1.05333 111:6.2 112:1.48 113:6.62667 114:1.48 115:7.05333 116:1.48 117:7.45333 118:1.48 119:7.24 120:1.8 121:6.78667 122:1.88 12
101 1:1 2:1 3:1.14318 4:1.0625 5:1.39481 6:1.10667 7:1.58031 8:1.23121 9:1.8062 10:1.20619 11:2.06667 12:1.26168 13:2.31844 14:1.32877 15:2.84 16:1.47489 17:3.13333 18
1.96 105:4.70667 106:1.96 107:4.70667 108:1.96 109:4.94667 110:1.96 111:5.26667 112:2.01333 113:5.69333 114:2.12 115:6.44 116:2.28 117:6.70667 118:2.22667 119:7 120:2
101 1:2.18121 2:1 3:2.52125 4:1 5:3.16104 6:1 7:3.54922 8:1 9:3.91473 10:1 11:4.43467 12:1 13:5 14:1 15:6.30667 16:1 17:6.89333 18:1 19:7.42667 20:1 21:7.93333 22:1 23
:1.16 123:5.77333 124:1.21333 125:6.52 126:1.24 127:7 128:1.24
101 1:1 2:1 3:1.42953 4:1 5:2.01818 6:1 7:2.53368 8:1.09249 9:3.00517 10:1.08247 11:3.624 12:1.14953 13:4.26257 14:1.18265 15:5.61333 16:1.21918 17:6.25333 18:1.29224
3333 103:1.64 104:1.50667 105:1.96 106:1.69333 107:2.65333 108:1.77333 109:3.34667 110:1.77333 111:4.04 112:1.77333 113:4.76 114:1.77333 115:5.42667 116:1.77333 117:6
101 1:1 2:1 3:1.21477 4:1 5:1.54026 6:1 7:1.84974 8:1 9:2.17829 10:1 11:2.49333 12:1 13:2.89944 14:1 15:3.64 16:1 17:3.96 18:1 19:4.38667 20:1 21:4.70667 22:1 23:5.08
17:7.32 118:1.32 119:7.61333 120:1.32 121:8.04 122:1.32 123:8.36 124:1.32 125:8.68 126:1.18667 127:9 128:1
101 1:1 2:2.03448 3:1.21477 4:1.175 5:1.54026 6:1.69333 7:1.80829 8:1.46243 9:2.03359 10:1.37113 11:2.36533 12:1.33645 13:2.65363 14:1.29224 15:3.42667 16:1.29224 17:3
3 108:1.4 109:5.34667 110:1.4 111:5.64 112:1.4 113:5.93333 114:1.4 115:6.25333 116:1.4 117:6.62667 118:1.4 119:7 120:1.4 121:7.26667 122:1.66667 123:7.56 124:1.96 125:
101 1:1 2:1 3:1.19687 4:1 5:1.47792 6:1 7:1.70466 8:1.04624 9:1.9509 10:1.08247 11:2.23733 12:1.07477 13:2.5419 14:1.07306 15:3.10667 16:1.07306 17:3.50667 18:1.07306
50667 104:1.88 105:5.29333 106:1.88 107:5.00 108:1.88 109:4.73333 110:1.88 111:4.36 112:1.88 113:3.93333 114:1.88 115:3.74667 116:1.88 117:3.98667 118:1.88 119:4.30667
101 1:6.38782 2:6.65517 3:6.63758 4:6.125 5:7.56623 6:5.37333 7:7.56995 8:4.79191 9:7.57364 10:4.38144 11:7.80533 12:4.06542 13:8.15884 14:3.99543 15:1.16 16:1 17:1 18:1
101 1:1 2:1 3:1.05369 4:1 5:1.14545 6:1 7:1.26943 8:1 9:1.33075 10:1 11:1.42667 12:1 13:1.55866 14:1 15:1.74667 16:1 17:1.88 18:1 19:1.98667 20:1 21:2.06667 22:1 23:2
101 1:1 2:2.24138 3:1.07159 4:1 5:1.16623 6:1 7:1.26943 8:1.83237 9:1.39276 10:1.74227 11:1.49067 12:1.6257 13:1.6257 14:1.65753 15:1.85333 16:1.65753 17:1.96 18:1
1.48 105:7.98667 106:1.48 107:8.09333 108:1.48 109:8.2 110:1.48 111:8.30667 112:1.48 113:8.44 114:1.48 115:8.57333 116:1.48 117:8.68 118:1.48 119:8.78667 120:1.48 121:
101 1:7.1745 2:5.62069 3:7.33557 4:5.1875 5:8.37662 6:4.57333 7:8.37824 8:4.09827 9:8.37984 10:3.76289 11:8.63733 12:3.50467 13:9 14:3.41096 15:1 16:1 17:1 18:1 19:1 2
7 99:8.06667 100:1.90667 101:8.17333 102:1.88 103:8.36 104:1.88 105:8.49333 106:1.88 107:8.6 108:1.88 109:8.73333 110:1.85333 111:8.86667 112:1.85333 113:8.86667 114:1
101 1:1 2:1 3:1.07159 4:1 5:1.16623 6:1 7:1.2487 8:1 9:1.33075 10:1 11:1.448 12:1 13:1.53631 14:1 15:1.74667 16:1 17:1.93333 18:1 19:2.01333 20:1 21:2.2 22:1 23:2.28 2
101 1:1 2:1 3:1.05369 4:1 5:1.14545 6:1.05333 7:1.2487 8:1.04624 9:1.35142 10:1.12371 11:1.448 12:1.11215 13:1.6257 14:1.10959 15:1.88 16:1.18265 17:2.01333 18:1.18265
100:1.69333 101:7.34667 102:1.69333 103:7.50667 104:1.69333 105:7.64 106:1.69333 107:7.82667 108:1.69333 109:7.96 110:1.69333 111:8.06667 112:1.69333 113:8.17333 114:1
101 1:1 2:1 3:1.07159 4:1 5:1.16623 6:1.05333 7:1.26943 8:1.09249 9:1.39276 10:1.16495 11:1.49067 12:1.14953 13:1.6257 14:1.14612 15:1.85333 16:1.14612 17:1.96 18:1.21
108:1.4 109:7.93333 110:1.4 111:8.04 112:1.4 113:8.14667 114:1.4 115:8.28 116:1.4 117:8.38667 118:1.4 119:8.52 120:1.4 121:8.65333 122:1.4 123:8.76 124:1.4 125:8.86667
101 1:1 2:2.17241 3:1.10738 4:2.0625 5:1.20779 6:1.96 7:1.31088 8:1.87861 9:1.47545 10:1.78351 11:1.55467 12:1.71028 13:1.58101 14:1.76712 15:1.82667 16:1.80365 17:2.0
98:1.45333 99:7.88 100:1.45333 101:8.04 102:1.45333 103:8.2 104:1.42667 105:8.41333 106:1.42667 107:8.54667 108:1.42667 109:8.73333 110:1.4 111:8.89333 112:1.37333 113
101 1:1 2:1.82759 3:1.07159 4:1.175 5:1.16623 6:1.69333 7:1.2487 8:1.60116 9:1.37209 10:1.53608 11:1.448 12:1.56075 13:1.58101 14:1.54795 15:1.85333 16:1.54795 17:1.96
:7.90667 104:1.29333 105:8.04 106:1.32 107:8.17333 108:1.26667 109:8.30667 110:1.26667 111:8.38667 112:1.21333 113:8.49333 114:1.13333 115:8.6 116:1 117:8.73333 118:1

```

Figure 3.5: Scaled file of pre-processed points

Parameter Selection:

The effectiveness of SVM machine depends on the Parameters used, *i.e.*, the selection of SVM's kernel, kernel parameter and the soft margin parameter C. In this thesis the kernel model used is linear and the kernel parameter's default value of epsilon value is 0.01 and gamma value is 0.005.

These parameters are passed for training the SVM machine with the scaled file of training set to get a model file which is further used for recognition. There are 77 strokes in the middle zone, which are considered in this work, has 100 samples each, they are after pre-processing saved in a text file named as training.txt, means there are total $75 \times 100 = 7500$ strokes in a file and each are resampled to 64 points after processing. This file is then scaled to be in upper and lower limit and a new file as training.txt.scale is generated when training.txt file is provided as input to the tool LIBSVM. This scaled file, training.txt.scale, was used for training the SVM machine to generate a model file. The model file was generated, on passing the scaled file with the required parameters, *i.e.*, kernel selection, Gamma value and epsilon value. A sample model file shown in Figure 3.6.

```

lower_pre_Training.txt.scale - Notepad
File Edit Format View Help
svm_type c_svc
kernel_type rbf
gamma 0.005
nr_class 7
total_sv 182
rho 0.465968 -0.588906 0.276966 0.363347 -0.0298383 -0.238845 -0.767241 -0.313047 -0.124932 -0.467855 -0.60702 0.751116 0.716808 0.467619 0.292787 0.263574 -0.171634
-0.389149 -0.390019 -0.562812 -0.212855
label 101 104 102 107 103 105 106
nr_sv 21 33 12 42 36 26 12
SV
0 0 0 0 0.02274184991538448 1:1 2:1 3:1.06667 4:1 5:1.13763 6:1 7:1.2069 8:1 9:1.32829 10:1 11:1.38013 12:1 13:1.63905 14:1 15:1.85619 16:1 17:1.93645 18:1
19:2.09333 20:1 21:2.17333 22:1 23:2.30667 24:1 25:2.3913 26:1 27:2.47157 28:1 29:2.57333 30:1 31:2.70667 32:1 33:2.84 34:1 35:2.94667 36:1 37:3.09396 38:1 39:3.18667
40:1 41:3.32 42:1 43:3.50172 44:1 45:3.7331 46:1 47:3.8169 48:1 49:3.83784 50:1 51:3.96907 52:1 53:4.05017 54:1 55:4.15152 56:1 57:4.31313 58:1 59:4.5122 60:1
61:4.86813 62:1 63:5.15385 64:1 65:5.19331 66:1 67:5.15884 68:1 69:5.33213 70:1 71:5.26481 72:1 73:5.24 74:1 75:5.32 76:1 77:5.48 78:1 79:5.56 80:1 81:5.72 82:1
83:5.85333 84:1 85:6.01333 86:1 87:6.2 88:1 89:6.33333 90:1 91:6.49333 92:1 93:6.65333 94:1 95:6.78667 96:1 97:6.97333 98:1 99:7.13333 100:1 101:7.26667 102:1
103:7.45333 104:1 105:7.61333 106:1 107:7.74667 108:1 109:7.93333 110:1 111:8.09333 112:1 113:8.22667 114:1 115:8.36 116:1 117:8.46667 118:1 119:8.57333 120:1 121:8.68
122:1 123:8.78667 124:1 125:8.89333 126:1 127:9.128:1
0.09762359782466458 0.03244392142710537 0 0 0.0556178366823569 0.03258613284284003 1:1 2:1 3:1.05 4:1 5:1.12043 6:1 7:1.22414 8:1 9:1.27646 10:1 11:1.34557 12:1
13:1.59172 14:1 15:1.74916 16:1 17:1.88294 18:1 19:1.98667 20:1 21:2.06667 22:1 23:2.25333 24:1 25:2.3913 26:1 27:2.55184 28:1 29:2.70667 30:1 31:2.81333 32:1 33:2.92
34:1 35:3.05333 36:1 37:3.28188 38:1 39:3.4 40:1 41:3.53333 42:1 43:3.74914 44:1 45:4.1032 46:1 47:4.23944 48:1 49:4.24324 50:1 51:4.5189 52:1 53:4.55853 54:1
55:4.77104 56:1 57:4.90572 58:1 59:5.20906 60:1 61:5.57143 62:1 63:5.92308 64:1 65:5.90706 66:1 67:5.99639 68:1 69:6.11191 70:1 71:6.12892 72:1 73:6.04 74:1 75:6.14667
76:1 77:6.28 78:1 79:6.41333 80:1 81:6.57333 82:1 83:6.68 84:1 85:6.81333 86:1 87:7.88 89:7.10667 90:1 91:7.24 92:1 93:7.4 94:1 95:7.53333 96:1 97:7.64 98:1
99:7.77333 100:1 101:7.96 102:1 103:8.06667 104:1 105:8.17333 106:1 107:8.28 108:1 109:8.41333 110:1 111:8.54667 112:1 113:8.65333 114:1 115:8.76 116:1 117:8.89333
118:1 119:8.86667 120:1 121:1.18667 122:1.10667 123:8.70667 124:1.21333 125:8.62667 126:1.34667 127:8.52 128:1.48
0.796848971044147 0.2052551650771105 0.6762946584505393 0.7019731300931834 0.511551590756876 0.4233960517765504 1:1 2:1 3:1.2 4:1 0.02941 5:1.44731 6:1.0613 7:1.46552
8:1.128 9:1.65659 10:1.16736 11:1.88121 12:1.24454 13:2.51479 14:1.30126 15:3.03344 16:1.3506 17:3.27425 18:1.39544 19:3.61333 20:1.40727 21:3.93333 22:1.48057 23:4.2
24:1.46259 25:4.55853 26:1.45485 27:4.90635 28:1.45485 29:5.16 30:1.45333 31:5.48 32:1.45638 33:5.82667 34:1.45485 35:6.14667 36:1.45485 37:6.44966 38:1.45333 39:6.76
40:1.45333 41:7.10667 42:1.45333 43:7.57045 44:1.45333 45:8.14591 46:1.45333 47:8.43662 48:1.45333 49:8.43243 50:1.45333 51:8.91753 52:1.45333 53:9.54 54:1.45333
55:8.81145 56:1.45333 57:8.56902 58:1.48 59:8.58188 60:1.50667 61:8.70696 62:1.50667 63:8.81538 64:1.50667 65:8.76208 66:1.56 67:8.27798 68:1.58667 69:8.81805
70:1.58667 71:7.52265 72:1.58667 73:7.18667 74:1.58667 75:7.10667 76:1.64 77:6.86667 78:1.64 79:6.62667 80:1.64 81:6.65333 82:1.69333 83:6.41333 84:1.69333 85:6.17333
86:1.74667 87:5.93333 88:1.77333 89:5.69333 90:1.82667 91:5.45333 92:1.90667 93:5.21333 94:1.90667 95:4.97333 96:1.90667 97:4.73333 98:1.90667 99:4.49333 100:1.90667
101:4.25333 102:1.90667 103:4.01333 104:1.90667 105:3.77333 106:1.90667 107:3.53333 108:1.90667 109:3.29333 110:1.90667 111:3.05333 112:1.90667 113:2.81333 114:1.88
115:2.57333 116:1.88 117:2.33333 118:1.90667 119:2.09333 120:1.85333 121:1.90667 122:1.8 123:1.8 124:1.45333 125:1.4 126:1.45333 127:1.128 128:1.45333
1 0.09380975185359319 1 0.4759611505096738 0.2081009631715083 0.2734679065782617 1:1 2:2.86572 3:1.1 4:2.94118 5:1.18925 6:3.02299 7:1.2931 8:3.112 9:1.36285
10:3.20921 11:1.44924 12:3.30568 13:1.73373 14:3.20921 15:2.07023 16:3.10359 17:2.20401 18:3.0076 19:2.33333 20:2.92 21:2.46667 22:2.86572 23:2.70667 24:2.79592
25:2.87291 26:2.76589 27:3.00669 28:2.76589 29:3.21333 30:2.76 31:3.34667 32:2.77181 33:3.53333 34:2.76589 35:3.66667 36:2.76589 37:3.84564 38:2.76 39:3.98667 40:2.76
41:4.12 42:2.76 43:4.35395 44:2.76 45:4.64413 46:2.76 47:4.77465 48:2.76 49:4.78378 50:2.76 51:4.98625 52:2.76 53:4.98662 54:2.76 55:5.14815 56:2.76 57:5.30976 58:2.76

```

Figure 3.6: Model file of pre-processed point

Now when a writer writes a stroke or character on writing pad or on an interface for recognition, the points will be pre-processed or 0.03258613284284003 and then stored in SVM format. This, pre-processed file, was scaled first, to the same parameters as it was used for scaling the training file. Then this scaled file will be passed to the SVM machine with the model file and we get the recognized stroke id which is identified by the SVM as an output.

3.3 Post-Processing

Printing of respective character, which is recognized by SVM, on the user interface of the system is done in this phase. After getting stroke-id which is recognized by the Support Vector Machine as discussed earlier in section 3.2, it is then passed to post-processing phase. This phase contains a database of the rules for creating a character (kumar and sharma, 2013). A particular character can be created though different strokes, the database contains all those rules of strokes that can make a character individually or by combination of stroke. For example stroke id 142 individually forms Ura (𑂣) and combination of strokes 141+121 also forms Ura (𑂣) as shown in Table 3.4. A character formation rule can be one or more than one is depending on the classification of stroke-id made. As shown in Table 3.4 some character has 2 rules and some has 4 rules of its formation and it can be more than ten also.

Table 3.4 shows some example of the rules that can form a character.

Character	Rule 1	Rule 2	Rule 3	Rule 4
ੳ	142	141+121		
ਅ	144	144+121	145+162+121	
ੲ	150	182+147	150+121	146+147+121
ਸ	153	217+154	151+154	151+162+121
ਹ	156	155+121		
ਕ	158	157+121		
ਖ	160	159+163+129	161+162+163	159+163
ਗ	164+154	164+162+121	165+154	165+162+121
ਘ	166+121	167+162+121	166	167+162
ਙ	169	168+121	168	

With the help of these rules, a Unicode is return by the program of post-processing and respective character of Gurmukhi is printed using the Unicode standards.

RESULTS AND DISCUSSIONS

Algorithms, discussed in Chapter-3, have been implemented on the raw data of strokes which was captured through an electronic device written by different writers. The co-ordinates of the strokes were captured and stored in a text file. On this data pre-processing algorithms were processed and saved in a text file in SVM format which is further passed to SVM for recognition of stroke-id. This data were divided in three datasets of 30, 50 and 70 for each zone, *i.e.*, lower zone, middle zone and upper zone. Dataset of 30 samples contains 30 samples of each stroke. Similar work was done for dataset of 50 and 70.

In this work, *k*-fold cross validation was used for testing. In *k*-fold cross validation we first divide the training set into *k* equal subsets, out of which 1 subset is tested by classifier SVM, and remaining subsets are used for training the SVM classifier. By cross validation, each subset of data is tested and the average percentage of correctly recognized stroke-id count is recorded.

In 3 fold testing the data set of 30 samples is divided into 3 equal datasets out of which 2 will be used for training the SVM machine and 1 will be used for testing and then accuracy for different parameters of gamma and epsilon (SVM parameters) value is recognized. Every divided dataset was tested and average accuracy was recorded. A similar process was followed for 4 and 5 fold testing.

Data was pre-processed using algorithms proposed in Agrawal, (2012) and proposed in this thesis. Testing of both result sets was performed using 3 schemes and their result is presented further using three dimensional graphs. The red color shows the resultant accuracy of algorithms proposed by Agrawal, (2012) while the green color shows the resultant accuracy of the proposed algorithms in this thesis work. In these graphs the X-axis represents epsilon value ranging from 0.1 to 0.5, the Y-axis represents the gamma value ranging from 0.001 to 0.02 and Z –axis represents the recognition accuracy in percentage.

4.1 SCHEME 1: 30 samples of lower, middle and upper zones each has been used for cross validation.

For the lower zone 30 samples of each stroke have been taken which are categorized in lower zone as described earlier. Maximum accuracy for 3, 4, 5 fold cross validation is given in Table 4.1.

Table 4.1: Cross validation accuracy of lower zone for 3, 4, and 5 cross validation for 30 samples

<i>k</i>-Fold Cross Validation (SVM)	Maximum accuracy of Algorithm proposed in Agrawal, (2012)	Maximum Accuracy of Algorithm proposed in this work
3	97.1%	95.7%
4	97.1%	95.7%
5	96.6%	96.1%

Three Dimensional graphs are shown in Figure 4.1, Figure 4.2 and Figure 4.3 for all the values of SVM parameters (gamma ranging from 0.001-0.02 and epsilon ranging from 0.1-0.5) have been drawn. These graphs show the variation between results obtained in Agrawal, (2012) and results obtained in this thesis work. It can be concluded from the graphs that the accuracy is better in Agrawal, (2012) as compared to results in this work for 3, 4 and 5 fold cross validation on a data set of 30 samples for lower zone.

It can be reasoned that the pre-processing phase contains the missing point interpolation step which needs at least 10 points to be implemented. In Agrawal, (2012) the Bindu of Gurmukhi character is bypassed and not processed but in this work it has been taken in the consideration. As a result, the accuracy rate becomes lower in the lower zone.

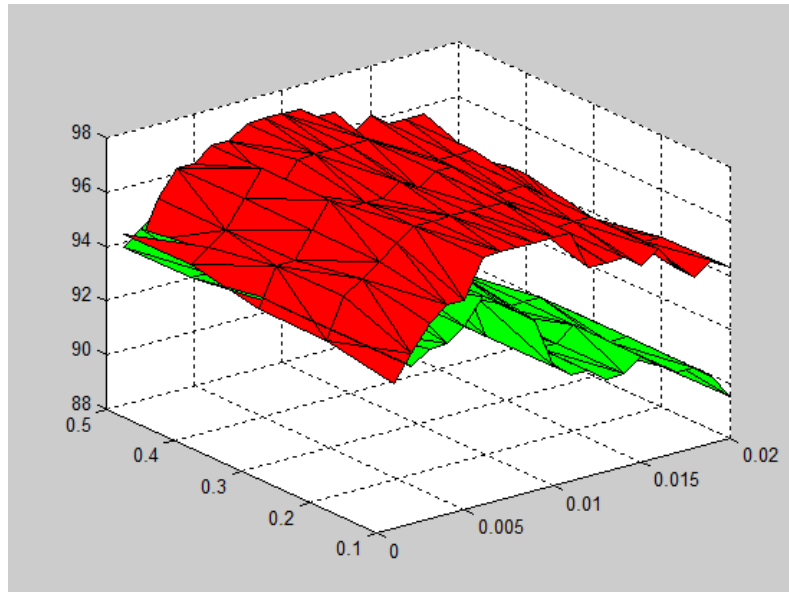


Figure 4.1: Surface graph showing variation in two methods for 3 fold cross validation on a data set of 30 samples for lower zone

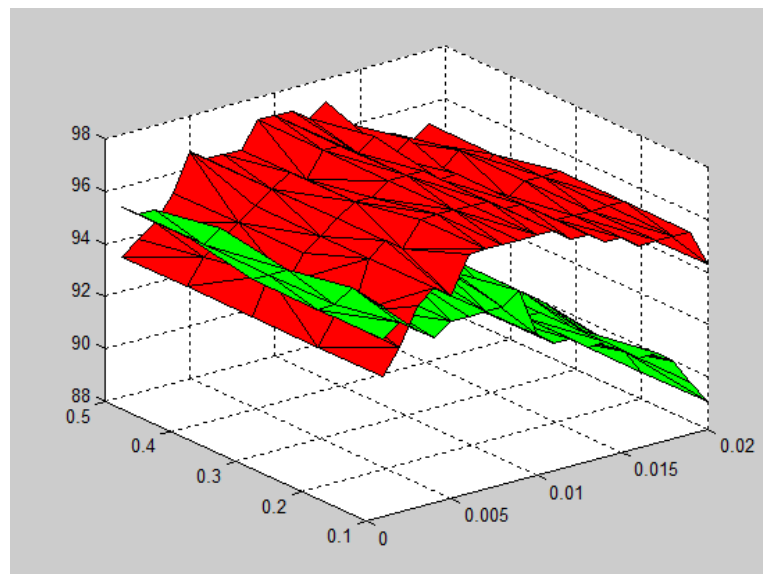


Figure 4.2: Surface graph showing variation in two methods for 4 fold cross validation on a data set of 30 samples for lower zone

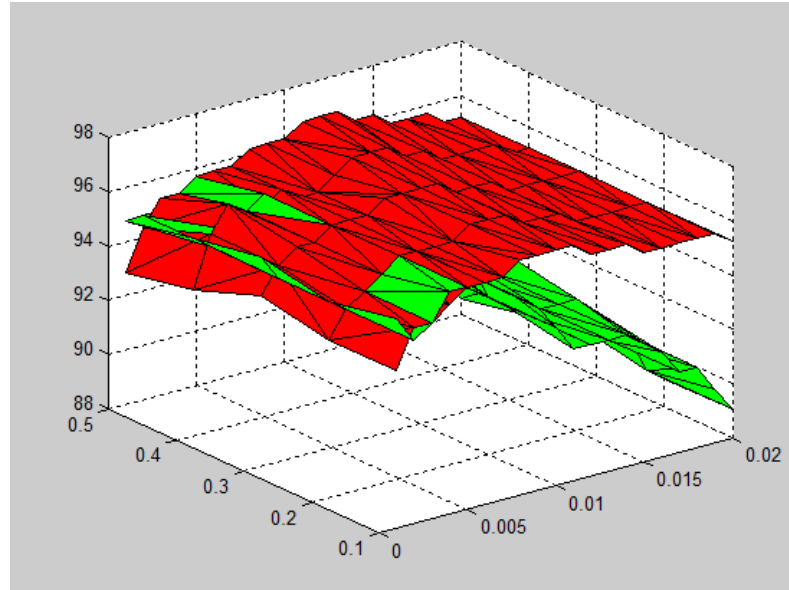


Figure 4.3: Surface graph showing variation in two methods for 5 fold cross validation on a data set of 30 samples for lower zone

For the middle zone 30 samples of each stroke have been taken. Maximum accuracy is given in the Table 5.1 for 3, 4 and 5 fold cross validation considering the parameters of the SVM.

Table 4.2: Cross validation accuracy of middle zone for 3, 4, and 5 cross validation for 30 samples.

<i>k</i>-Fold Cross Validation (SVM)	Maximum accuracy of Algorithm proposed in Agrawal, (2012)	Maximum Accuracy of Algorithm proposed in this work
3	83.9%	89.1%
4	84.1%	90.2%
5	85.2%	90.8%

Three Dimensional graphs are shown in Figure 4.4, Figure 4.5 and Figure 4.6 for all the values of SVM parameters (gamma ranging from 0.001-0.02 and epsilon ranging from 0.1-0.5) have been drawn. These graphs show the variation between results obtained in Agrawal, (2012) and results obtained in this thesis work. It can be concluded from the graphs that the accuracy is better in this work as compared to results in Agrawal, (2012) for 3, 4 and 5 fold cross validation on a data set of 30 samples for middle zone.

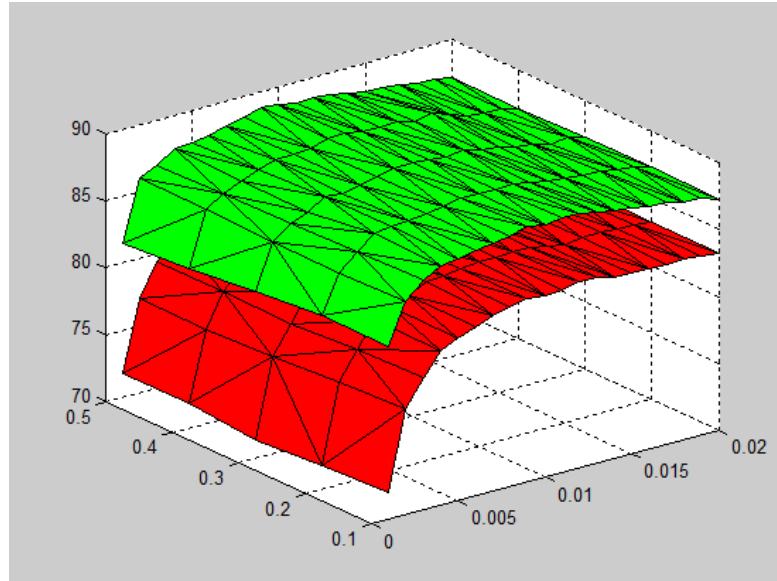


Figure 4.4: Surface graph showing variation in two methods for 3 fold cross validation on a data set of 30 samples for middle zone

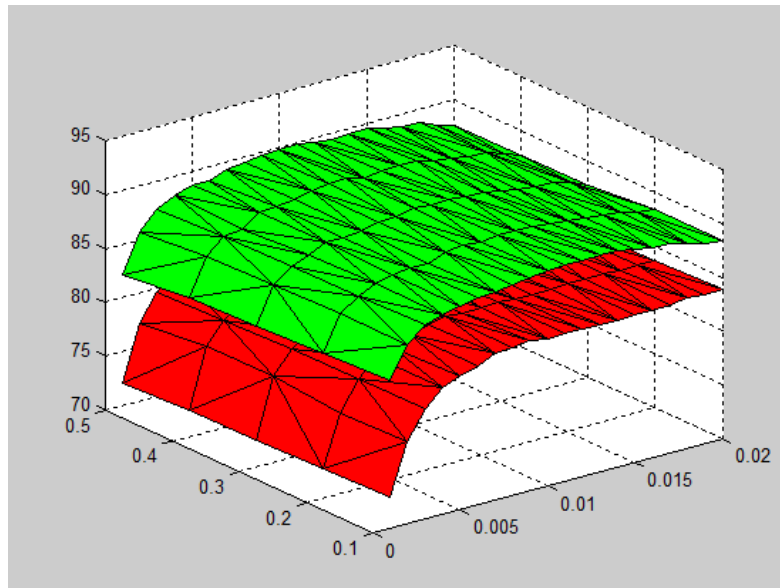


Figure 4.5: Surface graph showing variation in two methods for 4 fold cross validation on a data set of 30 samples for middle zone

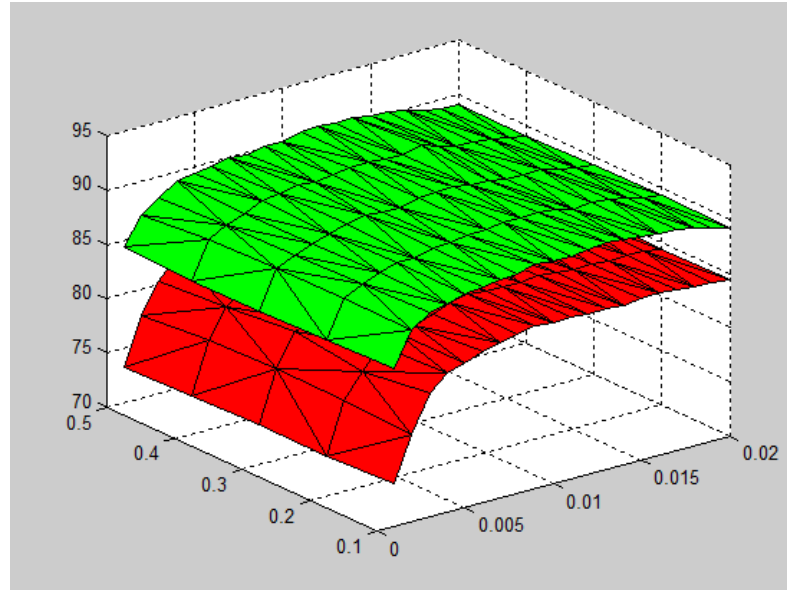


Figure 4.6: Surface graph showing variation in two methods for 5 fold cross validation on a data set of 30 samples for middle zone

In the upper zone 30 samples of each stroke have been taken. Maximum accuracy is given in the Table 5.1 for 3, 4, 5 fold cross validation.

Table 4.3: Cross validation accuracy of upper zone for 3, 4, and 5 cross validation for 30 samples

<i>k</i>-Fold Cross Validation (SVM)	Maximum accuracy of Algorithm proposed in Agrawal, (2012)	Maximum Accuracy of Algorithm proposed in this work
3	87.3%	88.1%
4	86.4%	88.8%
5	86.2%	88.6%

Three Dimensional graphs are shown in Figure 4.7, Figure 4.8 and Figure 4.9 for all the values of SVM parameters (gamma ranging from 0.001-0.02 and epsilon ranging from 0.1-0.5) have been drawn. These graphs show the variation between results obtained in Agrawal, (2012) and results obtained in this thesis work. It can be concluded from the graphs that the accuracy is better in this work as compared to results in Agrawal, (2012) for 3, 4 and 5 fold cross validation on a data set of 30 samples for upper zone.

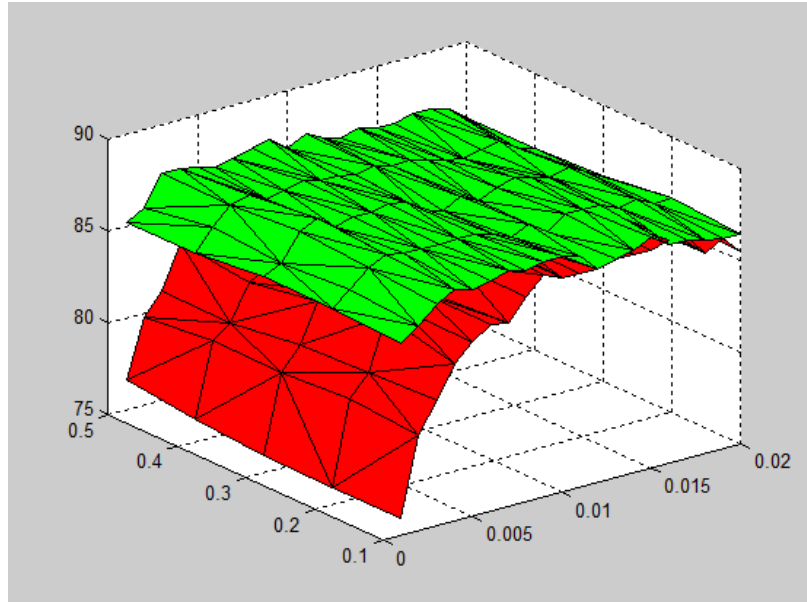


Figure 4.7: Surface graph showing variation in two methods for 3 fold cross validation on a dataset of 30 samples for upper zone

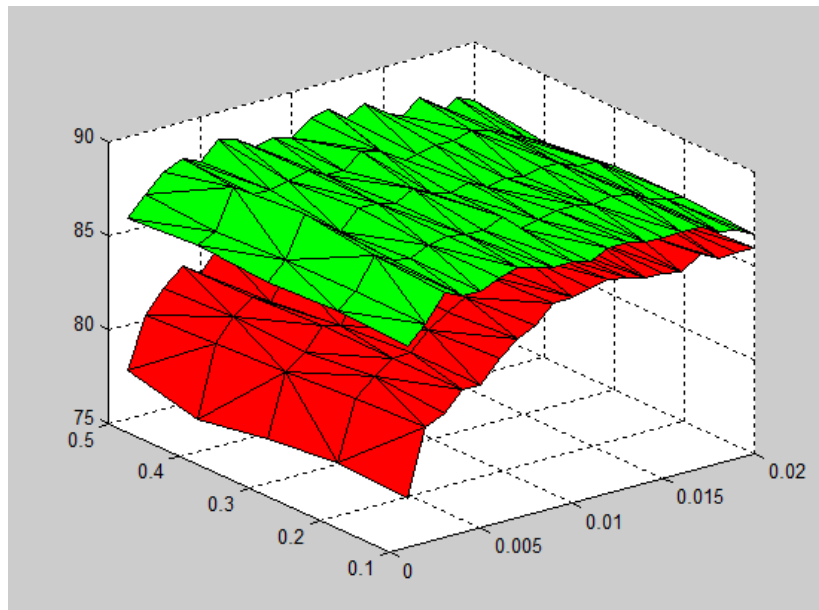


Figure 4.8: Surface graph showing variation in two methods for 4 fold cross validation on a dataset of 30 samples for upper zone

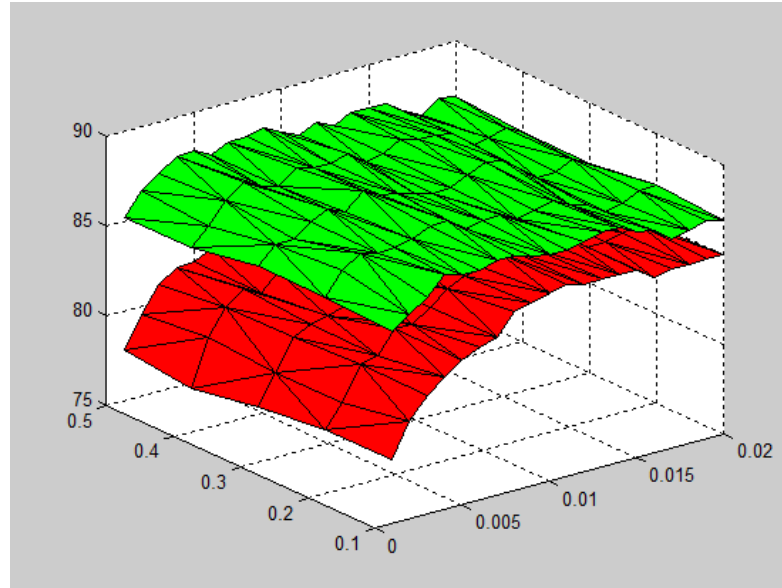


Figure 4.9: Surface graph showing variation in two methods for 5 fold cross validation on a dataset of 30 samples for upper zone

4.2 SCHEME 2: 50 samples in lower, middle and upper zones each has been taken for the cross validation. And the results are shown below.

For the lower zone 50 samples of each stroke have been taken. Maximum accuracy is given in the Table 5.1 for 3, 4, 5 fold cross validation.

Table 4.4: Cross validation accuracy of lower zone for 3, 4, and 5 cross validation for 50 samples

<i>k</i>-Fold Cross Validation (SVM)	Maximum accuracy of Algorithm proposed in Agrawal, (2012)	Maximum Accuracy of Algorithm proposed in this work
3	92.3%	96.3%
4	93.7%	96.3%
5	93.4%	96.3%

Three Dimensional graphs are shown in Figure 4.10, Figure 4.11 and Figure 4.12 for all the values of SVM parameters (gamma ranging from 0.001-0.02 and epsilon ranging from 0.1-0.5) have been drawn. These graphs show the variation between results obtained in Agrawal, (2012) and results obtained in this thesis work. It can be concluded from the graphs that the accuracy is better in this work as compared to

results in Agrawal, (2012) for 3, 4 and 5 fold cross validation on a data set of 50 samples for lower zone.

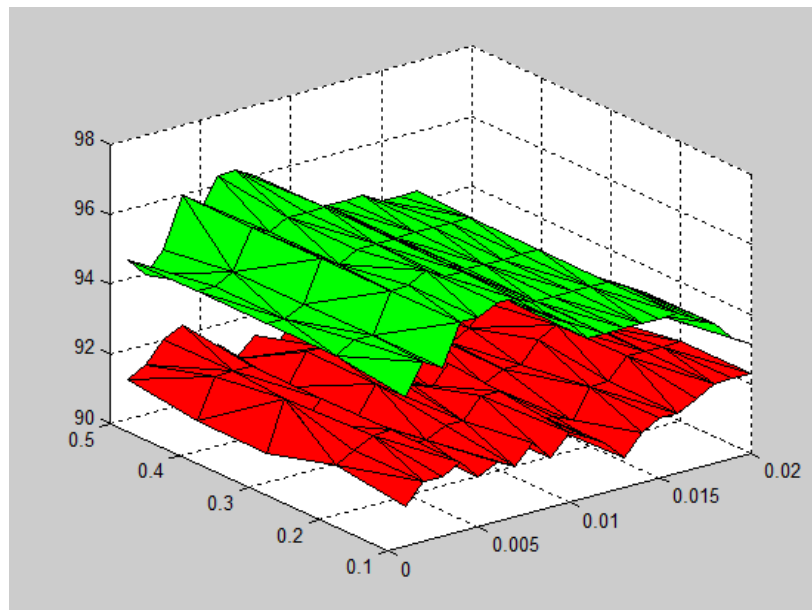


Figure 4.10: Surface graph showing variation in two methods for 3 fold cross validation on a dataset of 50 samples for lower zone

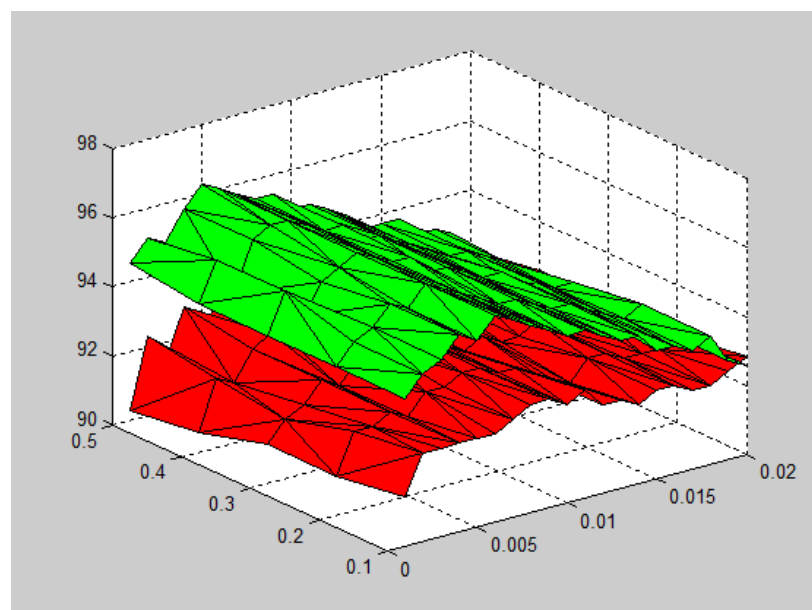


Figure 4.11: Surface graph showing variation in two methods for 4 fold cross validation on a dataset of 50 samples for lower zone

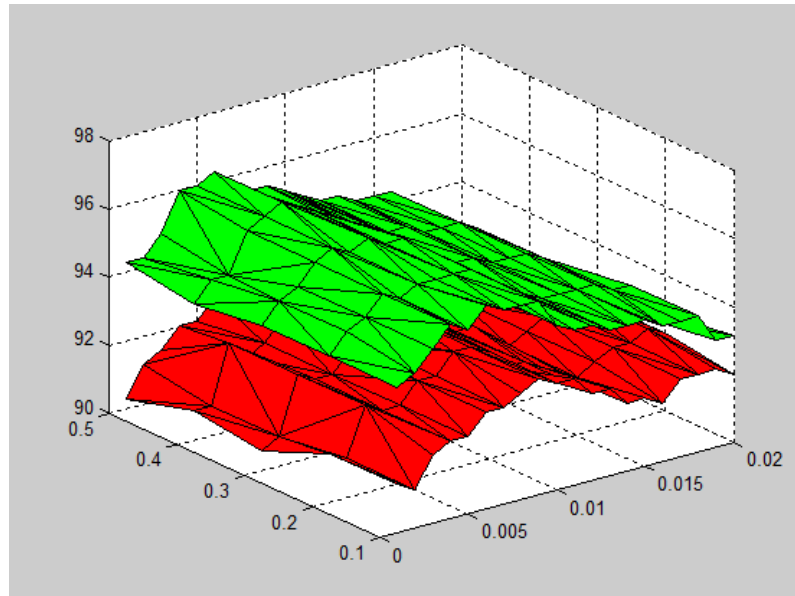


Figure 4.12: Surface graph showing variation in two methods for 5 fold cross validation on a dataset of 50 samples for lower zone

For the middle zone 50 samples of each stroke have been taken. Maximum accuracy is given in the Table 5.1 for 3, 4, 5 fold cross validation.

Table 4.5: Cross validation accuracy of middle zone for 3, 4, and 5 cross validation for 50 samples

<i>k</i>-Fold Cross Validation (SVM)	Maximum accuracy of Algorithm proposed in Agrawal, (2012)	Maximum Accuracy of Algorithm proposed in this work
3	87.2%	91.8%
4	87.5%	91.9%
5	88.5%	92.6%

Three Dimensional graphs are shown in Figure 4.13, Figure 4.14 and Figure 4.15 for all the values of SVM parameters (gamma ranging from 0.001-0.02 and epsilon ranging from 0.1-0.5) have been drawn. These graphs show the variation between results obtained in Agrawal, (2012) and results obtained in this thesis work. It can be concluded from the graphs that the accuracy is better in this work as compared to results in Agrawal, (2012) for 3, 4 and 5 fold cross validation on a data set of 50 samples for middle zone.

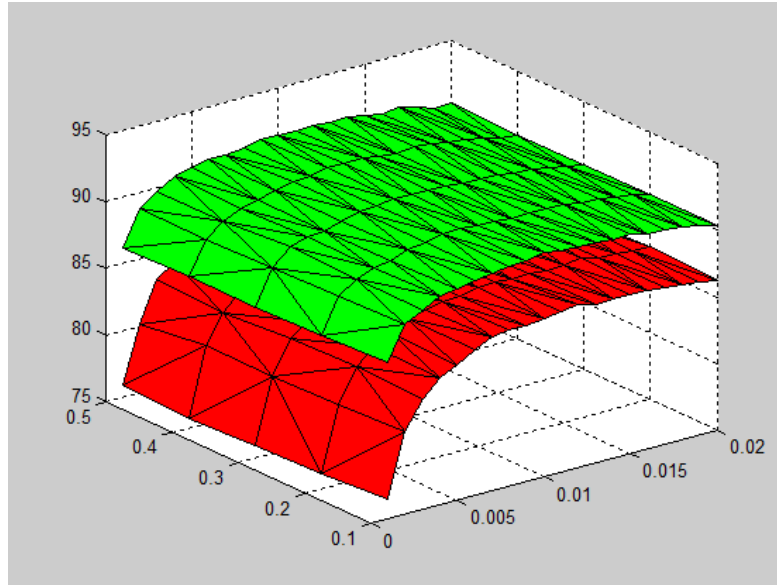


Figure 4.13: Surface graph showing variation in two methods for 3 fold cross validation on a dataset of 50 samples for middle zone

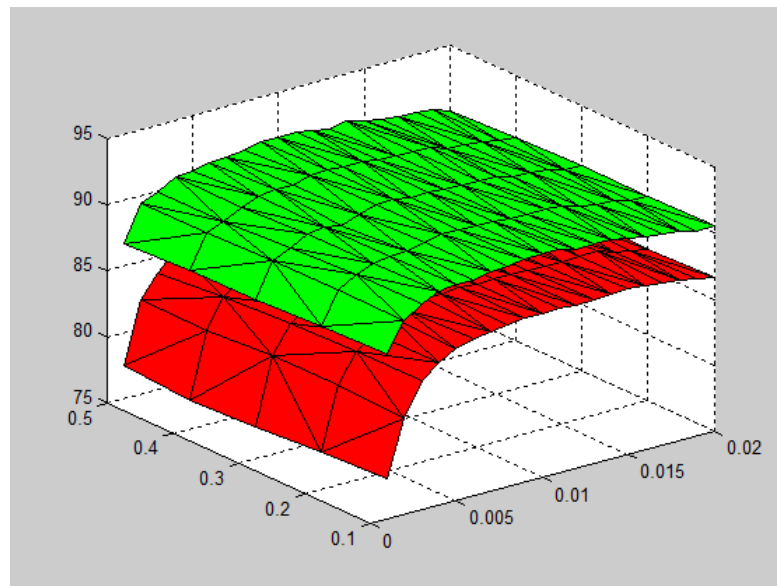


Figure 4.14: Surface graph showing variation in two methods for 4 fold cross validation on a dataset of 50 samples for middle zone

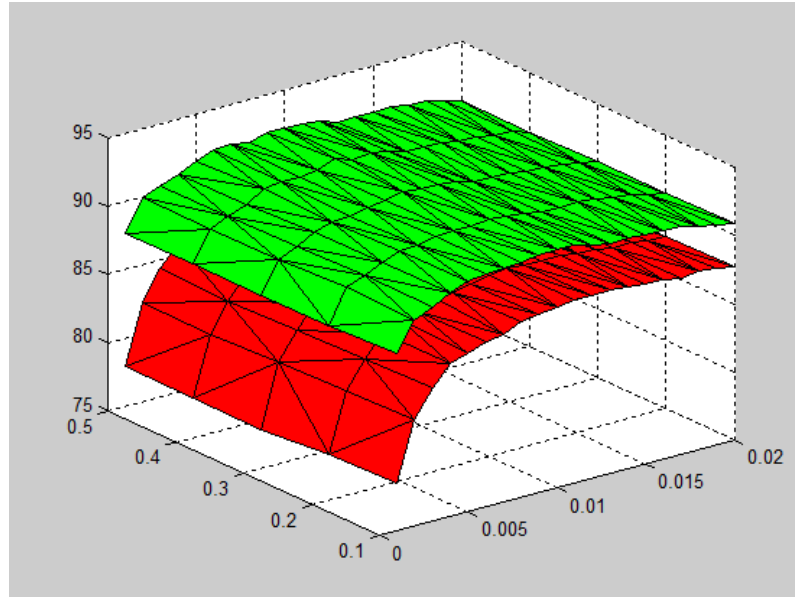


Figure 4.15: Surface graph showing variation in two methods for 5 fold cross validation on a dataset of 50 samples for middle zone

In the upper zone 50 samples of each stroke have been taken. Maximum accuracy is given in the Table 5.1 for 3, 4, 5 fold cross validation.

Table 4.6: Cross validation accuracy of upper zone for 3, 4, and 5 cross validation for 50 samples

<i>k</i>-Fold Cross Validation (SVM)	Maximum accuracy of Algorithm proposed in Agrawal, (2012)	Maximum Accuracy of Algorithm proposed in this work
3	85.7%	86.2%
4	85.4%	87.7%
5	85.3%	87.6%

Three Dimensional graphs are shown in Figure 4.16, Figure 4.17 and Figure 4.18 for all the values of SVM parameters (gamma ranging from 0.001-0.02 and epsilon ranging from 0.1-0.5) have been drawn. These graphs show the variation between results obtained in Agrawal, (2012) and results obtained in this thesis work. It can be concluded from the graphs that the accuracy is better in this work as compared to results in Agrawal, (2012) for 3, 4 and 5 fold cross validation on a data set of 50 samples for upper zone.

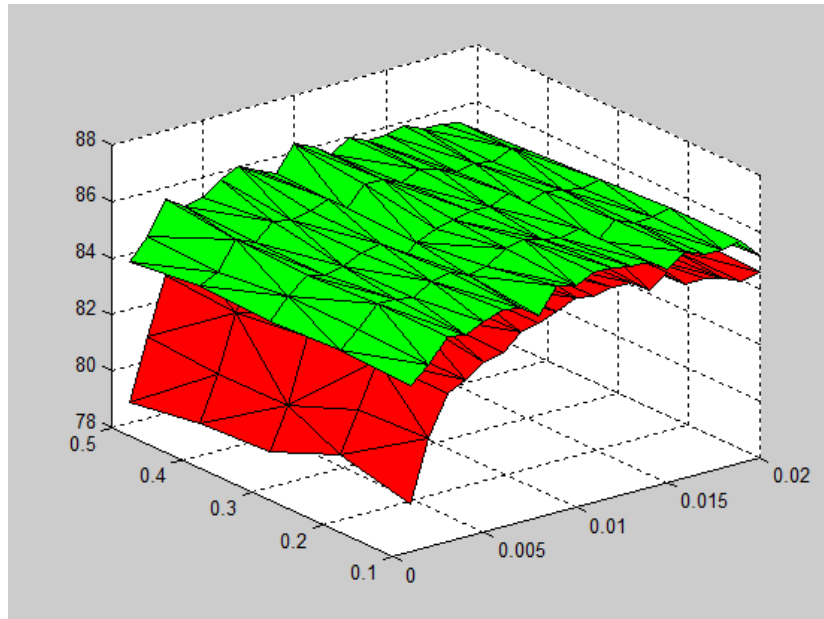


Figure 4.16: Surface graph showing variation in two methods for 3 fold cross validation on a dataset of 50 samples for upper zone

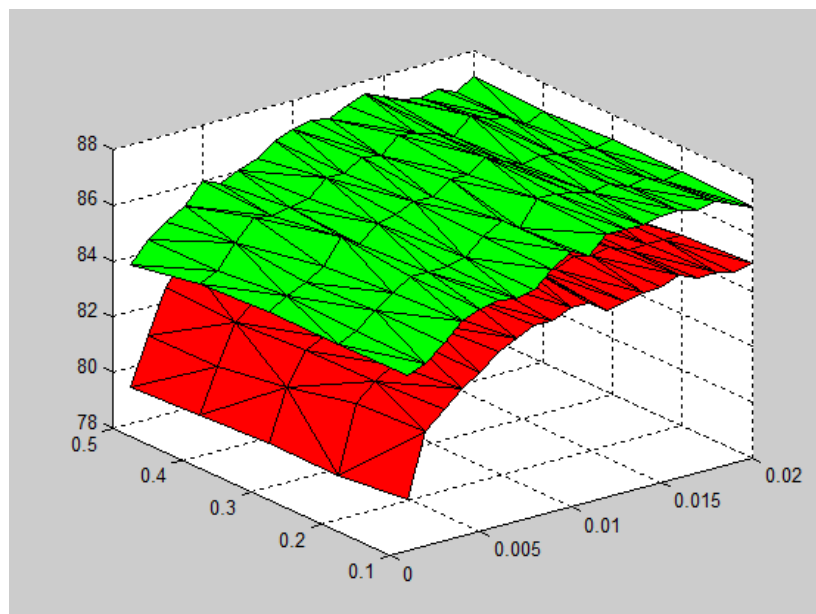


Figure 4.17: Surface graph showing variation in two methods for 4 fold cross validation on a dataset of 50 samples for upper zone

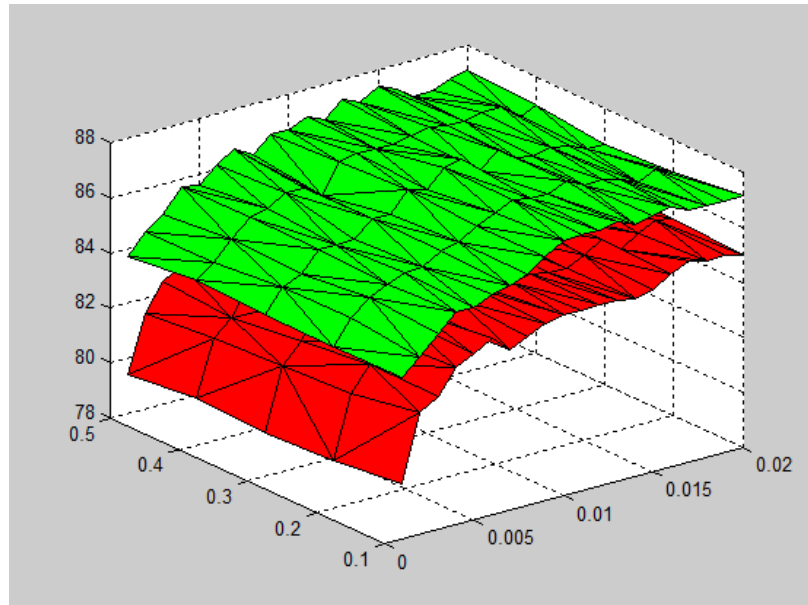


Figure 4.18: Surface graph showing variation in two methods for 5 fold cross validation on a dataset of 50 samples for upper zone

4.3 SCHEME 3: 70 samples in lower, middle and upper zones each has been taken for the cross validation. And the results are shown below.

For the lower zone 70 samples of each stroke have been taken. Maximum accuracy is given in the Table 5.1 for 3, 4, 5 fold cross validation.

Table 4.7: Cross validation accuracy of lower zone for 3, 4, and 5 cross validation for 70 samples

<i>k</i>-Fold Cross Validation (SVM)	Maximum accuracy of Algorithm proposed in Agrawal, (2012)	Maximum Accuracy of Algorithm proposed in this work
3	94.6%	96.7%
4	95.5%	96.5%
5	96.5%	95.3%

Three Dimensional graphs are shown in Figure 4.19, Figure 4.20 and Figure 4.21 for all the values of SVM parameters (gamma ranging from 0.001-0.02 and epsilon ranging from 0.1-0.5) have been drawn. These graphs show the variation between results obtained in Agrawal, (2012) and results obtained in this thesis work. It can be concluded from the graphs that the accuracy is better in this work as compared to

results in Agrawal, (2012) for 3, 4 and 5 fold cross validation on a data set of 70 samples for lower zone.

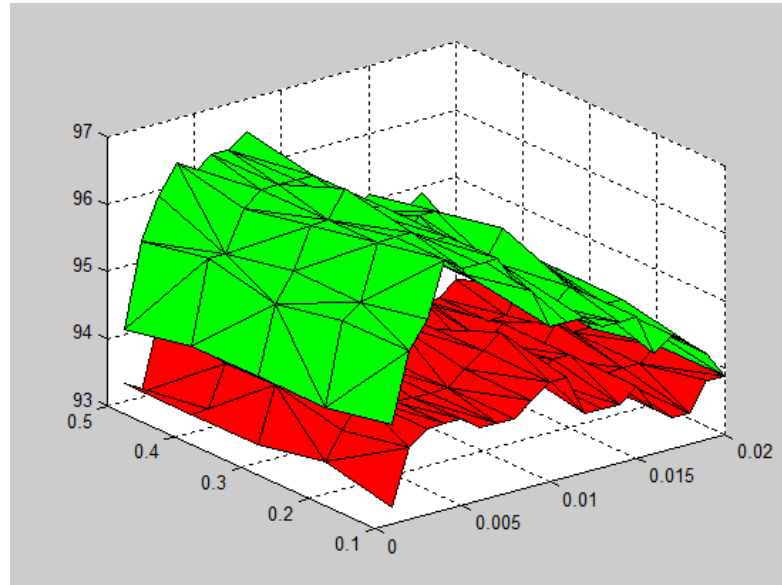


Figure 4.19: Surface graph showing variation in two methods for 3 fold cross validation on a dataset of 70 samples for lower zone

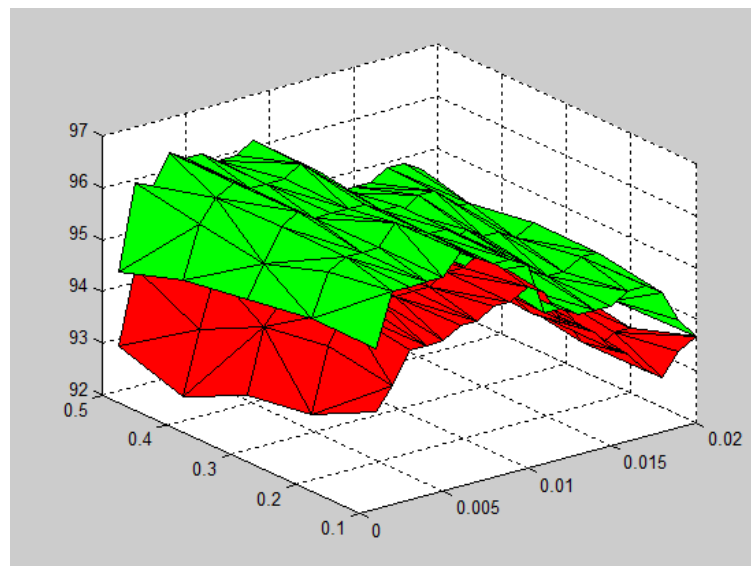


Figure 4.20: Surface graph showing variation in two methods for 4 fold cross validation on a dataset of 70 samples for lower zone

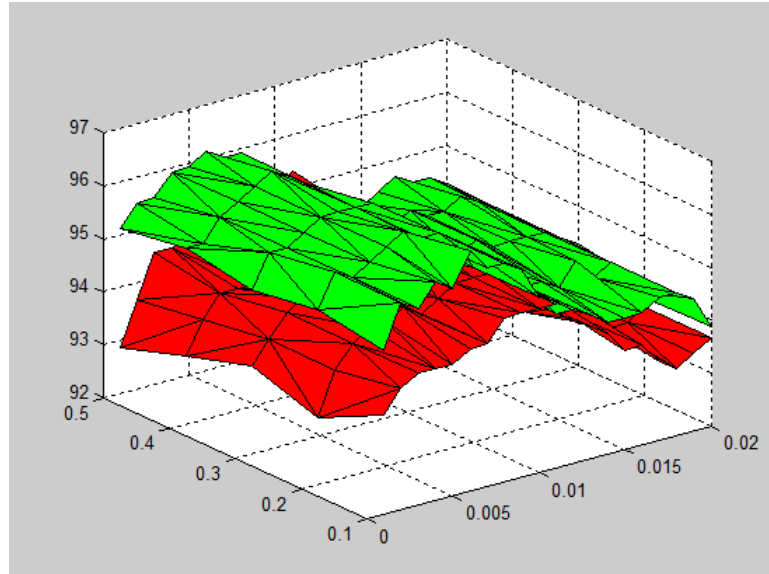


Figure 4.21: Surface graph showing variation in two methods for 5 fold cross validation on a dataset of 70 samples for lower zone

For the middle zone 70 samples of each stroke have been taken. Maximum accuracy is given in the Table 5.1 for 3, 4, 5 fold cross validation.

Table 4.8: Cross validation accuracy of middle zone for 3, 4, and 5 cross validation for 70 samples

<i>k</i>-Fold Cross Validation (SVM)	Maximum accuracy of Algorithm proposed in Agrawal, (2012)	Maximum Accuracy of Algorithm proposed in this work
3	89.3%	92.8%
4	89.9%	93.4%
5	90.3%	93.6%

Three Dimensional graphs are shown in Figure 4.22, Figure 4.23 and Figure 4.24 for all the values of SVM parameters (gamma ranging from 0.001-0.02 and epsilon ranging from 0.1-0.5) have been drawn. These graphs show the variation between results obtained in Agrawal, (2012) and results obtained in this thesis work. It can be concluded from the graphs that the accuracy is better in this work as compared to results in Agrawal, (2012) for 3, 4 and 5 fold cross validation on a data set of 70 samples for middle zone.

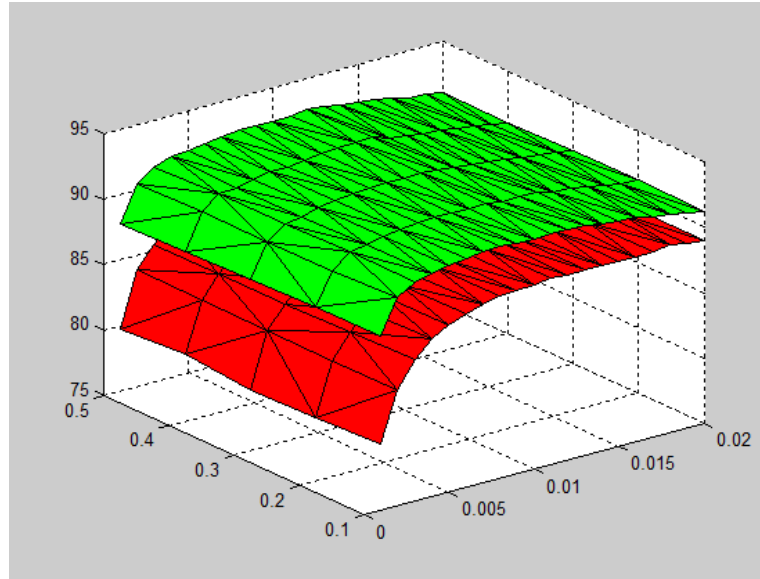


Figure 4.22: Surface graph showing variation in two methods for 3 fold cross validation on a dataset of 70 samples for middle zone

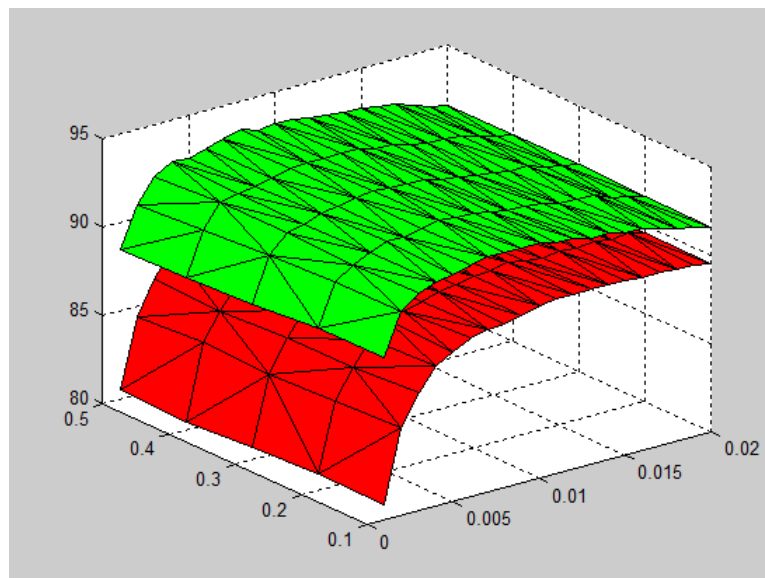


Figure 4.23: Surface graph showing variation in two methods for 4 fold cross validation on a dataset of 70 samples for middle zone

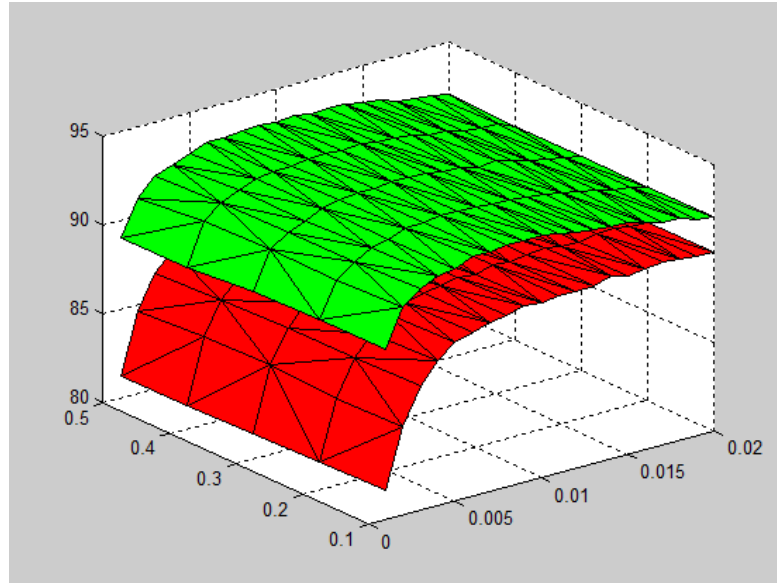


Figure 4.24: Surface graph showing variation in two methods for 5 fold cross validation on a dataset of 70 samples for middle zone

For the upper zone 70 samples of each stroke has been taken. Maximum accuracy is given in the Table 5.1 for 3, 4, 5 fold cross validation.

Table 4.9: Cross validation accuracy of upper zone for 3, 4, and 5 cross validation for 70 samples

<i>k</i>-Fold Cross Validation (SVM)	Maximum accuracy of Algorithm proposed in Agrawal, (2012)	Maximum Accuracy of Algorithm proposed in this work
3	85.6%	88.0%
4	86.1%	88.5%
5	86.0%	88.7%

Three Dimensional graphs are shown in Figure 4.25, Figure 4.26 and Figure 4.27 for all the values of SVM parameters (gamma ranging from 0.001-0.02 and epsilon ranging from 0.1-0.5) have been drawn. These graphs show the variation between results obtained in Agrawal, (2012) and results obtained in this thesis work. It can be concluded from the graphs that the accuracy is better in this work as compared to results in Agrawal, (2012) for 3, 4 and 5 fold cross validation on a data set of 70 samples for upper zone.

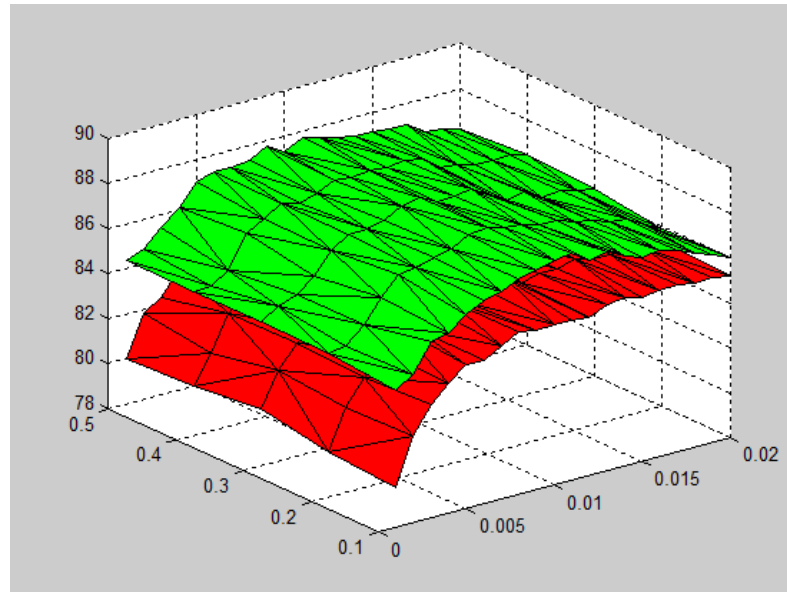


Figure 4.25: Surface graph showing variation in two methods for 3 fold cross validation on a dataset of 70 samples for upper zone

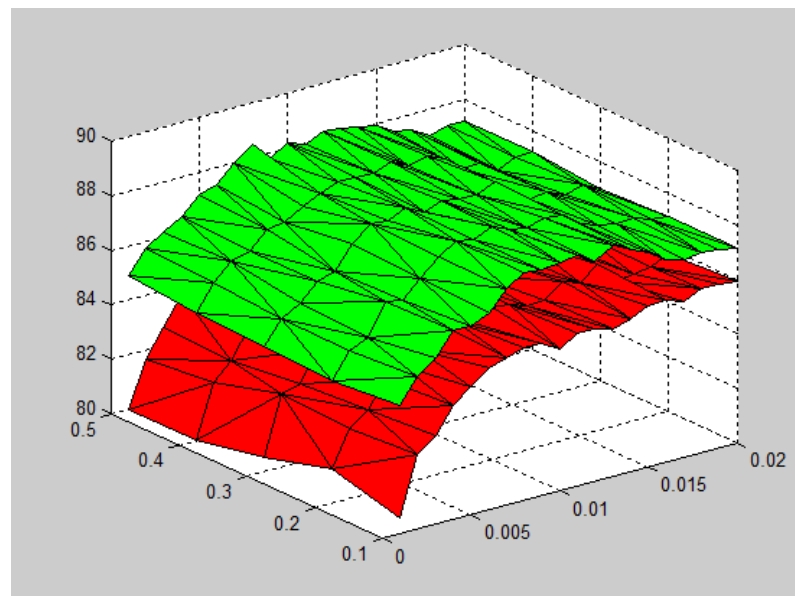


Figure 4.26: Surface graph showing variation in two methods for 4 fold cross validation on a dataset of 70 samples for upper zone

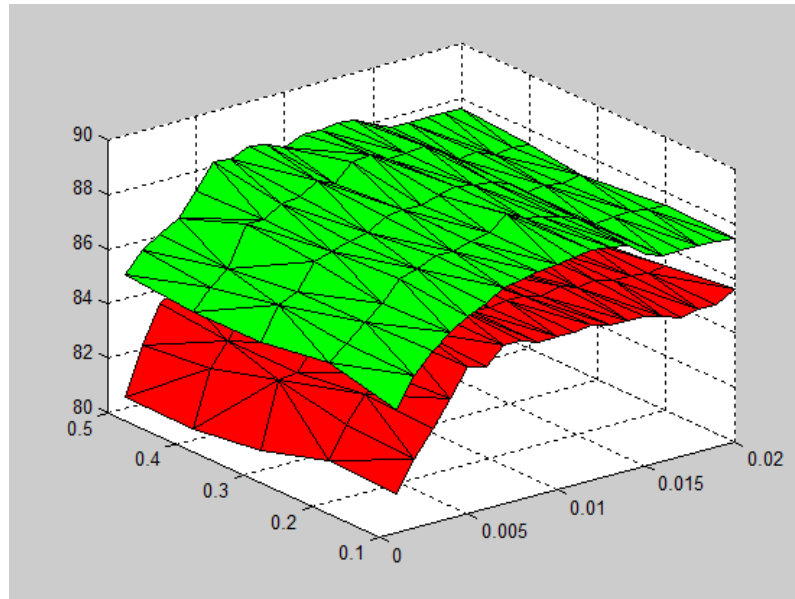


Figure 4.27: Surface graph showing variation in two methods for 5 fold cross validation on a dataset of 70 samples for upper zone

CONCLUSION AND FUTURE SCOPE

CONCLUSION

In this thesis, recognition accuracy of strokes has been improved by reconsidering the algorithms in pre-processing phase. The algorithm implemented in Agrawal, (2012) achieved a maximum accuracy of 90.4% but algorithms implemented in this thesis work have achieved an accuracy rate of 94.4% when a 75-25 ratio of training and testing data was considered. Algorithms implemented in this thesis are:

- Removal of duplicate points
- Normalization and centering
- Missing point interpolation
- Resampling of points

The main importance of the pre-processing phase is that after resampling, the 64 points should be at equal distance from each other and should retain its shape as it was written by the writer, which makes the difference of accuracy therefore in this thesis work we have got the higher accuracy rate. Thus, it can be concluded that pre-processing phase plays an important part in obtaining greater accuracy of the recognition system.

Future Scope

To increase the accuracy of recognition system we can add some new process in pre-processing phase as

- Smoothing of strokes can be done while the pen is moving from writing the stroke or after completion of stroke.
- Slant correction can also be added for increasing the accuracy as some writer has the habit of writing in cursive and in slant way.
- The complexity of the algorithms can be reduced to decrease the time of recognition.

- Strokes with less than 10 points captured can also be handled through extrapolation and through post-processing.
- We can add more features like curvature and direction for getting better accuracy.

REFERENCES

- 1) Aparna, K.H., Subramanian, V., Kasirajan, M., Prakash, G. V., Chakravarthy, V.S., and Madhvanath, S., 2004. Online handwriting recognition for Tamil. Ninth International Workshop on Frontiers in Handwriting Recognition, pp. 438-443.
- 2) Araki, N., Okuzaki M., Ishigaki, K. H., 2008. A Statistical Approach for Handwritten Character Recognition Using Bayesian Filter, 3rd International Conference on Innovative Computing Information and Control (ICICIC), pp. 194-198.
- 3) Arora, S., Bhattacharjee, D., Nasipuri, M., Basu, D. K. and Kundu, M., 2008. Combinig Multiple Feature Extraction Techniques for Handwriting Devnagari Character Recognition, Industrial and Information Systems, IEEE Region 10 Colloqium and the Third ICIIS, pp. 1-6.
- 4) Dungre, V. J. *et al.*, 2010. A review of Research on Devnagari Character recognition, International Journal of Computer Applications, vol. 12, No, 2, pp. 8-15.
- 5) Guerfali, W. and Plamondon, R., 1993. Normalizing and restoring online handwriting”, Pattern Recognition, Vol. 26, No. 3, pp. 419, 1993.
- 6) Homayoon, S.M., Beigi, K. N., Gregory J. Clary, and Subrahmonia. J., 1994. Size normalization in on-line unconstrained handwriting, IEEE Journal 0-8186-6950-0194.
- 7) Hosny, I., Abdou, S. and Fahmy, A., 2011. Using Advanced Hidden Markov Models for Online Arabic Handwriting Recognition, First Asian Conference on Pattern Recognition, pp.565-569.
- 8) Kumar, A. and Bhattacharya, S., 2010. Online Devnagari Isolated Character Recognition for the iPhone using Hidden Markov Model. Proceedings of the IEEE

Students' Technology Symposium, pp. 300-304.

- 9) Lehal, G. S., Singh, C., 2000. A Gurmukhi Script Recognition System, Proceeding of International Conference on Pattern Recognition, Vol. 2, pp. 557-560.
- 10) Lehal, G. S., Singh, C., 2002. A post Processor for Gurmukhi OCR, *Sadhana*, Vol. 27, Part 1, pp. 99-111.
- 11) Lehal, G. S., Singh, C., 2006. A Complete Machine printed Gurmukhi OCR, *Vivek*.
- 12) Mukherji, P., Rege, P., 2009. Shape Feature and Fuzzy Logic Based Offline Devnagari Handwritten Optical Character Recognition, *Journal of Pattern Recognition Research* 4, pp. 52-68.
- 13) Nair, A. and Leedham, C.G., 1991. Pre-processing of line codes for online recognition purposes, *Electronics Letters*, vol.2, no. 1, pp. 1-2.
- 14) Pal, U., Wakabayashi, Kimura. 2009. Comparative Study of Devnagari Handwritten Character Recognition using Different Feature and Classifier, 10th International Conference on Document Analysis and Recognition, pp. 1111-1115.
- 15) Pal, U. and Chaudhary, B.B., 2004. Indian Script Character Recognition, A survey, *Pattern Recognition*, Elsevier, pp. 1887-1899.
- 16) Pal, S., Mitra, J., Ghose, S., Banerjee P., 2007. A projection Based Statistical Approach for Handwritten Character Recognition, *Proceedings of International Conference on Computational Intelligence and Multimedia Applications*, vol. 2, pp. 404-408.
- 17) Sharma, A., Kumar, R. and Sharma, R. K., 2008. Online Handwritten Gurmukhi Character Recognition Using Elastic Matching, *Congress on Signal and Image Processing*, vol. 2, pp. 391-396.

- 18) Sharma, A., Kumar, R. and Sharma, R. K., 2009. Digit Extraction and Recognition from Machine Printed Gurumukhi Documents, Proceedings of the International Workshop on Multilingual OCR, MORC, Spain.
- 19) Kumar, R., Sharma, R. K., 2013. An Efficient Post-processing Algorithm for Online Handwritten Gurmukhi Character Recognition using Set Theory, International Journal for Pattern Recognition and Artificial Intelligence. Vol. 27, No. 4. Pp. 1353002 (17).
- 20) Singh, S., Aggarwal, A., and Dhir, R., 2012. Use of Gabor Filters for Recognition of Handwritten Gurmukhi Character, International Journal of Advance Research in Computer Science and Software Engineering. Vol. 2, Issue 5.
- 21) Tappert C.C., 1984. Adaptive on-line handwriting writing. Presented at The 7th International Conference on Pattern Recognition.
- 22) Tappert, C.C., Suen, C.Y. and Wakahara, T., 1988. Online handwriting recognition-A survey, Ninth International Conference on Pattern Recognition, vol.2, pp. 1123-1132.
- 23) Tappert, C.C. and Suen, C.Y., Wakahara, T., 1990. The state of the art in online handwriting recognition, IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 12, no. 8, pp. 787-808.
- 24) Vamvakas, G., Gatos, B., Petridis, S. and Stamatopoulos, N., 2007. An Efficient Feature Extraction and Dimensionality Reduction Scheme for Isolated Greek Handwritten Character Recognition, Ninth International Conference on Document Analysis and Recognition (ICDAR), Vol. 2, pp. 1073-1077.
- 25) Zheng, J., and Zhu, G., 2006. On-Line handwriting signature recognition based on wavelet energy feature matching, The Sixth World Congress on Intelligent Control and Automation, vol. 2, pp. 9885-9888.
- 26) Agrawal, R., 2012. Recognition of Online Handwritten Gurmukhi Stroke using

Support Vector Machine, M.tech. thesis, School of Mathematics and Computer Applications, Thapar University, Patiala.

- 27) Sharma, A., 2009. Online handwritten Gurmukhi Character Recognition, Ph.D. thesis, School of Mathematics and Computer Applications, Thapar University, Patiala.
- 28) Support Vector Machine (SVM), StatSoft website [Online]. Available: <http://www.statsoft.com/textbook/support-vector-machines/>
- 29) Support Vector Machine (SVM), Wikipedia website [Online]. Available: http://en.wikipedia.org/wiki/Support_vector_machine.
- 30) Unicode Inc., (1991-2013). Unicode 6.1 Character Code Charts [Online]. Available: <http://www.unicode.org/charts/PDF/U0A00.pdf>.
- 30) LIBSVM-Tools, A Library for Support Vector Machine [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>