

# Hand Gesture Recognition Using Neural Networks

by

**Klimis Symeonidis**

Submitted to the

School of Electronic and Electrical Engineering

On August 23, 2000

in partial fulfillment of the requirements for the degree of

Master of Science in Multimedia Signal Processing communications

## **Abstract**

Using orientation histograms a simple and fast algorithm will be developed to work on a workstation. It will recognize static hand gestures, namely, a subset of American Sign Language (ASL). Previous systems have used datagloves or markers for input in the system.

A pattern recognition system will be using a transform that converts an image into a feature vector, which will then be compared with the feature vectors of a training set of gestures. The final system will be implemented with a Perceptron network.

Thesis Supervisor: Terry Windeatt  
Centre for Vision, Speech and Signal Processing

Abstract	1
1 Introduction	5
2 Applications	7
<b>2.1 American Sign Language</b>	<b>7</b>
3 Background	9
<b>3.1 Object Recognition</b>	<b>10</b>
3.11 Large Object Tracking	10
3.12 Shape recognition	11
4 Goals	12
5 Neural Networks	13
<b>5.1 Neuron Model</b>	<b>14</b>
5.12 Simple Neuron	14
5.13 Transfer Functions	15
5.14 Supervised Learning	16
5.15 Unsupervised Learning	16
<b>5.2 Advantages of Neural Computing</b>	<b>17</b>
<b>5.3 Limitations of Neural Computing</b>	<b>18</b>
<b>5.4 Single Perceptron</b>	<b>19</b>
5.41 The Perceptron Convergence Algorithm	20
5.42 Linearly Separable Only	23
6 What Is MATLAB?	25
7 Approach	26
<b>7.1 Image database</b>	<b>26</b>
<b>7.2 Orientation Histograms</b>	<b>30</b>
7.21 Operation	31
<b>7.2 Perceptron implementation in MATLAB</b>	<b>35</b>
7.32 Linear Classification	36
7.33 Pre-processing	39
7.34 Input , Target , Test and Output files	40
7.35 Limitations and Cautions	42
8 RESULTS	44
Conclusion	56
References:	58
APPENDIX	61
<b>Part1 : Image Processing Program</b>	<b>61</b>
<b>Part2 : Neural Network Program</b>	<b>65</b>

## List of Figures

Figure 1: ASL examples	8
Figure 2: Neural Net block diagram	13
Figure 3: Neuron	14
Figure 4 : Transfer Functions	15
Figure 5 : Perceptron Schematic Diagram	19
Figure 6: Perceptron Signal Flow Graph	21
Figure 7 : Linear Separable	23
Figure 8: Train – Test images	27
Figure 9 : Pattern Recognition System	29
Figure 10 : Illumination Variance	29
Figure 11 : Orientation histogram	31
Figure 12: X-Y filters	32
Figure 13: Orientation histogram examples	34
Figure 14	34
Figure 15: Perceptron flow-chart	38
Figure 16: Train Error	41
Figure 17: Zero	47
Figure 18: One	48
Figure 19: L	49
Figure 20	49
Figure 21: H	50
Figure 22: A	51
Figure 23: V	52
Figure 24: W	53
Figure 25: Signs G, H, V, B	54

## List of Tables

<i>Table 1: Train Set sample</i>	28
<i>Table 2: 'x1 and not x2</i>	23
<i>Table 3: Target Vectors</i>	45
<i>Table 4: Zero Test Results</i>	47
<i>Table 5: One Test Results</i>	48
<i>Table 6: L Test Results</i>	49
<i>Table 7: H Test Results</i>	50
<i>Table 8: L Test Results</i>	51
<i>Table 9: L Test Results</i>	52
<i>Table 10: L Test Results</i>	53
<i>Table 11: G Test Results</i>	54
<i>Table 12: B Test Results</i>	54

# 1 Introduction

Since the introduction of the most common input computer devices not a lot have changed. This is probably because the existing devices are adequate. It is also now that computers have been so tightly integrated with everyday life, that new applications and hardware are constantly introduced. The means of communicating with computers at the moment are limited to keyboards, mice, light pen, trackball, keypads etc.

These devices have grown to be familiar but inherently limit the speed and naturalness with which we interact with the computer.

As the computer industry follows Moore's Law since middle 1960s, powerful machines are built equipped with more peripherals. Vision based interfaces are feasible and at the present moment the computer is able to "see". Hence users are allowed for richer and user-friendlier man-machine interaction. This can lead to new interfaces that will allow the deployment of new commands that are not possible with the current input devices. Plenty of time will be saved as well.

Recently, there has been a surge in interest in recognizing human hand gestures. Hand-gesture recognition has various applications like computer games, machinery control (e.g. crane), and thorough mouse replacement. One of the most structured sets of gestures belongs to sign language. In sign language, each gesture has an assigned meaning (or meanings).

Computer recognition of hand gestures may provide a more natural-computer interface, allowing people to point, or rotate a CAD model by rotating their hands. Hand gestures can be classified in two categories: static and dynamic. A static gesture is a particular hand configuration and pose, represented by a single image. A dynamic gesture is a moving gesture, represented by a sequence of images. We will focus on the recognition of static images.

Interactive applications pose particular challenges. The response time should be very fast. The user should sense no appreciable delay between when he or she makes a gesture or

motion and when the computer responds. The computer vision algorithms should be reliable and work for different people.

There are also economic constraints: the vision-based interfaces will be replacing existing ones, which are often very low cost. A hand-held video game controller and a television remote control each cost about \$40. Even for added functionality, consumers may not want to spend more. When additional hardware is needed the cost is considerable higher.

Academic and industrial researchers have recently been focusing on analyzing images of people. While researchers are making progress, the problem is hard and many present day algorithms are complex, slow or unreliable. The algorithms that do run near real-time do so on computers that are very expensive relative to the existing hand-held interface devices.

## 2 Applications

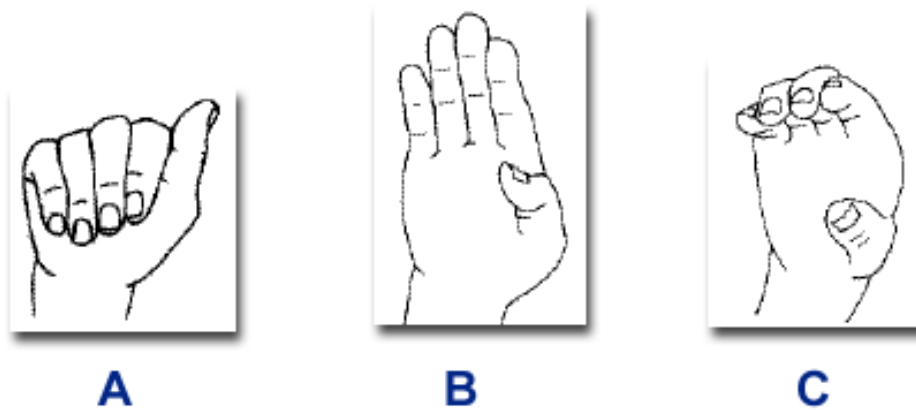
Creating a proper sign language (ASL – American Sign Language at this case) dictionary is not the desired result at this point. This would combine advanced grammar and syntax structure understanding of the system, which is outside the scope of this project. The American Sign Language will be used as the database since it's a tightly structured set. From that point further applications can be suited. The distant (or near?) future of computer interfaces could have the usual input devices and in conjunction with gesture recognition some of the user's feelings would be perceived as well.

Taking ASL recognition further a full real-time dictionary could be created with the use of video. As mentioned before this would require some Artificial Intelligence for grammar and syntax purposes.

Another application is huge database annotation. It is far more efficient when properly executed by a computer, than by a human.

### 2.1 American Sign Language

American Sign Language is the language of choice for most deaf people in the United States. It is part of the “deaf culture” and includes its own system of puns, inside jokes, etc. However, ASL is one of the many sign languages of the world. As an English speaker would have trouble understanding someone speaking Japanese, a speaker of ASL would have trouble understanding the Sign Language of Sweden. ASL also has its own grammar that is different from English. ASL consists of approximately 6000 gestures of common words with finger spelling used to communicate obscure words or proper nouns. Finger spelling uses one hand and 26 gestures to communicate the 26 letters of the alphabet. Some of the signs can be seen in **Fig(1)** below.



**Figure 1:** ASL examples

Another interesting characteristic that will be ignored by this project is the ability that ASL offers to describe a person, place or thing and then point to a place in space to temporarily store for later reference.

ASL uses facial expressions to distinguish between statements, questions and directives. The eyebrows are raised for a question, held normal for a statement, and furrowed for a directive. There has been considerable work and research in facial feature recognition, they will not be used to aid recognition in the task addressed. This would be feasible in a full real-time ASL dictionary.



### 3 Background

Research on hand gestures can be classified into three categories. The first category, glove based analysis, employs sensors (mechanical or optical) attached to a glove that transduces finger flexions into electrical signals for determining the hand posture. The relative position of the hand is determined by an additional sensor. This sensor is normally a magnetic or an acoustic sensor attached to the glove. For some dataglove applications, look-up table software toolkits are provided with the glove to be used for hand posture recognition.

The second category, vision based analysis, is based on the way human beings perceive information about their surroundings, yet it is probably the most difficult to implement in a satisfactory way. Several different approaches have been tested so far. One is to build a three-dimensional model of the human hand. The model is matched to images of the hand by one or more cameras, and parameters corresponding to palm orientation and joint angles are estimated. These parameters are then used to perform gesture classification. A hand gesture analysis system based on a three-dimensional hand skeleton model with 27 degrees of freedom was developed by *Lee and Kunii*. They incorporated five major constraints based on the human hand kinematics to reduce the model parameter space search. To simplify the model matching, specially marked gloves were used.

The third category, *analysis of drawing gestures*, usually involves the use of a stylus as an input device. Analysis of drawing gestures can also lead to recognition of written text.

The vast majority of hand gesture recognition work has used mechanical sensing, most often for direct manipulation of a virtual environment and occasionally for symbolic communication. Sensing the hand posture mechanically has a range of problems, however, including reliability, accuracy and electromagnetic noise. Visual sensing has the potential to make gestural interaction more practical, but potentially embodies some of the most difficult problems in machine vision. The hand is a non-rigid object and even worse self-occlusion is very usual.

Full ASL recognition systems (words, phrases) incorporate datagloves. *Takashi and Kishino* discuss a Dataglove-based system that could recognize 34 of the 46 Japanese gestures (user dependent) using a joint angle and hand orientation coding technique. From their paper, it seems the test user made each of the 46 gestures 10 times to provide data for principle component and cluster analysis. A separate test was created from five iterations of the alphabet by the user, with each gesture well separated in time. While these systems are technically interesting, they suffer from a lack of training.

Excellent work has been done in support of machine sign language recognition by *Sperling and Parish*, who have done careful studies on the bandwidth necessary for a sign conversation using spatially and temporally sub-sampled images. Point light experiments (where “lights” are attached to significant locations on the body and just these points are used for recognition), have been carried out by *Poizner*.

Most systems to date study isolate/static gestures. In most of the cases those are finger-spelling signs.

### **3.1 Object Recognition**

#### **3.11 Large Object Tracking**

In some interactive applications, the computer needs to track the position or orientation of a hand that is prominent in the image. Relevant applications might be computer games, or interactive machine control. In such cases, a description of the overall properties of the image may be adequate. Image moments, which are fast to compute, provide a very coarse summary of global averages of orientation and position. If the hand is on a uniform background, this method can distinguish hand positions and simple pointing gestures.

The large-object-tracking method makes use of a low-cost detector/processor to quickly calculate moments. This is called the artificial retina chip. This chip combines image detection with some low-level image processing (named artificial retina by analogy with

those combined abilities of the human retina). The chip can compute various functions useful in the fast algorithms for interactive graphics applications.

### **3.12 Shape recognition**

Most applications, such as recognizing particular static hand signal, require a richer description of the shape of the input object than image moments provide.

If the hand signals fell in a predetermined set, and the camera views a close-up of the hand, we may use an example-based approach, combined with a simple method to analyze hand signals called orientation histograms.

These example-based applications involve two phases; training and running. In the training phase, the user shows the system one or more examples of a specific hand shape. The computer forms and stores the corresponding orientation histograms. In the run phase, the computer compares the orientation histogram of the current image with each of the stored templates and selects the category of the closest match, or interpolates between templates, as appropriate. This method should be robust against small differences in the size of the hand but probably would be sensitive to changes in hand orientation.

## 4 Goals

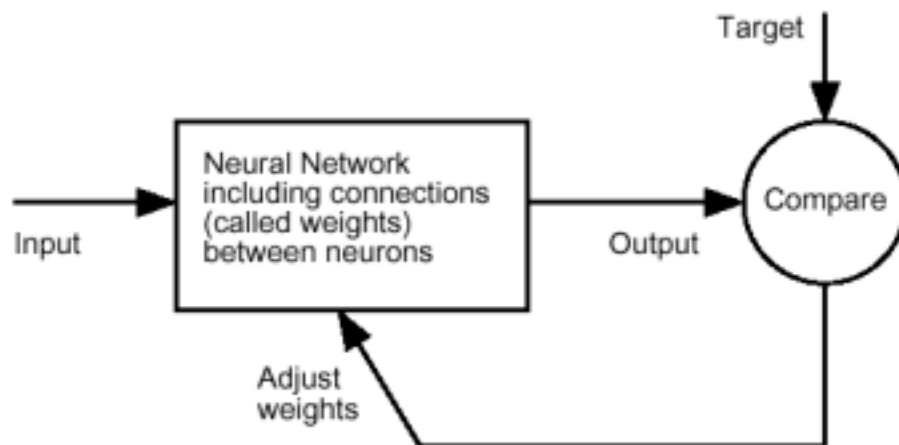
The scope of this project is to create a method to recognize hand gestures, based on a pattern recognition technique developed by *McConnell*; employing histograms of local orientation. The orientation histogram will be used as a feature vector for gesture classification and interpolation.

High priority for the system is to be simple without making use of any special hardware. All the computation should occur on a workstation or PC. Special hardware would be used only to digitize the image (scanner or digital camera

## 5 Neural Networks

Neural networks are composed of simple elements operating in parallel. These elements are inspired by biological nervous systems. As in nature, the network function is determined largely by the connections between elements. We can train a neural network to perform a particular function by adjusting the values of the connections (weights) between elements.

Commonly neural networks are adjusted, or trained, so that a particular input leads to a specific target output. Such a situation is shown in fig(2). There, the network is adjusted, based on a comparison of the output and the target, until the network output matches the target. Typically many such input/target pairs are used, in this *supervised learning* (training method studied in more detail on following chapter), to train a network.



**Figure 2:** Neural Net block diagram

Neural networks have been trained to perform complex functions in various fields of application including pattern recognition, identification, classification, speech, vision and control systems.

Today neural networks can be trained to solve problems that are difficult for conventional computers or human beings. The supervised training methods are commonly used, but other networks can be obtained from *unsupervised training* techniques or from direct *design* methods. Unsupervised networks can be used, for instance, to identify groups of

data. Certain kinds of linear networks and Hopfield networks are designed directly. In summary, there are a variety of kinds of design and learning techniques that enrich the choices that a user can make.

The field of neural networks has a history of some five decades but has found solid application only in the past fifteen years, and the field is still developing rapidly. Thus, it is distinctly different from the fields of control systems or optimization where the terminology, basic mathematics, and design procedures have been firmly established and applied for many years.

## 5.1 Neuron Model

### 5.12 Simple Neuron

A neuron with a single scalar input and no bias is shown on the left below.

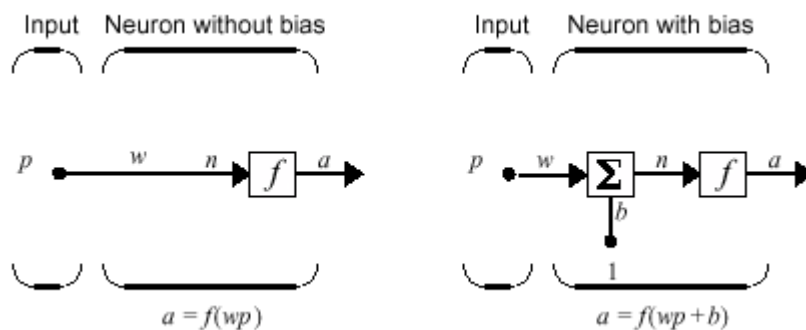


Figure 3: Neuron

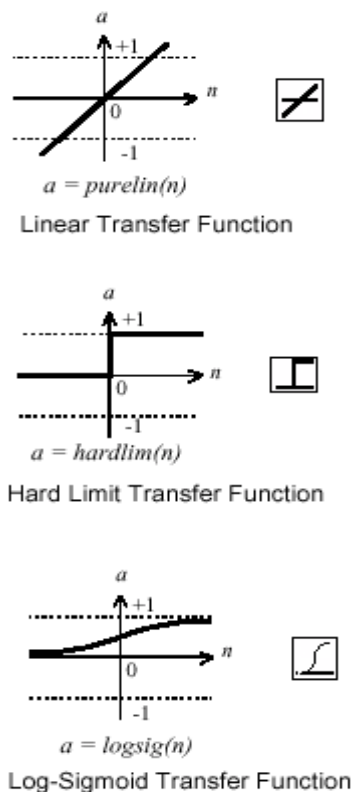
The scalar input  $p$  is transmitted through a connection that multiplies its strength by the scalar weight  $w$ , to form the product  $wp$ , again a scalar. Here the weighted input  $wp$  is the only argument of the transfer function  $f$ , which produces the scalar output  $a$ . The neuron on the right has a scalar bias,  $b$ . You may view the bias as simply being added to the product  $wp$  as shown by the summing junction or as shifting the function  $f$  to the left by an amount  $b$ . The bias is much like a weight, except that it has a constant input of 1. The transfer function net input  $n$ , again a scalar, is the sum of the weighted input  $wp$  and the bias  $b$ . This sum is the argument of the transfer function  $f$ . Here  $f$  is a transfer function,

typically a step function or a sigmoid function, that takes the argument  $n$  and produces the output  $a$ . Examples of various transfer functions are given in the next section. Note that  $w$  and  $b$  are both *adjustable* scalar parameters of the neuron. The central idea of neural networks is that such parameters can be adjusted so that the network exhibits some desired or interesting behavior.

Thus, we can train the network to do a particular job by adjusting the weight or bias parameters, or perhaps the network itself will adjust these parameters to achieve some desired end. All of the neurons in the program written in MATLAB have a bias. However, you may omit a bias in a neuron if you wish.

### 5.13 Transfer Functions

Three of the most commonly used transfer functions are shown in fig(5).



**Figure 4 :** Transfer Functions

The hard limit transfer function shown above limits the output of the neuron to either 0, if the net input argument  $n$  is less than 0, or 1, if  $n$  is greater than or equal to 0. This is the

function used for the Perceptron algorithm written in MATLAB to create neurons that make a classification decision.

There are two modes of learning: Supervised and unsupervised. Below there is a brief description of each one to determine the best one for our problem

### **5.14 Supervised Learning**

Supervised learning is based on the system trying to predict outcomes for known examples and is a commonly used training method. It compares its predictions to the target answer and "learns" from its mistakes. The data start as inputs to the input layer neurons. The neurons pass the inputs along to the next nodes. As inputs are passed along, the weighting, or connection, is applied and when the inputs reach the next node, the weightings are summed and either intensified or weakened. This continues until the data reach the output layer where the model predicts an outcome. In a supervised learning system, the predicted output is compared to the actual output for that case. If the predicted output is equal to the actual output, no change is made to the weights in the system. But, if the predicted output is higher or lower than the actual outcome in the data, the error is propagated back through the system and the weights are adjusted accordingly. This feeding errors backwards through the network is called "back-propagation." Both the Multi-Layer Perceptron and the Radial Basis Function are supervised learning techniques. The Multi-Layer Perceptron uses the back-propagation while the Radial Basis Function is a feed-forward approach which trains on a single pass of the data.

### **5.15 Unsupervised Learning**

Neural networks which use unsupervised learning are most effective for describing data rather than predicting it. The neural network is not shown any outputs or answers as part of the training process--in fact, there is no concept of output fields in this type of system. The primary unsupervised technique is the Kohonen network. The main uses of Kohonen and other unsupervised neural systems are in cluster analysis where the goal is to group



"like" cases together. The advantage of the neural network for this type of analysis is that it requires no initial assumptions about what constitutes a group or how many groups there are. The system starts with a clean slate and is not biased about which factors should be most important.

## **5.2 Advantages of Neural Computing**

There are a variety of benefits that an analyst realizes from using neural networks in their work.

- Pattern recognition is a powerful technique for harnessing the information in the data and generalizing about it. Neural nets learn to recognize the patterns which exist in the data set.
- The system is developed through learning rather than programming. Programming is much more time consuming for the analyst and requires the analyst to specify the exact behavior of the model. Neural nets teach themselves the patterns in the data freeing the analyst for more interesting work.
- Neural networks are flexible in a changing environment. Rule based systems or programmed systems are limited to the situation for which they were designed--when conditions change, they are no longer valid. Although neural networks may take some time to learn a sudden drastic change, they are excellent at adapting to constantly changing information.
- Neural networks can build informative models where more conventional approaches fail. Because neural networks can handle very complex interactions they can easily model data which is too difficult to model with traditional approaches such as inferential statistics or programming logic.
- Performance of neural networks is at least as good as classical statistical modeling, and better on most problems. The neural networks build models that are more reflective of the structure of the data in significantly less time.

- ✿ Neural networks now operate well with modest computer hardware. Although neural networks are computationally intensive, the routines have been optimized to the point that they can now run in reasonable time on personal computers. They do not require supercomputers as they did in the early days of neural network research.

### **5.3 *Limitations of Neural Computing***

There are some limitations to neural computing. The key limitation is the neural network's inability to explain the model it has built in a useful way. Analysts often want to know why the model is behaving as it is. Neural networks get better answers but they have a hard time explaining how they got there.

There are a few other limitations that should be understood. First, It is difficult to extract rules from neural networks. This is sometimes important to people who have to explain their answer to others and to people who have been involved with artificial intelligence, particularly expert systems which are rule-based.

As with most analytical methods, you cannot just throw data at a neural net and get a good answer. You have to spend time understanding the problem or the outcome you are trying to predict. And, you must be sure that the data used to train the system are appropriate and are measured in a way that reflects the behavior of the factors. If the data are not representative of the problem, neural computing will not produce good results. This is a classic situation where "garbage in" will certainly produce "garbage out."

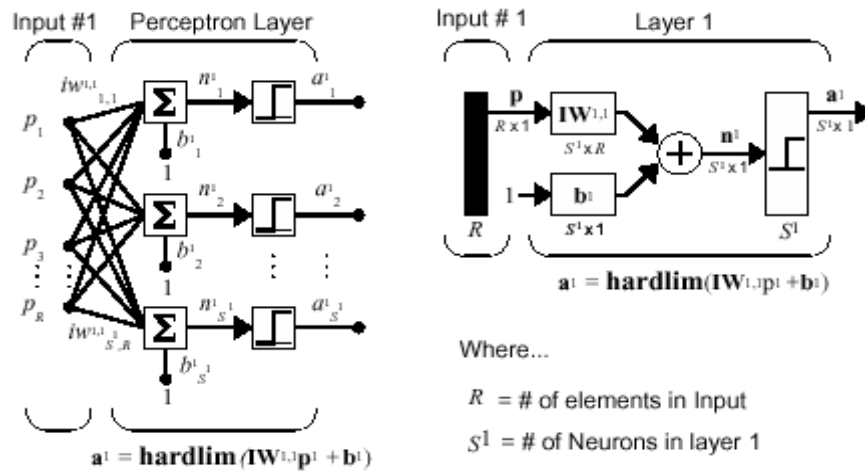
Finally, it can take time to train a model from a very complex data set. Neural techniques are computer intensive and will be slow on low end PCs or machines without math coprocessors. It is important to remember though that the overall time to results can still be faster than other data analysis approaches, even when the system takes longer to train. Processing speed alone is not the only factor in performance and neural networks do not require the time programming and debugging or testing assumptions that other analytical approaches do.

## 5.4 Single Perceptron

The perceptron is a program that learns **concepts**, i.e. it can learn to respond with *True* (1) or *False* (0) for inputs we present to it, by repeatedly "studying" examples presented to it.

The structure of a single perceptron is very simple. There are two inputs, a bias, and an output. A simple schematic diagram is shown in fig(5) .

Note that both the inputs and outputs of a perceptron are binary - that is they can only be 0 or 1.



**Figure 5 : Perceptron Schematic Diagram**

Each of the inputs and the bias is connected to the main perceptron by a weight. A weight is generally a real number between 0 and 1. When the input number is fed into the perceptron, it is multiplied by the corresponding weight. After this, the weights are all summed up and fed through a hard-limiter. Basically, a hard-limiter is a function that defines the threshold values for 'firing' the perceptron. For example, the limiter could be:

$$f(x) = \begin{cases} x \leq 0 \rightarrow 0 \\ x \geq 1 \rightarrow 1 \end{cases}$$

For example, if the sum of the input multiplied by the weights is -2, the limiting function would return 0. Or if the sum was 3, the function would return 1.

Now, the way a perceptron learns to distinguish patterns is through modifying its weights - the concept of a learning rule must be introduced. In the perceptron, the most common form of learning is by adjusting the weights by the *difference between the desired output and the actual output*. Mathematically, this can be written:

$$\Delta w_i = x_i \delta$$

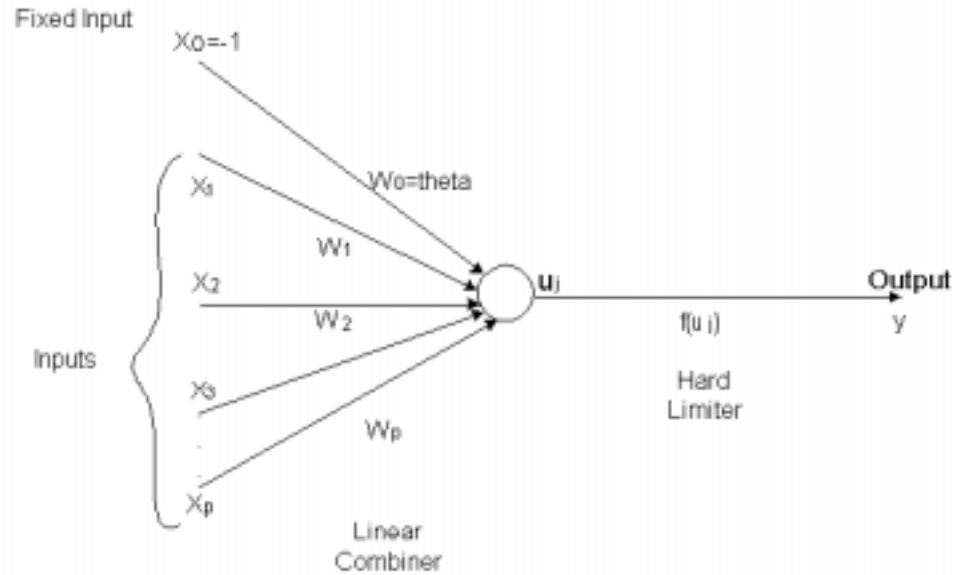
where  $\delta = (\text{desired output}) - (\text{actual output})$

Learning on a perceptron is guaranteed, as stated by the Perceptron Convergence Theorem which states that if a solution can be implemented on a perceptron, the learning rule will find the solution in a finite number of steps.

Proof of this theorem can be found in *Minsky and Papert's* 1989 book, *Perceptrons*.

### 5.41 The Perceptron Convergence Algorithm

For the development of the error-correction learning algorithm for a single-layer perceptron, we will work with the signal-flow graph shown in fig(6). In this case the threshold  $\theta(n)$  is treated as a synaptic weight connected to a fixed input equal to -1.



**Figure 6:** Perceptron Signal Flow Graph

We may then define the  $(p+1)$ -by-1 input vector:

$$x(n) = [-1, x_1(n), x_2(n), \dots, x_p(n)]^T$$

Correspondingly we define the  $(p+1)$ -by-1 weight vector:

$$w(n) = [\theta(n), w_1(n), w_2(n), \dots, w_p(n)]^T$$

Below are some variable and parameters used in the convergence algorithm

$\theta(n)$  = threshold

$y(n)$  = actual response

$d(n)$  = desired response

$\eta$  = learning rate parameter,  $0 < \eta < 1$

So lets see the 4-step algorithm in greater detail:

### Step 1: Initialization

Set  $w(0)=0$ . Then perform the following computations for time  $n=1,2,...$ .

### Step 2: Activation

At time  $n$ , activate the perceptron by applying continuous-valued input vector  $x(n)$  and desired response  $d(n)$ .

### Step 3: Computation of Actual Response

Compute the actual response of the perceptron:

$$y(n) = \text{sgn}[w^T(n)x(n)]$$

The linear output is written in the form:

$$u(n) = w^T(n)x(n)$$

where :

$$\begin{aligned} \text{sgn}(u) &= +1 && \text{if } u > 0 \\ \text{sgn}(u) &= -1 && \text{if } u < 0 \end{aligned}$$

### Step 4: Adaptation of Weight Vector

$$w(n+1) = w(n) + \eta[d(n) - y(n)]x(n)$$

where  $d(n) = +1$  if  $x(n)$  belongs to class  $C_1$

where  $d(n) = -1$  if  $x(n)$  belongs to class  $C_2$

### Step 5

Increment time  $n$  by one unit and go back to step 2

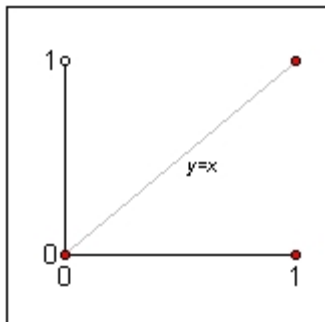
## 5.42 Linearly Separable Only

Perceptrons can only solve problems where the solutions can be divided by a line (or hyperplane) - this is called linear separation. To explain the concept of linear separation further, let us look at the function shown to the left. The function reads '*x1 and (not x2)*'. Let us assume that we run the function through a perceptron, and the weights converge at 0 for the bias, and 2, -2 for the inputs. If we calculate the *net* value (the weighted sum) we get:

$y = x_1 \wedge \neg x_2$		
$x_1$	$x_2$	$y$
0	0	0
0	1	0
1	0	1
1	1	0

**Table 1:** '*x1 and not x2*'

$$net = \sum_{i=0}^2 w_i x_i = (0)x_0 + (2)x_1 + (-2)x_2 = 2x_1 - 2x_2$$



**Figure 7 :** Linear Separable

So the perceptron correctly draws a line that divides the two groups of points. Note that it doesn't only have to be a line, it can be a hyperplane dividing points in 3-D space, or beyond. This is where the power of perceptrons lies, but also where its limitations lie. For example, perceptrons cannot solve the XOR binary functions.

Non-linear separable problems can be solved by a perceptron provided an appropriate set of first-layer processing elements exists. The first layer of fixed processing units computes a set of functions whose values depend on the input pattern. Even though the data set of input patterns may not be separable, when viewed in the space of original

input variables, it can easily be the case that the same set of patterns becomes linearly separable. We will later see how such a technique is used in the perceptron's development so that the chance that the preprocessing layer is maximized.



## 6 What Is MATLAB?



The name MATLAB stands for *matrix laboratory*.

MATLAB is a high-performance language for technical computing. It integrates computation, visualization, and programming in an easy-to-use environment where problems and solutions are expressed in familiar mathematical notation. Typical uses include:

- Math and computation
- Algorithm development
- Modeling, simulation, and prototyping
- Data analysis, exploration, and visualization
- Scientific and engineering graphics
- Application development, including Graphical User Interface building

MATLAB is an interactive system whose basic data element is an array that does not require dimensioning. This allows you to solve many technical computing problems, especially those with matrix and vector formulations, in a fraction of the time it would take to write a program in a scalar non-interactive language such as C or Fortran.

MATLAB has evolved over a period of years with input from many users. In university environments, it is the standard instructional tool for introductory and advanced courses in mathematics, engineering, and science. In industry, MATLAB is the tool of choice for high-productivity research, development, and analysis.

The reason that I have decided to use MATLAB for the development of this project is its toolboxes. Toolboxes allow you to *learn* and *apply* specialized technology. Toolboxes are comprehensive collections of MATLAB functions (M-files) that extend the MATLAB environment to solve particular classes of problems. It includes among others image processing and neural networks toolboxes.

## 7 Approach

### 7.1 Image database

The starting point of the project was the creation of a database with all the images that would be used for training and testing.

The image database can have different formats. Images can be either hand drawn, digitized photographs or a 3D dimensional hand. Photographs were used, as they are the most realistic approach.

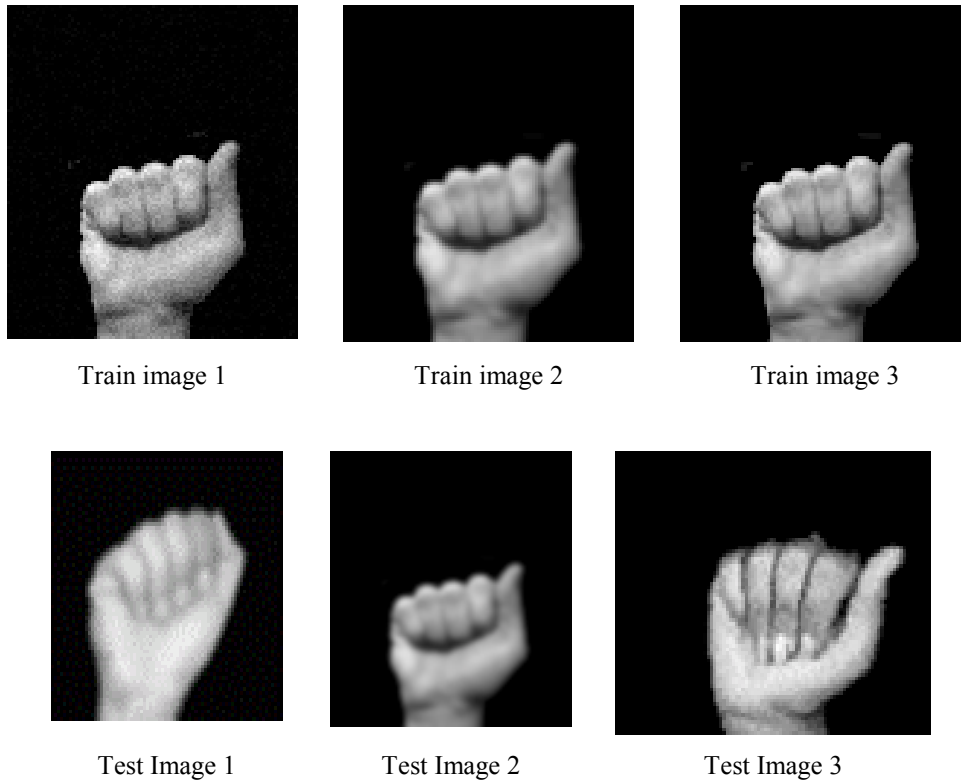
Images came from two main sources. Various ASL databases on the Internet and photographs I took with a digital camera. This meant that they have different sizes, different resolutions and some times almost completely different angles of shooting. Images belonging to the last case were very few but they were discarded, as there was no chance of classifying them correctly. Two operations were carried out in all of the images. They were converted to grayscale and the background was made uniform. The internet databases already had uniform backgrounds but the ones I took with the digital camera had to be processed in *Adobe Photoshop*.

Drawn images can still simulate translational variances with the help of an editing program (e.g. *Adobe Photoshop*).

The database itself was constantly changing throughout the completion of the project as it was it that would decide the robustness of the algorithm. Therefore, it had to be done in such way that different situations could be tested and thresholds above which the algorithm didn't classify correct would be decided.

The construction of such a database is clearly dependent on the application. If the application is a crane controller for example operated by the same person for long periods the algorithm doesn't have to be robust on different person's images. In this case noise and motion blur should be tolerable. The applications can be of many forms and since I wasn't developing for a specific one I have tried to experiment for many alternatives.

We can see an example below. In the first row are the training images. In the second the testing images.



**Figure 8:** Train – Test images

For most of the gestures the training set originates from a single gesture. Those were enhanced in *Adobe Photoshop* using various filters. The reason for this is that I wanted the algorithm to be very robust for images of the same database. If there was a misclassification to happen it would be preferred to be for unknown images.

The final form of the database is this.

**Train set:**

Eight training sets of images, each one containing three images. Each set originates from a single image for testing.

**Test Set:**

The number of test images varies for each gesture. There is no reason for keeping those on a constant number. Some images can tolerate much more variance and images from

new databases and they can be tested extensively, while other images are restricted to fewer testing images.

### Train Set sample

0	<i>Original</i>	<i>Blur radius (pixels)</i>	<i>Noise</i>
1	✓	×	×
2	×	0.5	×
3	×	×	10 Uniform

1	<i>Original</i>	<i>Blur radius (pixels)</i>	<i>Noise</i>
1	✓	×	×
2	×	0.4	×
3	×	×	10 Gaussian

H	<i>Original</i>	<i>Blur radius (pixels)</i>	<i>Noise</i>
1	✓	×	×
2	×	0.6	×
3	×	×	6 Gaussian

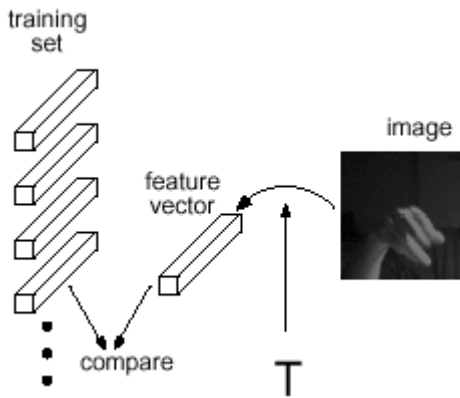
**Table 2:** Train Set sample

The databases are shown in more detail in the results section where all the training and testing images are listed.

The system could be approached either in high or low-level. The former would employ models of the hand, finger, joints and perhaps fit such a model to the visual data. This approach offers robustness, but at the expense of speed.

A low-level approach would process data at a level not much higher than that of pixel intensities.

Although this approach would not have the power to make inferences about occluded data, it could be simple and fast. The pattern recognition system that will be used can be seen in Fig (9). Some transformation  $T$ , converts an image into a feature vector, which will be then compared with feature vectors of a training set of gestures.



**Figure 9 :** Pattern Recognition System

We will be seeking for the simplest possible transformation  $T$ , which allows gesture recognition.

Histogram orientation has the advantage of being robust in lighting change conditions. If we follow the pixel-intensities approach certain problems can arise for varying illumination. Taking a pixel-by-pixel difference of the same photo under different lighting conditions would show a large distance between these two identical gestures. For the pixel-intensity approach no transformation  $T$  has been applied. The image itself is used as the feature vector. In Fig (10) we can see the same hand gesture under different lighting conditions.



**Figure 10 :** Illumination Variance

Another important aspect of gesture recognition is translation invariance. The position of the hand within the image should not affect the feature vector. This could be enforced by

forming a local histogram of the local orientations. This should treat each orientation element the same, independent of location.

Therefore, orientation analysis should give robustness in illumination changes while histogramming will offer translational invariance. This method will work if examples of the same gesture map to similar orientation histograms, and different gestures map to substantially different histograms.

## 7.2 *Orientation Histograms*

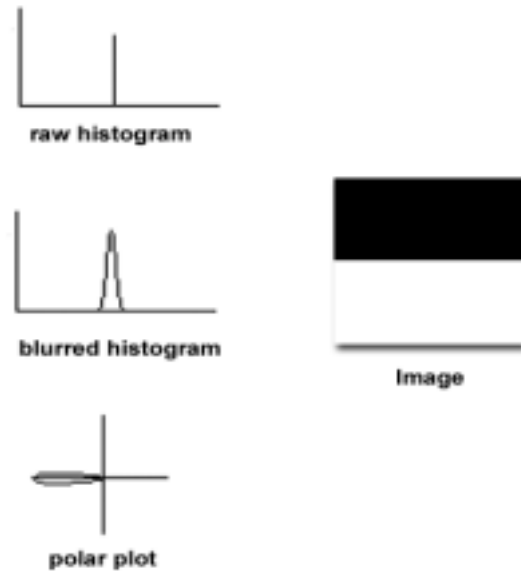
We want gestures to be the same regardless of where they occur with the images borders. To achieve this we will ignore position altogether, and tabulate a histogram of how often each orientation element occurred in the image. Clearly, this throws out information and some distinct images will be confused by their orientation histograms. In practice, however, one can choose a set of training gestures with substantially different orientation histograms from each other.

One can calculate the local orientation using image gradients. I used two 3 – tap x and y derivative filters. The outputs of the x and y derivative operators will be dx and dy. Then the gradient direction is **atan (dx, dy)**. I had decided to use the edge orientation as the only feature that will be presented to the neural network. The reason for this is that if the edge detector was good enough it would have allowed me to test the network with images from different databases. Another feature that could have been extracted from the image would be the gradient magnitude using the formula below

$$\sqrt{dx^2 + dy^2}$$

This would lead though to testing the algorithm with only similar images. Apart from this the images before resized should be of approximately the same size. This is the size of the hand itself in the canvas and not the size of the canvas. Once the image has been processed the output will be a single vector containing a number of elements equal to the number of bins of the orientation histogram.

Figure 11 shows the orientation histogram calculation for a simple image. Blurring can be used to allow neighboring orientations to sense each other.



**Figure 11** : Orientation histogram

## 7.21 Operation

The program can be 'divided' in 6 steps. Lets examine them one by one.

### Step1

The first thing for the program to do is to read the image database. A *for* loop is used to read an entire folder of images and store them in MATLAB's memory. The folder is selected by the user from menus. A menu will firstly pop-up asking you whether you want to run the algorithm on test or train sets. Then a second menu will pop-up for the user to choose which ASL sign he wants to use.

### Step2

Resize all the images that were read in Step1 to 150x140 pixels. This size seems the optimal for offering enough detail while keeping the processing time low.

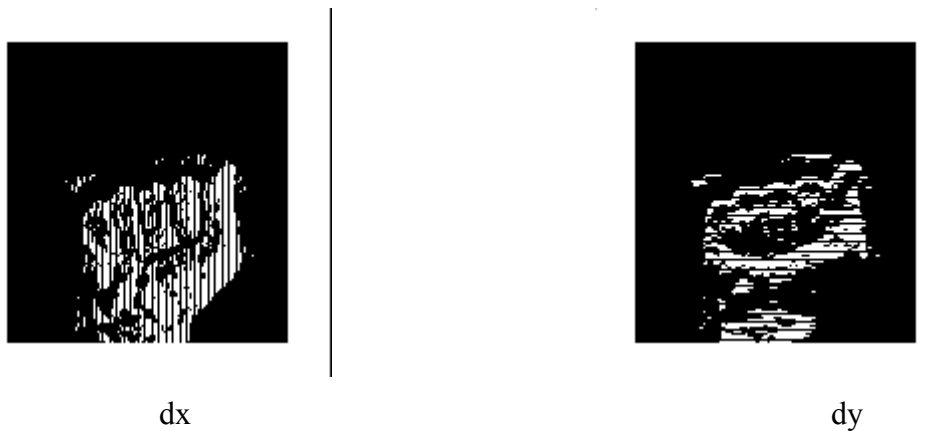
### Step3.

Next thing to do is to find the edges. As mentioned before 2 filters were used.

For the x direction  $x = [0 \ -1 \ 1]$

For the y direction  $y = \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}$ , which is the same as x but transposed and multiplied by  $-1$ .

In figure 12 you can see two images of the result of convolving an ASL sign of B with the x-filter and y-filter.



**Figure 12:** X-Y filters

As this would mean the only feature extracted from the images it had to offer enough discrimination among them. From the images above it doesn't seem like a good edge detector. In fact it doesn't look like an edge detector that much. During the beginning of the algorithm's writing I have experimented with all the known edge detectors (Sobel, Prewitt, Roberts and Canny). I have also tried various line detectors and combination of those two. An interesting way of testing various edge detectors and changing their values is an *Adobe Photoshop* filter which lets you input the values of the mask in a matrix and see the result on the image in real time. Even though some combinations of Canny (or Sobel) with line detection algorithms produced a very good result – that is visual- it wouldn't offer too much differentiation between the images. Of course for all those operators threshold values and some other parameters can be set as well. Therefore there was a lot of combinations to be tried but as mentioned before either the differentiation wasn't enough or in extreme conditions it was too much. In this case every image would classify in a class of each own or at all.

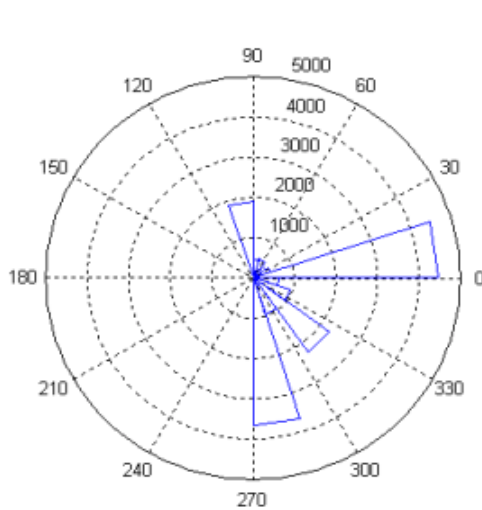


**Step 4**

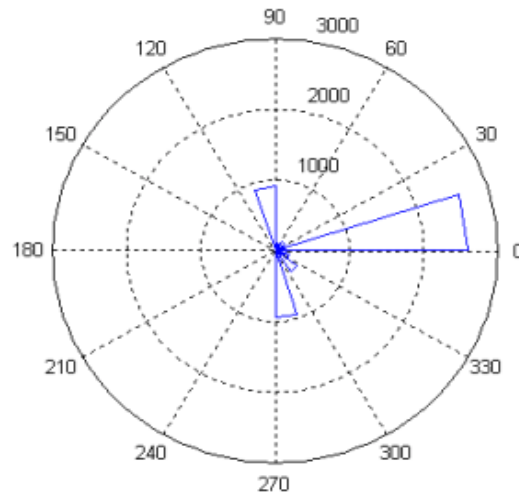
Dividing the two resulting matrices (images)  $dx$  and  $dy$  element by element and then taking the  $\text{atan}(\tan^{-1})$ . This will give the gradient orientation.

**Step 5**

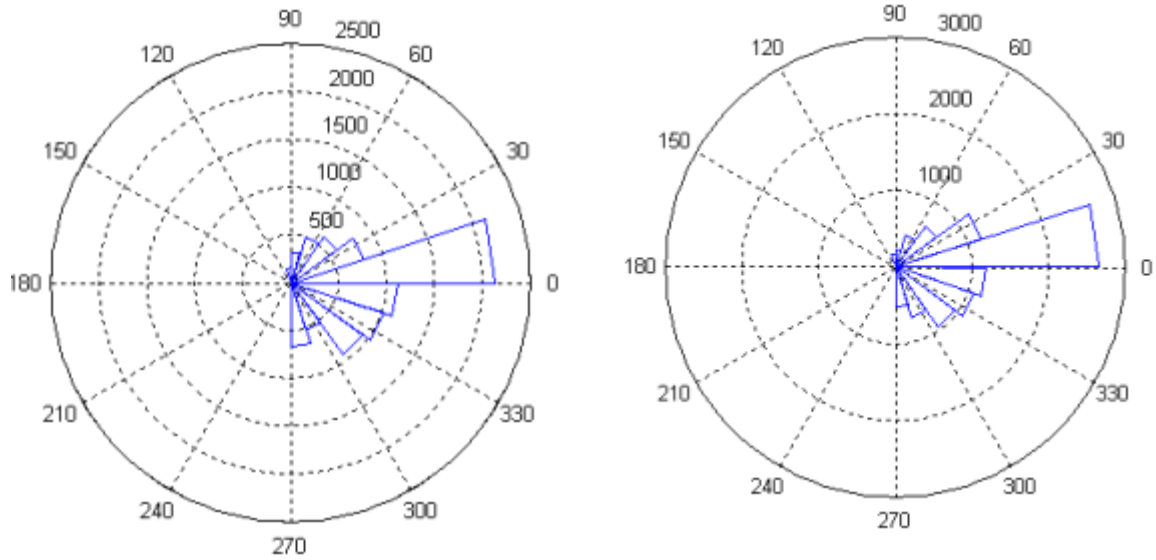
Then the MATLAB function *im2col* is called to rearrange the image blocks into columns. This is not a necessary step but it has to be done if we want to display the orientation histogram. *Rose* creates an angle histogram, which is a polar plot showing the distribution of values grouped according to their numeric range. Each group is shown as one bin. Below we can see some examples. While developing the algorithm those histograms are the fastest way of getting a good idea how good the detection is done.



**Orientation histogram of a\_1**



**Orientation histogram of a\_2**



**Orientation histogram of 5\_1**

**Orientation histogram of 5\_2**

**Figure 13:** Orientation histogram examples

Here you can see the original images that generated the histograms above in the same order



**Original image a\_1**



**Original image a\_2**



**Original image 5\_1**



**Original image 5\_2**

**Figure 14**

## Step 6

Converting the column matrix with the radian values to degrees. This way we can scan the vector for values ranging from  $0^\circ$  to  $90^\circ$ . This is because for real elements of  $X$ ,  $\text{atan}(X)$  is in the range  $[-\pi/2, \pi/2]$ .

This can also be seen from the orientation histograms where values come up only on the first and last quarter.

Determining the number of the histogram bins was another issue that was solved by experimenting with various values. I have tried with 18 20 24 and 36 bins. What I was looking for was the differentiation (or not) among the images. At the same time I was thinking of the neural network itself as this vector would be the input to the network. The smaller the vector the faster the processing. Finally, the actual resolution of each bin was set to  $10^\circ$ , which means 19 bins.

The algorithms development was organized having in mind MATLAB weaknesses. The major one is speed. MATLAB is perfect for speeding up the development process but it can be very slow on execution when bad programming practices have been employed. Nested loops slow down the program considerably. It is probably because MATLAB is built on loops. Therefore, unnecessary back-tracking was avoided and even some routines were written in full instead of using *for* loops. The code is not much in any case.

The same techniques were put in practice for the following program as well.

## 7.2 Perceptron implementation in MATLAB

### 7.31 Learning Rules

We will define a *learning rule* as a procedure for modifying the weights and biases of a network. (This procedure may also be referred to as a training algorithm.) The learning rule is applied to train the network to perform some particular task. Learning rules in the MATLAB toolbox fall into two broad categories: supervised learning and unsupervised learning.

Those two categories were described in detail in previous chapter. The algorithm has been developed using supervised learning.

In *supervised learning*, the learning rule is provided with a set of examples (the *training set*) of proper network behavior: where is an input to the network, and is the corresponding correct (*target*) output. As the inputs are applied to the network, the network outputs are compared to the targets. The learning rule is then used to adjust the weights and biases of the network in order to move the network outputs closer to the targets. The perceptron learning rule falls in this supervised learning category.

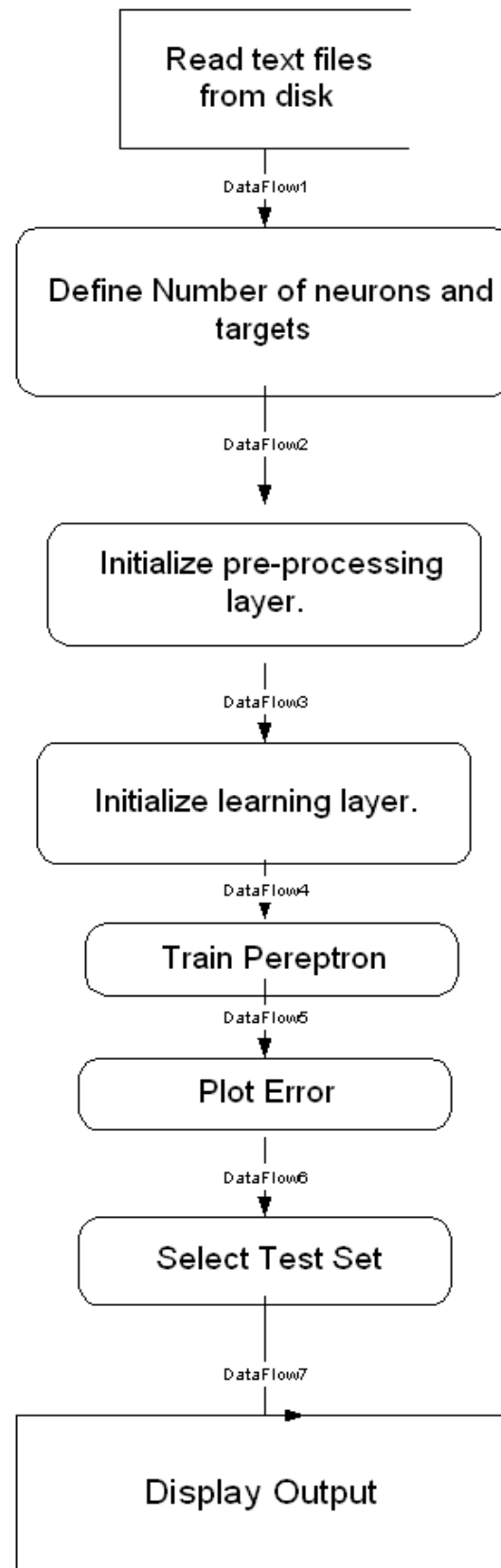
### 7.32 Linear Classification

Linear networks can be trained to perform linear classification with the function **train**. This function applies each vector of a set of input vectors and calculates the network weight and bias increments due to each of the inputs. Then the network is adjusted with the sum of all these corrections. We will call each pass through the input vectors an *epoch*. This contrasts with **adapt**, which adjusts weights for each input vector as it is presented.

**Adapt** is another function in MATLAB for training a neural network. I was using this at the first stages when I was using a back-propagation network. Their main difference is that with **train** only batch training (updating the weights after each presenting the complete data set) can be used, while with **adapt** you have the choice of batch and incremental training (updating the weights after the presentation of each single training sample). **Adapt** supports far less training functions. Since I didn't have a very good reason to go for incremental training I decided to use **train** which is more flexible.

Finally, **train** applies the inputs to the new network, calculates the outputs, compares them to the associated targets, and calculates a mean square error. If the error goal is met, or if the maximum number of epochs is reached, the training is stopped and **train** returns the new network and a training record. Otherwise **train** goes through another epoch. The LMS algorithm converges when this procedure is executed.

Below in figure 15 we can see a flow chart of the perceptron algorithm. It is always operating on the same order. There is a graphical interface only for selecting the test set as it is the only user input. The error is also displayed graphically for each epoch that is through.

**Figure 15:** Perceptron flow-chart

### 7.33 Pre-processing

We know that a feed-forward can represent an arbitrary functional mapping so it is possible to map raw data directly to the required outputs. In practice though it is preferable to apply pre-processing transformations to the input vectors before they are presented to the neural network.

```
%Initialize pre-processing layer.  
W1,b1] = initp(P,S1);%  
%Initialize learning layer.  
[W2,b2] = initp(S1,T);  
%TRAINING THE PERCEPTRON  
%The first layer is used to preprocess the input vectors:  
A1 = simup (P,W1,b1);
```

As we can see from the MATLAB source code above the first layer is used to preprocess the input data.

The network takes  $P$  as the input data processes and transforms it to  $A1$ .  $A1$  is  $85 \times 24$  matrix. The optimal number of neurons on the network has been decided to be 85. Hence, it transforms the decimal values of  $P$  to an  $85 \times 24$  binary matrix.

$A1$  will then be the input to the hidden layer that will train the network.

Reduction in the dimensionality of the input raw data is probably the most important reason for pre-processing. Input pre-processing can take various forms. The input vectors are pre-processed in the image-processing program as explained before. This was the feature extraction part of the program, which in a way is pre-processing step of great importance. In a different case if the images were presented to the network in their original form, the training and simulation would have taken hours or days for greater databases and high-resolution images. Instead we present to the network 19-element vectors, each one representing a single image. A dataset of 24 images –after the first layer preprocessing- contains  $85 \times 24 = 2040$  elements. A single image  $256 \times 256$  consists of 65536 pixels.

Therefore pre-processing either in the network itself or outside can save of computational resources and time. Of course we must be careful when we extract information this way because at the same time we lose information. Consequently, a good pre-processing strategy will retain enough critical information for correct training and classification.

### 7.34 Input , Target , Test and Output files

The input, test and target vectors are saved on the hard drive in text files. All data is stored in a single column. MATLAB can tell where one vector ends and another starts simply by writing so in the *fscanf* command as shown below.

```
TS1 = fscanf(fid,'%f',[19,inf]);  
fid = fopen('target.txt','rt');
```

Formatting those text files can be time consuming but once the training and target files are ready the test files can be created when needed.

The output is displayed on the command line.

The error is calculated by subtracting the output A from target T. Then the sum-squared error is calculated. Below some error graphs of the network training are shown.



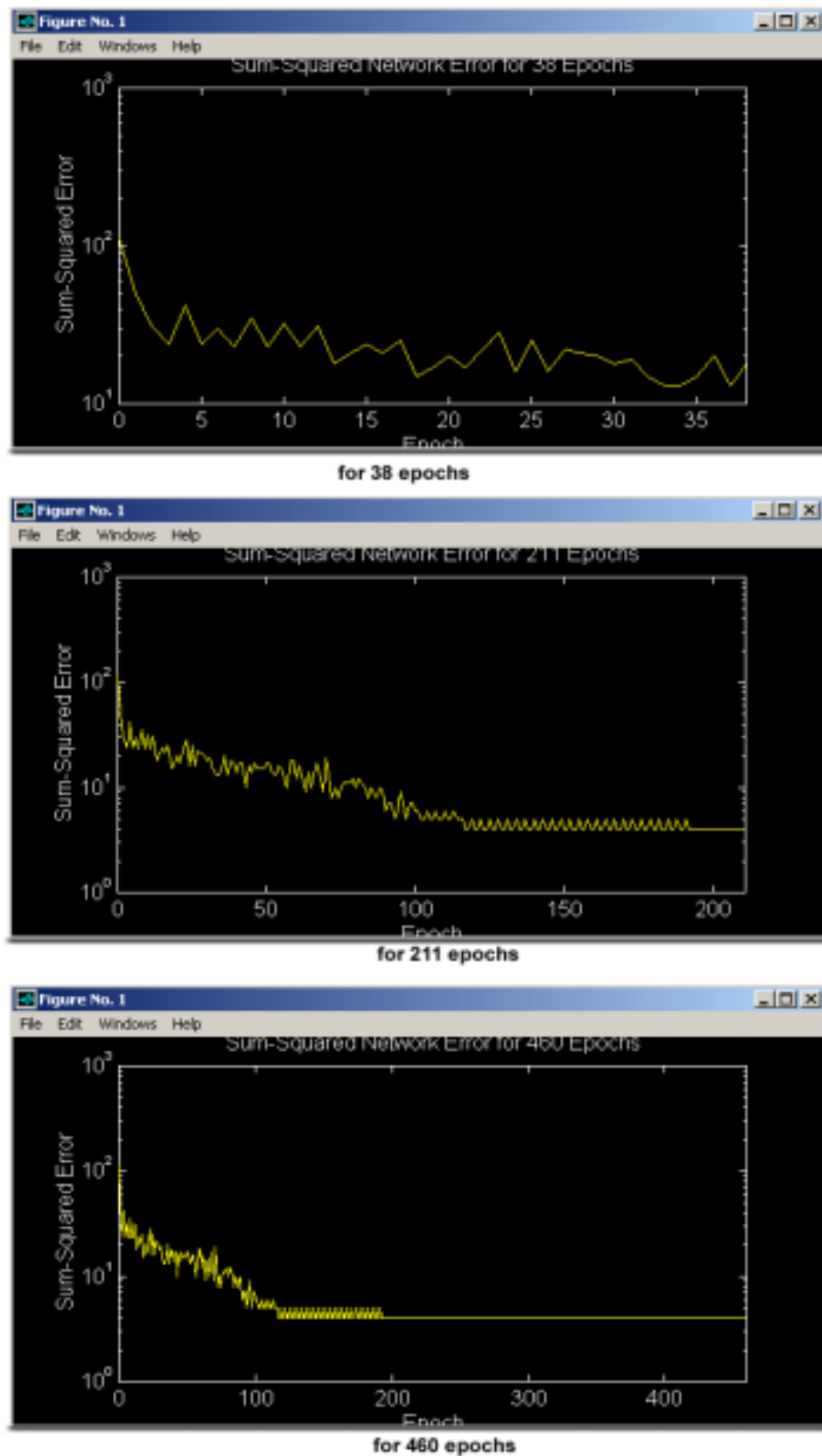


Figure 16: Train Error

Here is the error for the first 20 epochs of training

#### Training starts

```
[W2,b2,epochs,errors] = trainp (W2,b2,A1,T,TP);
```

TRAINP: 0/500 epochs, SSE = 86.

TRAINP: 1/500 epochs, SSE = 42.

TRAINP: 2/500 epochs, SSE = 24.

TRAINP: 3/500 epochs, SSE = 64.

TRAINP: 4/500 epochs, SSE = 42.

TRAINP: 5/500 epochs, SSE = 21.

TRAINP: 6/500 epochs, SSE = 46.

TRAINP: 7/500 epochs, SSE = 28.

TRAINP: 8/500 epochs, SSE = 21.

TRAINP: 9/500 epochs, SSE = 20.

TRAINP: 10/500 epochs, SSE = 36.

TRAINP: 11/500 epochs, SSE = 27.

TRAINP: 12/500 epochs, SSE = 20.

TRAINP: 13/500 epochs, SSE = 24.

TRAINP: 14/500 epochs, SSE = 22.

TRAINP: 15/500 epochs, SSE = 28.

TRAINP: 16/500 epochs, SSE = 23.

TRAINP: 17/500 epochs, SSE = 23.

TRAINP: 18/500 epochs, SSE = 15.

TRAINP: 19/500 epochs, SSE = 14.

TRAINP: 20/500 epochs, SSE = 19.

The error is constantly reducing until it converges to 0. When it does it will stop and start the testing process. In the case of reaching the 500 predetermined epochs without coming down to 0 it will test the network as is.

### 7.35 Limitations and Cautions

Perceptron networks should be trained with *adapt*, which presents the input vectors to the network one at a time and makes corrections to the network based on the results of each presentation. Use of *adapt* in this way guarantees that any linearly separable problem will be solved in a finite number of training presentations. Perceptrons can also be trained with the function *train*. When *train* is used for perceptrons, it presents the inputs to the

network in batches, and makes corrections to the network based on the sum of all the individual corrections. Unfortunately, there is no proof that such a training algorithm converges for perceptrons. On that account the use of *train* for perceptrons is not recommended. The algorithm developed does not always converge. After experimenting with various perceptron parameters, I realized that the number of neurons in the processing layer is altering the algorithms performance in terms of converging or not. Usually the most nodes the faster it converges. This is not simply proportional though as a very big number of nodes could lead to larger training errors and therefore less efficient classification.

When it doesn't converge the sum of squared errors reaches a maximum value of 2 (0 when it converges).

Therefore, convergence is not the qualifying criteria of the algorithm since with it can work perfectly with such a small error.

## 8 RESULTS

Here are some tables displaying the results obtained from the program. Sign images of the same letter are grouped together on every table. The table gives us information about the pre-processing operations that took place (i.e. blurring, noise, translation) and also if the image belongs to the same database with the training images.

The amount of each filter is also recorded so we can estimate the maximum values of noise the network can tolerate. This of course varies from image to image. The result also varies for every time the algorithm is executed. The variance is very small but it is there. So we cannot easily draw conclusions and set a certain threshold above which we can tell that the network will not classify correctly. It all comes down to the application again.

### Form of results:

The results come out of the network in column format. Each column is a classified image vector. The position of the '1' in the vector among the '0s' indicates which sign it is. Therefore there should be only one '1' in every vector, but this is not always the case. As you will see from the tables below there are situations that the perceptron cannot converge to a single solution but it gives two possible solutions. In almost all of those cases one of the classifications is correct. There are few others though that the vector is not classified at all.

This is a table with all the target vectors that have used to train the network and confirm correct classification or not.

0	1	5	A	H	L	V	W
0	0	0	0	0	0	0	1
0	0	0	0	0	0	1	0
0	0	0	0	0	1	0	0
0	0	0	0	1	0	0	0
0	0	0	1	0	0	0	0
0	0	1	0	0	0	0	0
0	1	0	0	0	0	0	0
1	0	0	0	0	0	0	0

**Table 3:** Target Vectors

Below you can see some result vectors as they come out of MATLAB: This is a test set for 'L'.

a2 =

```

0  0  0  0  0  0  0  0  0  0  0  0
0  0  0  0  0  0  0  0  0  0  0  0
1  1  1  1  1  1  1  1  1  1  1  1
0  0  0  0  0  0  0  0  0  0  0  0
0  1  1  0  0  1  0  0  1  0  0  0
0  0  0  0  0  0  0  0  0  0  0  0
0  0  0  0  0  0  0  0  0  0  0  0
0  0  0  0  0  0  0  0  0  0  0  0

```

Therefore the 1<sup>st</sup>, 4<sup>th</sup>, 5<sup>th</sup>, 7<sup>th</sup>, 8<sup>th</sup>, 10<sup>th</sup>, 11<sup>th</sup> and 12<sup>th</sup> vectors are correctly classified. The rest are classified as either 'L' or 'A'.

All the filtering operations were performed in *Adobe Photoshop 5.0*. hence the amount of each one was measured.

- Blur – Measured in pixel radius
- Motion Blur – Like taking photo of moving car. Specify angle and blur radius
- Noise - Either Uniform or Gaussian (stronger)
- Same Image – From the same image database
- Translated – move the signaling hand in the canvas to test translational invariance
- Classified
  - ✓ : correctly classified
  - ✗ : not classified
  - W** : classified as ‘**W**’ (wrongly)
  - L – a** : classified as either an ‘**L**’ or ‘**A**’ (Upper case correct)

I would like to clarify here that for any test set, image that is noted as coming from other databases is the only one. Therefore every image that is indicated so represents a specific database on each own.

Below are the result tables. They start with digit ‘0’ and continue with the rest of the digits and letters.

## Testing – 0 (Zero)



Figure 17: Zero

Image	Blur(Radius)	Noise	Same Image	Translated	Classified
1	-	-	NO	-	✓
2	-	18 u	Yes	-	✓
3	-	25 u	Yes	-	✓
4	-	25 g	Yes	-	✓
5	-	28 g	Yes	Yes	✓
6	1.0	-	Yes	-	✓
7	1.2	-	Yes	-	✓
8	1.4	-	Yes	-	✓
9	1.8	-	Yes	Yes	✓
10	-	15 u	Yes	-	✓
11	-	8 g	Yes	-	✓
12	0.8	-	Yes	-	✓
13	1.0	-	Yes	Yes	✓
14	2.0	-	Yes	-	✗
15	1.0	-	NO	-	1
16	-	-	NO	-	✓
17	-	-	NO	-	1

Table 4: Zero Test Results

For '0' the classification error is very small. For the 5<sup>th</sup> image the amount of noise (28 gaussian) is very high but it still classifies correctly. It will not tolerate blurring above a pixel radius of 2.0 though.

On the other hand translation doesn't seem to cause any problems.

## Testing – 1 (One)



Figure 18: One

Image	Blur	Noise	Same Image	Translated	Classified
1	Sharpen	-	Yes	Yes	✓
2	Sharpen More	-	Yes	-	✗
3	1.2	-	Yes	-	✓
4	1.5	-	Yes	-	0
5	2.0	10 g	Yes	Yes	5
6	-	15 g	Yes	-	✓
7	-	20 g	Yes	-	0
8	-	22 u	Yes	-	✓
9	-	25 u	Yes	-	0
10	1.0	-	Yes	-	✓
11	-	20 u	Yes	-	✓
12	Motion blur 90degrees, 2	-	NO	Yes	✓
13	-	-	NO	Yes	✓
14	-	-	NO	-	0

Table 5: One Test Results

Sharpening causes plenty of distortion on the images, especially on the edges. That is way not much sharpening will be tolerated. We can see from the table above that for great amounts of noise or blurring the perceptron classifies '1' signs as '0s'.

On image 5 a combination of blurring and gaussian noise is causing the network to classify the sign completely wrongly.



## Testing – L

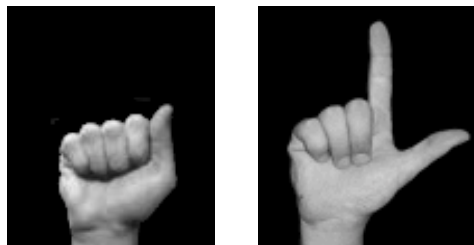


Figure 19: L

Image	Blur	Noise	Same Image	Translated	Classified
1	-	12 g	Yes	-	✓
2	-	15 g	Yes	-	L-a
3	-	20 g	Yes	-	L-a
4	-	10 u	Yes	-	✓
5	-	15 u	Yes	Yes	✓
6	-	20 u	Yes	Yes	L-a
7	1.2	-	Yes	-	✓
8	1.5	-	Yes	-	✓
9	2.0	-	Yes	-	✓
10	-	5 g	Yes	Yes	✓
11	-	10 g	Yes	-	L-a
12	0.7	-	Yes	-	✓
13	0.9	-	Yes	-	✓
14	-	-	NO	-	✓

Table 6: L Test Results

‘L’ classification is pretty good apart from confusion with ‘A’. Below you can see the two signs in their original form. Comparing images in their raw form with their vectors representations is not very safe. Roughly though we may be able to see where the problem starts from. In this case the fist and finger arrangements are quite similar apart from the index finger. Fingers are critical because they generate many edges.



‘A’

‘L’

Figure 20

## Testing – H



Figure 21: H

Image	Blur	Noise	Same Image	Translated	Classified
1	Motion Blur 45 degrees, 1.5 r	-	Yes	-	✓
2	Motion blur 90degrees 2r	-	Yes	-	✓
3	-	22 u	Yes	Yes	1
4	-	25 u	Yes	Yes	✗
5	-	10 g	Yes	-	✓
6	-	15 g	Yes	-	✗
7	-	20 g	Yes	-	✗
8	-	22 g	Yes	-	✓
9	Dust : 2r, T: 126 levels	-	Yes	-	✓
10	-	20 u	Yes	Yes	✓
11	-	-	NO	Yes	H-a
12	-	-	NO	Yes	H-a
13	-	-	NO	-	H-v
14	-	-	NO	-	H-a
15	-	-	NO	Yes	H-a

Table 7: H Test Results

An interesting result here is 3 and 4. For 22 levels of uniform noise the vector is wrongly classified as a '1'. With a little bit of extra noise the vector becomes unknown to the network, which doesn't classify it at all.

Image 7 doesn't classify at all with 20 levels of gaussian noise. The same image with more noise classifies correctly. There is no logical explanation for that. For 10 consecutive tests of the same set the probability of getting all 10 times the exact output is very low. The initialization of the weights and biases are left on MATLAB which means the network cannot be 100% stable.

## Testing – A

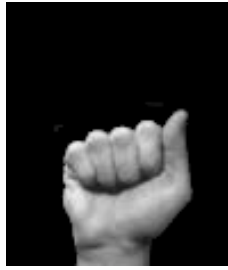


Figure 22: A

Image	Blur	Noise	Same Image	Translated	Classified
1	1.0	5 g	Yes	Yes	✓
2	1.5	10 g	Yes	Yes	✓
3	-	-	NO	-	✓
4	-	-	NO	-	✓
5	-	-	NO	Yes	✓
6	-	-	NO	-	✗
7	-	-	Yes	-	✓
8	-	-	NO	Yes	✓
9	2.0	-	NO	-	✓
10	-	-	NO	-	A-v
11	-	-	NO	-	✓

Table 8: L Test Results

While I was testing ‘A’ I noted that it would never classify wrongly. This was not expected because ‘L’ classified as ‘A’ many times. Hence I anticipated the reverse for ‘A’. So most of the test vectors are coming from various databases and the result is still acceptable.

## Testing – V



Figure 23: V

Image	Blur	Noise	Same Image	Translated	Classified
1	-	5 u	Yes	Yes	✓
2	-	10 u	Yes	-	✓
3	-	15 u	Yes	-	✓
4	-	20 u	Yes	-	✗
5	0.5	-	Yes	-	✓
6	1.0	-	Yes	Yes	1
7	1.2	-	No	-	V-1
8	-	-	No	Yes	✓

Table 9: L Test Results

‘V’ classifies correctly until the uniform noise level reaches 20. quite sensitive to blurring too. The only difference between a ‘V’ and a ‘I’ sign the middle finger. That is why for blurring over 1.0 pixels the network identifies it as ‘I’.

## Testing – W

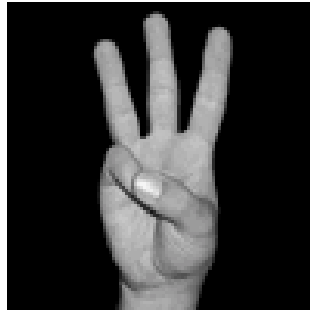


Figure 24: W

Image	Blur	Noise	Same Image	Translated	Classified
1	0.8	-	Yes	Yes	✓
2	-	-	No	-	W-v
3	-	0.8 u	No	-	L
4	Motion blur 45degrees – 3 radius	-	Yes	-	✓
5	-	Dust Radius 4, T:126	Yes	Yes	W-v

Table 10: L Test Results

The third vector is classified as ‘L’. The original image comes from a different database. What happened here is that the angle between the hand and the camera is slightly different from the usual one. ‘W’ and ‘V’ are very similar as you can see from the images. In the previous test set (‘V’) V was sometimes coming up as 1 and never as W. It is easier for the network to completely discard a blurred finger than come up with one that doesn’t exist.

## Testing – G

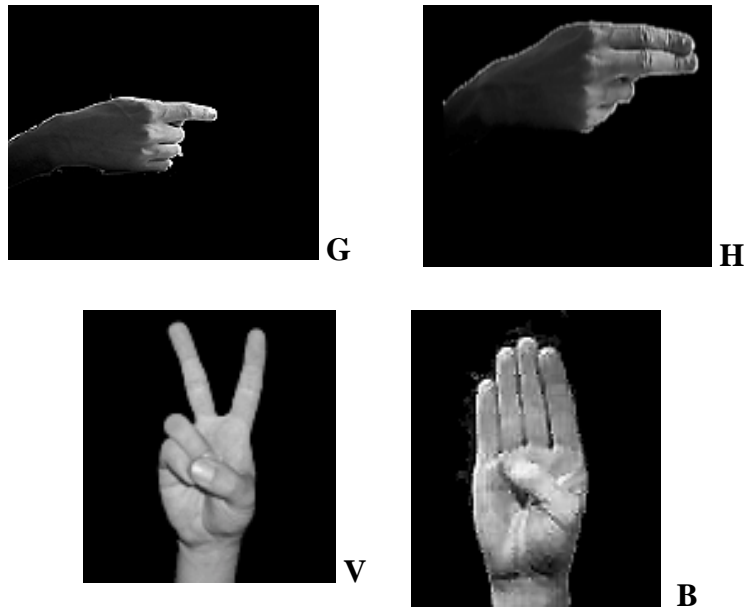
Image	Classified
1	×
2	H
3	×

**Table 11: G Test Results**

## Testing – B

Image	Classified
1	V
2	×
3	×

**Table 12: B Test Results**



**Figure 25: Signs G, H, V, B**

The two tables above display two signs ‘G’ and ‘B’ that have not been used to train the network but only to test it. Below the tables we can see the signs that the network has classified them to. (when classified). The fact that it didn’t classify the vectors in two out of the three cases is encouraging. This may not be a full test set of ASL but it proves that if a test vector has not been included in the train set as well it will either not classify at all

or it will to the nearest match. '***H***' is very similar visually to '***G***' and it would be very hard to distinguish between those two in any case.

## Conclusion

The idea of the project got started from a *McConnel's* idea of orientation histograms. Many researchers found the idea interesting and tried to use it in various applications. From hand recognition to cat recognition and geographical statistics.

My supervisor and I had the idea of trying to use this technique in conjunction with Neural Networks. In other approaches of pattern recognition that orientation histograms have been used different ways of comparing and classifying were employed. Euclidean distance is a straight forward approach to it. It is efficient as long as the data sets are small and not further improvement is expected. Another advantage of using neural networks is that you can draw conclusions from the network output. If a vector is not classified correct we can check its output and work out a solution.

As far as the orientation algorithm is concerned it can be further improved. The main problem is how good differentiation one can achieve. This of course is dependent upon the images but it comes down to the algorithm as well. Edge detection techniques are keep changing while line detection can solve some problems. One of the ideas that I had lately is the one of tangents but I don't know if it is feasible and there is not time of developing it.

To say that I have come to robust conclusions at the end of the project is not safe. This is possible only for the first part of the project. Regardless of how many times you run the program the output vector will always be the same. This is not the case with the perceptron. Apart from not being 100% stable there are so many parameters (e.g. number of layers, number of nodes) that one can play with that finding the optimal settings is not that straight forward. As mentioned earlier it all comes down to the application. If there is a specific noise target for example you can work to fit this specifications.

My major goal was speed and the avoidance of special hardware. This was achieved although it would be faster if written in C / C++ but the complexity of the design and



















implementation would have been much higher. MATLAB is slower but allows its users to work faster and concentrate on the results and not on the design.

The thesis and the MATLAB source code will be available for download from the Surrey server:

[http://www.ee.surrey.ac.uk/People/T.Windeatt/student\\_projects/klimis/](http://www.ee.surrey.ac.uk/People/T.Windeatt/student_projects/klimis/)

## References:

- 📖 <http://www.cs.rug.nl/~peterkr/FACE/face.html>
- 📖 <http://www.hav.com/>
- 📖 [http://www.dacs.dtic.mil/techs/neural/neural\\_ToC.html](http://www.dacs.dtic.mil/techs/neural/neural_ToC.html)
- 📖 <http://www.tk.uni-linz.ac.at/~schaber/ogr.html>
- 📖 <http://vismod.www.media.mit.edu/vismod/classes/mas622/projects/hands/>
- 📖 <http://aiintelligence.com/aii-info/techs/nn.htm#What>
- 📖 Christopher M. Bishop, “Neural networks for Pattern Recognition” Oxford, 1995.
- 📖 William T. Freeman, Michael Roth, “ Orientation Histograms for Hand Gesture Recognition” *IEEE Intl. Wkshp. On Automatic Face and Gesture Recognition, Zurich, June, 1995.*
- 📖 Maria Petrou, Panagiota Bosdogianni, “Image Processing, The Fundamentals”, Wiley
- 📖 Vladimir I. Pavlovic, Rajeev Sharma, Thomas S Huang, “ Visual Interpretation of Hand Gestures for Human-Computer Interaction : A review” *IEEE Transactions of pattern analysis and machine intelligence, Vol 19, NO 7, July 1997*
- 📖 Srinivas Gutta, Ibrahim F. Imam, Harry Wechsler, “ Hand gesture Recognition Using Ensembles of Radial Basis Functions (RBF) Networks and Decision Trees” *International Journal of Pattern Recognition and Artificial Intelligence, Vol 11 No.6 1997.*
- 📖 Simon Haykin, “Neural Networks, A comprehensive Foundation”, Prentice Hall
- 📖 Duane Hanselman, Bruce Littlefield, “Mastering MATLAB, A comprehensive tutorial and reference”, Prentice Hall

-  <http://www.cs.cmu.edu/afs/cs/project/cil/ftp/html/vision.html>
-  <http://www.merl.com/people/freeman/>
-  <http://deafness.about.com/health/deafness/gi/dynamic/offsite.htm?site=http://www.lifeprint.com/ASL101/welcome.htm>
-  <http://deafness.about.com/health/deafness/>
-  <http://www.cis.ohio-state.edu/~szhu/SCTV99.html>
-  <http://hebb.cis.uoguelph.ca/~skremer/Teaching/27642/BP/node3.html>
-  [http://www.math.tau.ac.il/~stainvas/inna\\_norm.html](http://www.math.tau.ac.il/~stainvas/inna_norm.html)
-  [http://www.cat.csiro.au/automation/staff/nil/csiro/publications/rep\\_edge/tpap.html](http://www.cat.csiro.au/automation/staff/nil/csiro/publications/rep_edge/tpap.html)
-  <http://www.ncrg.aston.ac.uk/~woonw/subpages/fusion.html>
-  <http://classes.monterey.edu/CST/CST332-01/world/Mat/Edge.htm>
-  <http://www.generation5.org/perceptron.shtml>
-  <http://www.cs.bgu.ac.il/~omri/Perceptron/>
-  [http://www.phy.syr.edu/courses/modules/MM/Neural/neur\\_train.html](http://www.phy.syr.edu/courses/modules/MM/Neural/neur_train.html)
-  <http://hebb.cis.uoguelph.ca/~skremer/Teaching/27642/BP/node2.html>
-  <http://www.willamette.edu/~gorr/classes/cs449/Classification/perceptron.html>
-  <http://www.bus.olemiss.edu/jjohnson/share/mis695/perceptron.htm>



## APPENDIX

### *Part1 : Image Processing Program*

```
%%%Klimis Symeonidis - Msc signal Processing Communications - Surrey  
Univertsity%%%%  
%%Orientation Histograms%%%%  
clc
```

```
% Select from menu Test or Train Sets  
F = MENU('Choose a database set','Test Set','Train Set');  
if F==1
```

```
%% Select Test Set  
K = MENU('Choose a file','Test A','Test V','Test W','Test 0','Test  
1','Test From other DBs');  
if K == 1  
    loop=4  
    for i=1:loop  
        string = ['test\a\' num2str(i) '.tif'];  
        Rimages{i} = imread(string);  
    end
```

```
elseif K == 2  
    loop=7  
    for i=1:loop  
        string = ['test\v\' num2str(i) '.tif'];  
        Rimages{i} = imread(string);  
    end
```

```
elseif K == 3  
    loop=5  
    for i=1:loop  
        string = ['Test\W\' num2str(i) '.tif'];  
        Rimages{i} = imread(string);  
    end
```

```
elseif K == 4  
    loop=5  
    for i=1:loop  
        string = ['test\0\' num2str(i) '.tif'];  
        Rimages{i} = imread(string);  
    end
```

```
elseif K == 5  
    loop=4  
    for i=1:loop  
        string = ['test\1\' num2str(i) '.tif'];  
        Rimages{i} = imread(string);  
    end
```

```
elseif K == 6  
    loop=13  
    for i=1:loop  
        string = ['test\otherdb\' num2str(i) '.tif'];
```

```

        Rimages{i} = imread(string);
    end
end
end;
%% Select Train Set
if F==2
    loop=3    %Set loop to 3. All train sets have 3 images
    L = MENU('Choose a file', 'Train A', 'Train V', 'Train W', 'Train
0', 'Train 1');
    if L == 1

        for i=1:loop
            string = ['train\a\' num2str(i) '.tif'];
            Rimages{i} = imread(string);
        end

    elseif L == 2
        for i=1:loop
            string = ['train\v\' num2str(i) '.tif'];
            Rimages{i} = imread(string);
        end

    elseif L == 3
        for i=1:loop
            string = ['train\W\' num2str(i) '.tif'];
            Rimages{i} = imread(string);
        end

    elseif L == 4
        for i=1:loop
            string = ['train\0\' num2str(i) '.tif'];
            Rimages{i} = imread(string);
        end

    elseif L == 5
        for i=1:loop
            string = ['train\1\' num2str(i) '.tif'];
            Rimages{i} = imread(string);
        end
    end
end
end

% Resize all images to 150x140
T{i}=imresize(Timages{i},[150,140]);

    x = [0 -1 1];                %x-derivative filter
    y = [0 1 -1]';              %y-derivative filter

% returns only those parts of the convolution that can be computed
without assuming that the array A is zero-padded
dx{i} = convn(T{i},x,'same');

% returns only those parts of the convolution that can be computed

```

```
without assuming that the array A is zero-padded
dy{i} = convn(T{i},Y,'same');

% divide the two matrices element by element to find gradient
orientation
gradient{i} = dy{i} ./dx{i};
theta{i} = atan(gradient{i});           %find the atan of 'gradient'

% turn the matrix into a column vector

c1{i}= im2col(theta{i},[1 1],'distinct');

% convert radians to degrees
N{i} = (c1{i}*180)/3.14159265359;

% read and store the orientation magnitude every 10 degrees
c1{i}=(N{i}>0)&(N{i}<10);
s1{i}=sum(c1{i});

c2{i}=(N{i}>10.0001)&(N{i}<20);
s2{i}=sum(c2{i});

c3{i}=(N{i}>20.0001)&(N{i}<30);
sum(c3{i});
s3{i}=sum(c3{i});

c4{i}=(N{i}>30.0001)&(N{i}<40);
sum(c4{i});
s4{i}=sum(c4{i});

c5{i}=(N{i}>40.0001)&(N{i}<50);
sum(c5{i});
s5{i}=sum(c5{i});

c6{i}=(N{i}>50.0001)&(N{i}<60);
sum(c6{i});
s6{i}=sum(c6{i});

c7{i}=(N{i}>60.0001)&(N{i}<70);
sum(c7{i});
s7{i}=sum(c7{i});

c8{i}=(N{i}>70.0001)&(N{i}<80);
sum(c8{i});
s8{i}=sum(c8{i});

c9{i}=(N{i}>80.0001)&(N{i}<90);
sum(c9{i});
s9{i}=sum(c9{i});

c10{i}=(N{i}>90.0001)&(N{i}<100);
sum(c10{i});
s10{i}=sum(c10{i});

c11{i}=(N{i}>-89.9)&(N{i}<-80);
sum(c11{i});
s11{i}=sum(c11{i});
```

```
c12{i}=(N{i}>-80.0001)&(N{i}<-70);
sum(c12{i});
s12{i}=sum(c12{i});

c13{i}=(N{i}>-70.0001)&(N{i}<-60);
sum(c13{i});
s13{i}=sum(c13{i});

c14{i}=(N{i}>-60.0001)&(N{i}<-50);
sum(c14{i});
s14{i}=sum(c14{i});

c15{i}=(N{i}>-50.0001)&(N{i}<-40);
sum(c15{i});
s15{i}=sum(c15{i});

c16{i}=(N{i}>-40.0001)&(N{i}<-30);
sum(c16{i});
s16{i}=sum(c16{i});

c17{i}=(N{i}>-30.0001)&(N{i}<-20);
sum(c17{i});
s17{i}=sum(c17{i});

c18{i}=(N{i}>-20.0001)&(N{i}<-10);
sum(c18{i});
s18{i}=sum(c18{i});

c19{i}=(N{i}>-10.0001)&(N{i}<-0.0001);
sum(c19{i});
s19{i}=sum(c19{i});

D{i}= [s1{i} s2{i} s3{i} s4{i} s5{i} s6{i} s7{i} s8{i} s9{i} s10{i}
s11{i} s12{i} s13{i} s14{i} s15{i} s16{i} s17{i} s18{i} s19{i}];
end;
end;
end;

%close the waiting bar
close(w);
```



## Part2 : Neural Network Program

```
%%%Klimis Symeonidis - Msc signal Processing Communications - Surrey
Univertsity%%%
```

```
%%%Perceptron network for hand gesture classification%%%
```

```
% Turn on echoing of commands inside the script-file.
```

```
echo on
```

```
% Clear the workspace (all variables).
```

```
% clear all
```

```
% load perf24
```

```
% Clear command window.
```

```
clc
```

```
% =====
```

```
% =====
```

```
% CLASSIFICATION WITH A 2-LAYER PERCEPTRON:
```

```
% The first layer acts as a non-linear preprocessor for
```

```
% the second layer. The second layer is trained as usual.
```

```
pause % Strike any key to continue...
```

```
clc
```

```
% DEFINING THE CLASSIFICATION PROBLEM
```

```
% =====
```

```
% A matrix P defines 24 19-element input (column) vectors:
```

```
% There are 3 examples of % each character, 8 characters, so 3 x 8
= 24 input
```

```
% patterns.
```

```
% A matrix T defines the categories with target (column)
```

```
% vectors. There are 3 numerals and 5 characters so, 8 target
vectors in total.
```

```
pause % Strike any key to continue...
```

```
clc
```

```
% Open the files with the input vectors.
```

```
fid = fopen('train8.txt','rt');
```

```
P1 = fscanf(fid,'%f',[19,inf]);
```

```
P=P1;
```

```
fid = fopen('testA.txt','rt');
```

```
TS1 = fscanf(fid,'%f',[19,inf]);
```

```
fid = fopen('test0.txt','rt');
```

```
TS2 = fscanf(fid,'%f',[19,inf]);
```

```
fid = fopen('test5.txt','rt');
```

```
TS3 = fscanf(fid,'%f',[19,inf]);
```

```
fid = fopen('testL.txt','rt');
```

```
TS4 = fscanf(fid,'%f',[19,inf]);
```

```
fid = fopen('testV.txt','rt');
```

```
TS5 = fscanf(fid,'%f',[19,inf]);
```

```
fid = fopen('testW.txt','rt');
```

```
TS6 = fscanf(fid,'%f',[19,inf]);
```

```

fid = fopen('testH.txt','rt');
TS7 = fscanf(fid,'%f',[19,inf]);
fid = fopen('test1.txt','rt');
TS8 = fscanf(fid,'%f',[19,inf]);
fid = fopen('testGB.txt','rt');
TS9 = fscanf(fid,'%f',[19,inf]);

%   Open the file with the target vectors.

fid = fopen('target8.txt','rt');
T = fscanf(fid,'%f',[8,inf]);

%clc
%   DEFINE THE PERCEPTRON
%   =====

%   P has 19 elements in each column,
%   so each neuron in the hidden layer
%   needs 19 inputs.

%R1;

%   To maximize the chance that the preprocessing layer
%   finds a linearly separable representation for the
%   input vectors, it needs a lot of neurons.
%   After trying a lot of different network architectures,
%   it has been found that the optimal number of neurons for
%   the hidden layer is 85.

S1 = 85;

%   T has 5 elements in each column,
%   so 5 neurons are needed.

S2 = 5;

%   INITP generates initial weights
%   and biases for the network:

%   Initialize pre-processing layer.

[W1,b1] = initp(P,S1);

%   Initialize learning layer.

[W2,b2] = initp(S1,T);

pause % Strike any key to train the perceptron...
clc

%   TRAINING THE PERCEPTRON
%   =====

%   TRAINP trains perceptrons to classify input vectors.

%   The first layer is used to preprocess the input vectors:

```

```
A1 = simup(P,W1,b1);

% TRAINP is then used to train the second layer to
% classify the preprocessed input vectors A1.
% The TP parameter is needed by the TRAINP function
% to define the number of epochs used.
% The first argument is the display frequency and
% the second is the maximum number of epochs.

TP = [1 500];

pause % Strike any key to start the training...

%Delete everything and also reset all figure properties,
%except position, to their default values.
clf reset

%Open a new Figure (graph window), and return
%the handle to the current figure.
figure(gcf)

%Set figure size.
setfsize(600,300);

% Training begins...

[W2,b2,epochs,errors] = trainp(W2,b2,A1,T,TP);

% ...and finishes.

pause % Strike any key to see a plot of errors...
clc

% PLOTTING THE ERROR CURVE
% =====

% Here the errors are plotted with
% respect to training epochs:

ploterr(errors);

% If the hidden (first) layer preprocessed the original
% non-linearly separable input vectors into new linearly
% separable vectors, then the perceptron will have 0 errors.

% If the error never reached 0, it means a new
% preprocessing layer should be created
% (perhaps with more neurons).
pause % Strike any key to use the classifier...
clc

% USING THE CLASSIFIER
% =====

% IF the classifier WORKED we can now try to classify
```

```
% the input vectors we like using SIMUP. Lets try the
% input vectors that we have used for training.

% Create a menu, so the user can select a test set.

K = MENU('Choose a file resolution','Test A','Test 0','Test 5','Test
L','Test V','Test W','Test H','Test 1','Test GB');

% According to the choice use the appropriate variables.

if K == 1
    TS = TS1;
elseif K == 2
    TS = TS2;
elseif K == 3
    TS = TS3;
elseif K == 4
    TS = TS4;
elseif K == 5
    TS = TS5;
elseif K == 6
    TS = TS6;
elseif K == 7
    TS = TS7;
elseif K == 8
    TS = TS8;
elseif K == 9
    TS = TS9;

else
    P = 0;
    R1 = 0;
end

a1 = simup(TS,W1,b1);    % Preprocess the vector

a2 = simup(a1,W2,b2)    % Classify the vector

echo off
disp('End of Hand Gesture Recognition')
```