

# Automatic Human Emotion Detection

Retshidisitswe Lehata

Thesis presented in fulfilment  
of the requirements for the degree of  
Bachelor of Science Honours  
at the University of the Western Cape

Supervisor: Mehrdad Ghaziasgar  
Co-supervisor: Reg Dodds

This version September 23, 2017



# Declaration

I, RETSHIDISITSWE LEHATA, declare that this thesis “*Automatic Human Emotion Detection*” is my own work, that it has not been submitted before for any degree or assessment at any other university, and that all the sources I have used or quoted have been indicated and acknowledged by means of complete references.

Signature: .....

Date: .....

RETSHIDISITSWE LEHATA.



# Abstract

Customer service is a large revenue stream for some companies, therefore ensuring that they provide the best quality service is likely to be their main priority. Not all customers readily express their emotions, verbally, regarding the quality of the service they are provided. Mehrabian 1980, states that 55% of communication is facial expression. The motivation for this project is to apply an Automatic Human Emotion Detection(AHED) system in cases where an employee interacts with a customer. The AHED system focuses on emotion recognition using facial expressions. The proposed method for this project, first captures an image from the camera, the image is processed using the Viola Jones Algorithm to detect the face and the Histogram of Oriented Gradients(HOG) to extract the features from the face. Lastly the AHED system is trained using Support Vector Machines(SVM), to classify each emotion. There are six universal human expressions described by Ekman and Friesen, namely Surprise, Fear, Disgust, Anger, Happiness and Sadness. Grayscale frontal images will be used as input for the AHED system.

# Acknowledgment

This thesis is a compilation of the efforts of many people that helped and me through the years. I would first like to thank my supervisor Mehrdad Ghazi-asgar and co-supervisor Regg Dodds for encouraging me during my study. Without our weekly meetings, this work would not have been possible. At this time I would like to extend a very special thanks to Waleed Deaney for being my mentor, without his help I would certainly not be where I am today.

# Contents

Declaration . . . . .	iii
Abstract . . . . .	v
Acknowledgment . . . . .	vi
List of Figures . . . . .	ix
1. Introduction . . . . .	1
1.1 Problem Statement . . . . .	1
1.2 Proposed Solution . . . . .	2
1.3 Proposed Method . . . . .	2
2. Related Work . . . . .	3
2.1 Capture Image from the Camera and Image Processing . . . . .	3
2.2 Face Detection - Finding the Face . . . . .	4
2.3 Feature Extraction - Extract the Features from the Face . . . . .	5
2.4 Training the machine learning technique . . . . .	6
2.5 Results . . . . .	7
3. Image Processing Techniques . . . . .	8
3.1 Introduction . . . . .	8
3.2 Viola-Jones Object Detection . . . . .	8
3.3 Image Preprocessing . . . . .	10
3.3.1 Resizing . . . . .	10
3.3.2 Gray Scaling . . . . .	10
3.4 Histogram of Oriented Gradients . . . . .	10
3.4.1 Input Image . . . . .	11
3.4.2 Normalize Gamma & Colour . . . . .	11
3.4.3 Gradients . . . . .	11
3.4.3.1 Example: . . . . .	12
3.4.4 Weighted Vote in Spatial & Oriented Cells . . . . .	13
3.4.5 Contrast Normalize over Overlapping spatial cells . . . . .	14
3.4.6 Collect HOGs over Detection Window . . . . .	14

4.	Implementation . . . . .	15
4.1	Introduction . . . . .	15
4.2	High-Level View of the System . . . . .	15
4.3	Low-Level View of the System . . . . .	16
4.3.1	Low-Level View of Image Processing . . . . .	17
4.3.2	Low-Level View of SVM Model Testing & Training . . . . .	17
4.3.3	Low-Level View of Final System . . . . .	22
4.3.4	Optimizing HOG features . . . . .	22
4.4	Code Documentation . . . . .	23
4.5	Conclusion . . . . .	23
	Bibliography . . . . .	24



# List of Figures

2.1	The four stages for facial expression recognition . . . . .	3
2.2	Histogram Equalization . . . . .	4
2.3	Preprocessed Image with Oval Masks . . . . .	4
2.4	Architecture of an artificial neuron and a multilayered neural network . . . . .	6
3.1	Integral Image[6] . . . . .	9
3.2	Haar-like features[6] . . . . .	10
3.3	HOG feature extraction chain[10] . . . . .	11
3.4	Region of interest . . . . .	11
3.5	Image window . . . . .	11
3.6	One dimensional masks,(a)X-direction and (b) Y-direction . . . .	12
3.7	Image window . . . . .	12
3.8	The result of a one dimensional mask applied in the X-direction	12
3.9	The result of a one dimensional mask applied in the Y-direction	12
3.10	$16 \times 16$ pixel image . . . . .	13
3.11	Image divided into $4 \times 4$ pixel cells . . . . .	13
3.12	$\theta$ as an angle . . . . .	13
3.13	Unsigned gradients with 9 orientation bins . . . . .	13
3.14	Blocks of $B_x \times B_y$ cells . . . . .	14
3.15	Block A & B with a 50% overlap . . . . .	14
3.16	Cell histograms for contrast normalization in a block . . . . .	14
4.1	High-Level View of System . . . . .	15
4.2	Low-Level View of Image Processing . . . . .	17
4.3	Low-Level View of SVM Model Testing & Training . . . . .	17
4.4	Key Functions of a SVM [12] . . . . .	19
4.5	Confusion Matrix . . . . .	21

4.6	Confusion Matrix of SVM Model Classification . . . . .	21
4.7	Low-Level View of Final System . . . . .	22



# Chapter 1

## Introduction

Communication plays a large role, in daily human interaction. It can take the form of verbal or non-verbal communication, Mehrabian found that over 93% of verbal communication is conveyed through ones tone of voice, 38%, and non-verbal cues, 55% [1]. Understanding non-verbal communication is a valuable skill, in that it is a universal form of communication. Non-verbal communication is a combination of body language, physical gestures and facial expressions. Ekman & Friesen found six facial expressions that are universally identifiable in recognizing Fear, Anger, Disgust, Surprise, Happiness and Sadness[2]. A facial expression is made up of the changes in facial muscles mainly the mouth, eyes and eyebrows. These changes help to reflect ones current state of mind.

The rest of this chapter is organised as follows: Section 1.1 describes the problem statement; Section 1.2 provides the overview of the solution proposed in this project and Section 1.3 outlines the method of implementing the proposed solution.

### 1.1 Problem Statement

Companies use feedback from their customers as a metric to measure their customer satisfaction rate. Customer feedback can be initiated by the customer as a compliment or criticism, based on the service they were given. Alternatively, a company can offer their customers optional online or physical surveys to be completed. The problem is that customers are more likely to refrain from commenting on a service, unless provoked to do so. Customer service is a large revenue stream for companies whose core business is based around the customer. Ensuring that the customer stays happy is key, for their success.

## 1.2 Proposed Solution

An Automatic Human Emotion Detection(AHED) system can be applied in any environment that benefits from understanding facial expressions and human emotion. The proposed solution combines customer satisfaction with an automated system. This is done by using face detection to find the customers face in an image and facial feature extraction to identify the dominant customer emotion features in that image. To get the best results possible the training and classification of the system will be done using a machine learning technique. Companies can later incorporate the results of the system in improving their customer service, ensuring that their customers stay satisfied.

## 1.3 Proposed Method

A grayscale frontal image of the customer is used as input for the system. The Viola Jones Algorithm is then used to detect the location of the face in the image and the Histogram of Oriented Gradients(HOG)is used to extract the features. Lastly the AHED system is trained using Support Vector Machines(SVM), to classify each emotion.

# Chapter 2

## Related Work

There are a wide variety of methods that are used in the field of emotion recognition using facial expressions. However, a large number of those methods are implemented using a similar process. Shown in Figure 2.1, initially the image is captured and processed. Thereafter the face is detected and features are extracted. Then a machine learning technique is trained to do the classification of the emotions.

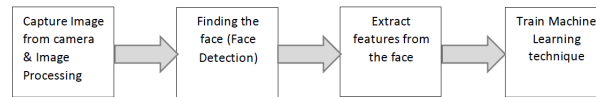


Figure 2.1: The four stages for facial expression recognition

The related work looks at all four stages in the implementation process and different methods used by other researchers. The rest of this chapter is organised as follows: Section 2.1 describes the process of obtaining the image from the camera and image processing; Section 2.2 explains how face detection functions; Section 2.3 explains the feature extraction process; Section 2.4 the training of the machine learning technique and Section 2.5 outlines the results achieved by other researchers.

### 2.1 Capture Image from the Camera and Image Processing

An image is taken from the camera and image processing tools are used to help normalize and standardize the input image.

- Goyal and Mittal, achieved the desired resolution and colour for their images by adjusting the brightness and contrast of the image[3].

- Reddy and Srinivas, scaled and cropped their images to  $(150 \times 120)$ , and ensured that the location of the eyes stayed the same in each image. The image was processed further, using an average combination of all the input image histograms. This process is called histogram equalization, see Figure 2.2, and helps in decreasing variation in an image. Histogram equalization is a technique for stretching out the intensity range of an image to enhance the contrast of the image[4].

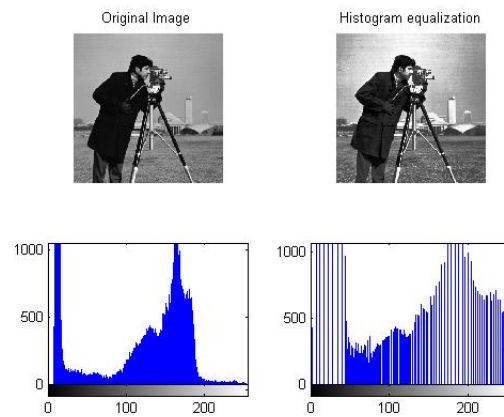


Figure 2.2: Histogram Equalization

- Boubenna and Lee, scaled their images to  $(100 \times 100)$  pixels[5].

## 2.2 Face Detection - Finding the Face

Finding the location of the face helps identify the region that contains all the features required to continue with facial expression recognition. The rest of the image is not important for this purpose.

- Reddy and Srinivas applied a fixed oval shaped mask over the image to extract the face region[4]. The images they used only contained faces, making it easier to apply the masks.

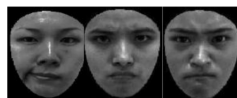


Figure 2.3: Preprocessed Image with Oval Masks

- Boubenna and Lee used the Viola and Jones algorithm to detect the location of the face in an image[6]. This algorithm uses Haar-like features

to help find the facial features, such as the eyes, nose and mouth. The Ada Boost algorithm is used to reduce the number of features, if there are too many. They used the canny edge detection operator to detect the edges of the face[5].

## 2.3 Feature Extraction - Extract the Features from the Face

Once the face has been detected, it is important to identify which features will be used for feature extraction. Either the full-face, or individual features from the face can be used as part of the feature set. These individual features can be the eyes, nose, mouth and eyebrows. The feature extraction algorithm can be applied based on its compatibility with the features chosen.

- Goyal and Mittal extracted the nose, mouth and eyes using the Viola and Jones Haar classifier[3].
- Reddy and Srinivas considered the entire face for the feature extraction not just the eyes, mouth and nose individually . First, they used Gabor filters to generate a bank of filters at 5 spatially varying frequencies and 8 orientations. The filtered outputs were then concatenated. Principle Component Analysis(PCA) was used to reduce dimensionality. PCA is a statistical technique that reduces the dimensions of feature vectors. The high dimensionality of feature vectors can cause over-fitting during classification. The PCA algorithm generates the eigenfaces for each image of dimension( $N \times N$ ). From this their system generated the eigenvector of dimension ( $N^2$ ) for each image. The vectors that relay the distribution of the face images the best are selected. These vectors are used to define the subspace(face space) of the face images[4]. The face image subspace represents a lower dimensional space( $N^2$ ) of the original image with dimentions( $N \times N$ ).
- Boubenna and Lee used Pyramid Histogram of Oriented Gradients(PHOG) to extract features. PHOG represents an image by its local shape and the spatial layout. The local shape of an image is represented by a histogram of edge orientations within an image sub-region, which are



divided into  $K$  bins. The spatial layout is represented by tiling the image into regions at different levels. Each image is divided into a sequence of increasingly finer spatial grids by repeatedly doubling the number of divisions in each axis direction[7]. The parameters of PHOG were set as follows: 3 for number of levels, 360 degrees for the number of dimensions and 16 for the number of bins. To decrease the number of features, a Genetic Algorithm(GA) was used, which resembles natural selection to find optimal features[5].

## 2.4 Training the machine learning technique

The training of the machine learning technique based on supervised learning whereby the machine learning technique is given labelled images (Happy, Sad, Anger, Disgust, Surprise, Fear and Neutral) and is required to learn them. Once the machine learning technique has completed its training, it can then be fed unlabelled images, and the result would be a prediction of which label best suits the given image.

- Goyal and Mittal used an Artificial Neural Network, with one hidden layer. The neural network architecture has three layers: input, hidden and output layers. Figure 2.4 provides a visual layout of the architecture of an individual neuron and a ANN with multiple layers. Feed-forward ANNs allows the signal to travel in one direction from the input layer to the output layer. Recurrent networks contain feedback connections, where the signal moves in both directions. To get accurate results from the ANN, the weights can be set explicitly using prior knowledge, or the ANN can be trained to help find the optimal weights[3, 8].

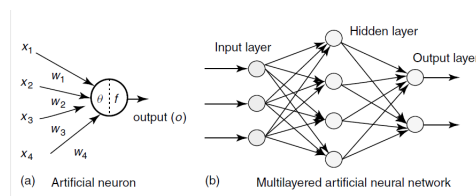


Figure 2.4: Architecture of an artificial neuron and a multilayered neural network

- Reddy and Srinivas, used an Artificial Neural Network, with two hidden layers[4].
- Boubenna and Lee, used Linear Discriminants Analysis(LDA) and K Nearest Neighbours(KNN). LDA finds the maximum distances within classes, to obtain maximum class separation. LDA only uses up to second order moments, such as the covariance and mean, of the class distribution. KNN classifies unlabelled samples according to the training samples. KNN finds the nearest K in the labelled samples and set them to the closest group, for the unlabelled samples. One distance measure is required, and [5] used the cosine distance measure.

## 2.5 Results

The results from the three studies are as follows, with [5] having the best overall results for their facial expression recognition system.

- Goyal and Mittal achieved an 80% classification accuracy, using a confusion matrix and a regression plot to verify the results[3].
- Reddy and Srinivas achieved an 85.7% classification accuracy using the JAFFE database[4].
- Boubenna and Lee achieved a 99.33% accuracy, using the Radboud Faces Database(RaFD)[5].

# Chapter 3

## Image Processing Techniques

### 3.1 Introduction

This chapter looks at image processing techniques used in obtaining the features needed to do the final classification. The Viola and Jones algorithm, developed by Viola and Jones, is used to detect the location of the frontal face in the image. Once we have this location we then extract the face, which represents our region of interest. The region of interest is then Gray-scaled and is now ready for feature extraction. The Histogram of Oriented Gradients(HOG) is used for the feature extraction process.

The rest of this chapter is organised as follows: Section 3.2 provides details on Viola-Jones Object Detection and it's key concepts; Section 3.3 covers the image pre-processing techniques used and Section 3.4 explains how Histogram of Oriented Gradients are used for feature extraction.

### 3.2 Viola-Jones Object Detection

The Viola-Jones algorithm [6] is an object detection method that uses Haar-like features. For this project the Viola-Jones object detection is used to find the location of frontal faces in images. The algorithm uses three concepts to effectively detect objects with certain features:

The first concept is the integral image, which allows for the features in the image to be evaluated much faster. This is also known as an intermediate view of the image. At each point( $x, y$ ) in the integral image, there is the sum of the pixels above and to the left of the point( $x, y$ ), inclusive.

Referring to Figure 3.1: In order to calculate the sum of the pixels within the rectangle D, only four original image references are needed.

- At point 1 in the integral image, the sum of all the pixels in rectangle A in the original image are used.

- At point 2 in the integral image the sum of all the pixels in rectangles A and B in the original image are added together ( $A+B$ ).
- At point 3 in the integral image the sum of all the pixels in rectangles A and C in the original image are added together ( $A+C$ ).
- Lastly, at point 4 in the integral image of all the rectangles are added together ( $A+B+C+D$ ) in the original image.

Thus, to get the sum of the pixels in rectangle D will result in  $4+1-(2+3)$ .

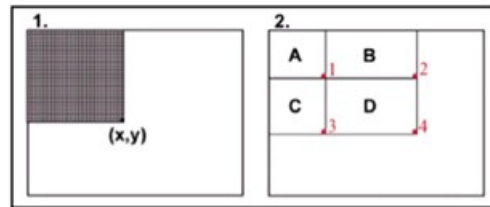


Figure 3.1: Integral Image[6]

The second concept in the Viola-Jones framework[6] is a classifier based on reducing a large feature set down to a smaller set of important features. This is done by using Ada-Boost. Ada-Boost finds a weak classifier and forces it to depend on a single feature, resulting in a stronger classifier. A weak classifier is selected at each stage of the boosting process, or feature selection process. The third concept is a method that combines weak classifiers in a rejection cascade. This increases the speed of detection as the focus is now on promising areas of the image. Each stage in the cascade is formed using Ada-Boost. The Viola-Jones algorithm uses many Haar-like features. I will describe Three of these Haar-like features.

- Two-rectangle feature - is calculated by subtracting the sum of the pixels of one rectangle from the sum of the other. The rectangles need to be the same size, shape and need to be vertically or horizontally adjacent.
- Three-rectangle feature - is calculated by subtracting the sum of the two outside rectangles from the middle rectangle.

- Four-rectangle feature - is calculated by subtracting the sum of the pixels of one diagonal pair from the other.

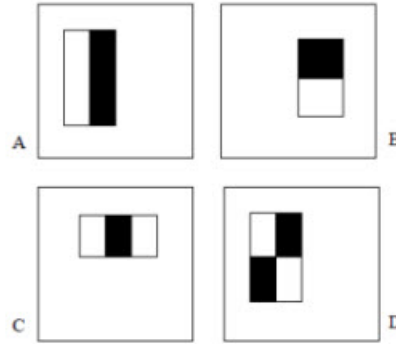


Figure 3.2: Haar-like features[6]

### 3.3 Image Preprocessing

#### 3.3.1 Resizing

The image is resized once we obtain our region of interest using Viola-Jones to detect the location of the face. The main benefit of resizing the image is to maintain uniformity in our feature set. Another benefit is that it scales down the number of pixels in a image, resulting in a smaller feature set[9].

#### 3.3.2 Gray Scaling

Converting an image from RGB(colour) to Grayscale helps in reducing the number of colour channels down to a single color channel. A commonly used method is the standard NTSC conversion formula, that calculates the luminance of a pixel[9]:

$$\text{Luminance of a pixel} = (0.2989 \times \text{red}) + (0.5870 \times \text{green}) + (0.1140 \times \text{blue})$$

### 3.4 Histogram of Oriented Gradients

The Histogram of Oriented Gradients is a feature extraction method for images. Where a image is divided into cells, of  $C_x \times C_y$  pixels, that form a grid over the image. Histograms are calculated for each of these cells based on the orientation of the gradients of the pixels in the cell. This is followed by the image further being divided into a grid, of  $B_x \times B_y$  cells, which are called

blocks. Each block is used to contrast-normalize the histograms(cells) present in the block. The dimensions of the final feature vector calculated by:

$$\text{total number of blocks} \times \text{number of cells in each block} \times \text{number of orientation bins}$$

Further details of the Histogram of Oriented Gradients will be discussed following Dalal & Triggs HOG feature extraction chain, see Figure 3.3, excluding the linear SVM and classification[10].

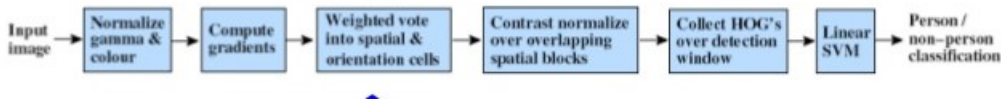


Figure 3.3: HOG feature extraction chain[10]

### 3.4.1 Input Image

Given an image, first identify the region of interest in the image. This region then forms your image window.



Figure 3.4: Region of interest



Figure 3.5: Image window

### 3.4.2 Normalize Gamma & Colour

Normalizing the image window is an optional addition to the HOG. Dalal & Triggs found that normalizing the image pixels( $p$ ) at this stage did not have a noticeable impact on the performance at the detection stage of their research. However when choosing to normalize the image window Gamma (power law) normalization had a negative impact on the results, while Square-root normalization had more of a positive impact on the results.

Gamma Normalization:  $\log(p)$

Square-root Normalization:  $\sqrt{p}$

### 3.4.3 Gradients

The gradient for the image window is computed by applying a one dimensional mask in both X ( $G_x$ ) and Y ( $G_y$ ) directions.

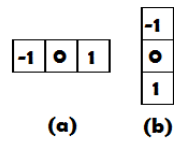


Figure 3.6: One dimensional masks, (a) X-direction and (b) Y-direction

The mask convolves over the image window. At each point where the mask is placed the pixels are multiplied by the mask. After that the two outer pixels are added together and the result is placed in the position of the center pixel. The mask is not able to compute the gradients on the pixels around the edge of the image. Unless extra pixels are added to the edges of the image before hand. The example below shows how the loss in pixels affects the resulting gradient image.

### 3.4.3.1 Example:

22	24	18	11	23
22	24	18	11	23
22	99	18	11	23
22	24	18	11	23
22	24	18	11	23

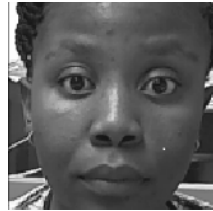


Figure 3.7: Image window

	-4	-13	5	
	-4	-13	5	
	-4	-88	5	
	-4	-13	5	
	-4	-13	5	



Figure 3.8: The result of a one dimensional mask applied in the X-direction

0	75	0	0	0
0	0	0	0	0
0	-75	0	0	0

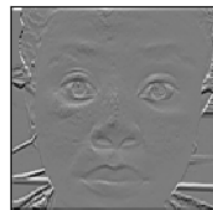


Figure 3.9: The result of a one dimensional mask applied in the Y-direction

Now that we have the gradients, we can compute the magnitude and orientation of the gradients from  $G_x$  &  $G_y$ .

$$\text{Magnitude: } G = \sqrt{G_x^2 + G_y^2}$$

$$\text{Orientation: } \theta = \arctan\left(\frac{G_y}{G_x}\right)$$

### 3.4.4 Weighted Vote in Spatial & Oriented Cells

We can now decide on the dimensions of each cell, before calculating the HOGs. In their research Dalal & Triggs found that the size of the cell are dependent on the size of the features you need to extract (e.g. eyes, nose, mouth).

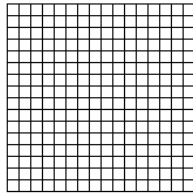


Figure 3.10:  $16 \times 16$  pixel image

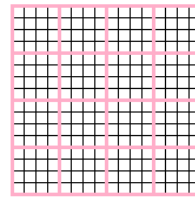


Figure 3.11: Image divided into  $4 \times 4$  pixel cells

The next parameter is the number of orientation bins. The orientation of the gradient can be described as the angle of the gradient. There are two options available when choosing the range of the gradient angle:

- Signed  $[0, 360]$  degrees
- Unsigned  $[0, 180]$  degrees

Unsigned gradients in the range  $[0, 180]$  degrees, with the number of orientation bins in the range  $[9, 12]$  are the preferred values for the orientation bins.



Figure 3.12:  $\theta$  as an angle

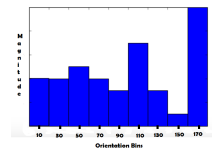


Figure 3.13: Unsigned gradients with 9 orientation bins

Looking at a single cell. Each pixel of the gradient magnitude image contributes to a orientation bin of a cells histogram. The value of the same pixel in the gradient orientation image helps you identify which orientation bin to place the gradient magnitude of the pixel.



### 3.4.5 Contrast Normalize over Overlapping spatial cells

Contrast normalization is used to ensure that the cells are not affected vastly by changes illumination and contrast in the image. Starting with dividing the image into blocks that can fit at least 2-3 features, these blocks are allowed to overlap one another for more detailed feature set. Contrast normalization works by taking the sum of the histograms in a block  $S_b$  and dividing each of the histograms  $H_{hist}$  by  $\sqrt{S_b^2 + \epsilon^2}$ . The result is a normalized histogram  $H_{norm}$  in each cell.

$$\text{Contrast normalization : } H_{norm} = \frac{H_{hist}}{\sqrt{S_b^2 + \epsilon^2}}$$

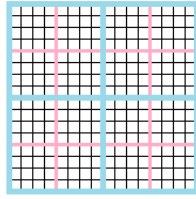


Figure 3.14: Blocks of  $B_x \times B_y$  cells

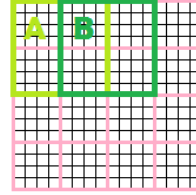


Figure 3.15: Block A & B with a 50% overlap

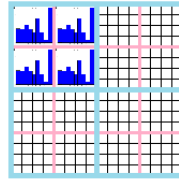


Figure 3.16: Cell histograms for contrast normalization in a block

### 3.4.6 Collect HOGs over Detection Window

The final step is to concatenate all the normalized histograms to form a one dimensional feature vector  $[H_{norm}, H_{norm}, H_{norm} \dots]$ . This feature vector is then used for the classification and training of the system.

# Chapter 4

## Implementation

### 4.1 Introduction

This chapter looks at the high-level and low-level views of the system and code documentation. The high-level view in Section 4.2 provides an outline of the processes followed during the implementation of the system, while the low-level view in Section 4.3 provides a more detailed description of the implementation of the system. The physical code documentation is not provide in this documentation, however a link to the GitHub repository is provided in Section 4.4.

### 4.2 High-Level View of the System

The high-level view of the system provides an overview of all the stages that the sytem follows when classifying an image given as input. These stages include: Capture Frame, Face Detection, Feature Extraction, Train Machine Learning Technique and Emotion Classification. Figure 4.1 serves as a visual aid for the content that follows.

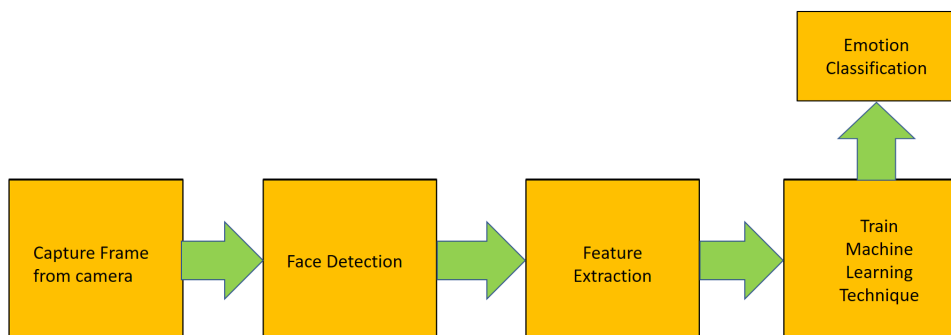


Figure 4.1: High-Level View of System

Looking at Figure 4.1, the High-level view is explained as:

- **Capture Frame** - The web camera records a constant stream of video input. The video input consists of a sequence of multiple image frames. The system captures each frame for processing as it is displayed on the video feed.
- **Face Detection** - Now that we have captured a single frame, we need to check if there is a face present in the frame. This is done using a face detection algorithm. If a face is present in the frame, the location of the face is extracted. The rest of the image is disregarded at this point.
- **Feature Extraction** - Every emotion displayed facially has it's own set of unique identifying features. By applying feature extraction we are able to represent these features in a way that a computer can understand and process. The feature extraction method is applied to the region of the image that contains the face.
- **Train Machine Learning Technique** - Machine learning is a method used by computers to learn how to identify patterns in a given set of features. This process is called training. When we train the system, our features are labelled (Happy, Sad, Angry etc.). Labelling the features helps guide the computer in the learning process.
- **Emotion Classification** - When the training is complete, classification helps to test the accuracy of the trained model. At this point the model should be able to identify emotions given unlabelled features.

### 4.3 Low-Level View of the System

The high-level view of the system dives deeper into the details of the components used to implement the system. This is done following the same stages used in the high-level view of the system. In this section we will look at three conceptual low-level views that relate to our system. These views are aligned to the processes followed in image processing, Support Vector Machine(SVM) model testing & training and implementing the final system.

### 4.3.1 Low-Level View of Image Processing

The first low-level view is a visual representation of how the high-level view relates to the image processing techniques discussed in Chapter 3 is presented in Figure 4.2.



Figure 4.2: Low-Level View of Image Processing

### 4.3.2 Low-Level View of SVM Model Testing & Training

The second low-level view relates to the training and testing of the SVM model used to classify the emotions. Where the 'Capture Frame' stage is replaced by 'Get Images from Dataset'.

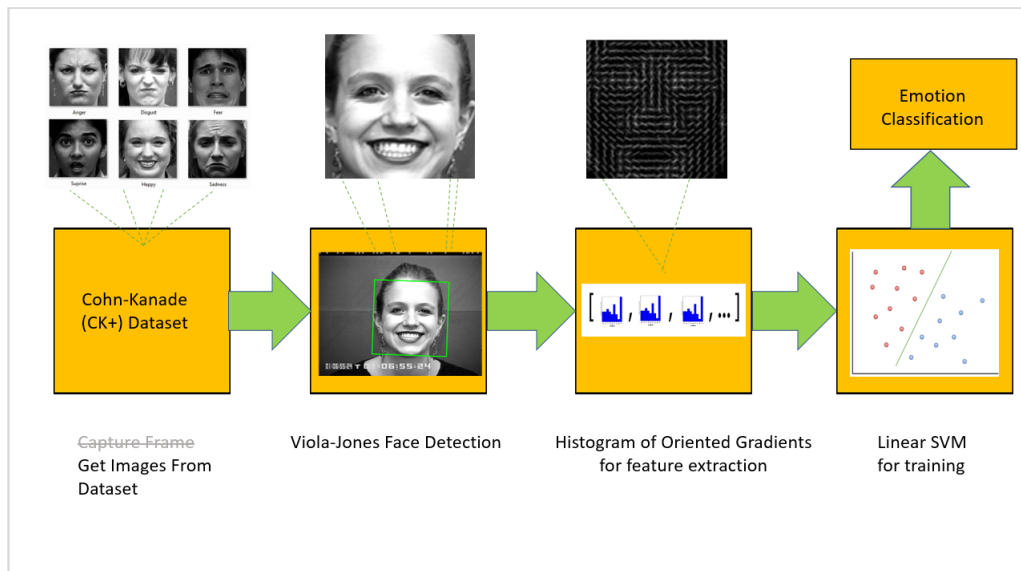


Figure 4.3: Low-Level View of SVM Model Testing & Training

Looking at Figure 4.3, the Low-level view for SVM model testing and training is explained as:

- **Get Images From Dataset** - The Extended Cohn-Kanade Dataset[11] is used as data for the training and testing of our SVM model. The images are divided into seven groups, whose labels are: Angry, Disgust, Fear, Happy, Sadness, Surprise and Neutral. The Table 4.1 provides the total number of images present for each label.

Emotion Label	N
Angry	45
Disgust	59
Fear	25
Happy	69
Sadness	28
Surprise	83
Neutral	35

Table 4.1: CK+ Image Labels and Totals

- **Face Detection** - The Viola-Jones face Detection is used to extract the face from each image in the CK+ dataset. The face region of the image is stored as a  $56 \times 56$  pixel grayscale image.
- **Feature Extraction** - The resulting images from the face detection are used as inputs for the Histogram of Oriented Gradients(HOG) which gives a one-dimensional feature vector for each image. Each image is given a numerical label from 1 to 7 based on the emotion displayed in the image, see Table 4.2 for the labels. The feature vector is stored in the feature dataset with the corresponding numerical emotion label(e.g. [Numerical Label][Feature Vector]).

Emotion Label	Angry	Disgust	Fear	Happy	Neutral	Sadness	Surprise
Numerical Label	1	2	3	4	5	6	7

Table 4.2: Emotions with corresponding Labels for the feature vectors

The parameters used for the HOG are listed in Table 4.3.

Parameter	Size	Type
Image	(56, 56)	Pixels
Cell	(4, 4)	Pixels
Block	(3, 3)	Cells
Overlap	66.66%	Blocks
Bins	9	(0 – 180 <i>degrees</i> ) Unsigned Gradients
Feature Vector Size: 11664		

Table 4.3: HOG Parameters

- Train Machine Learning Technique** - The machine learning technique used to do the classification for our system is Support Vector Machines(SVMs). [11]Used SVMs to test the accuracy of their CK+ dataset due to its proven accuracy with face and facial action detection. An SVM attempts to separate the closest negative and positive points in each class from each other. Once this separation is achieved it makes it easier to separate negative and positive points that are further away from each other, as the similarities in these points are less than those in the points that are closer. The distance between these points is calculated by subtracting the positive and negative points from each other. The positive and negative points closest to each other are considered our support vectors, these help in determining the separating hyperplane. The SVM needs to ensure that the distance between the support vectors is maximized. Figure 4.4 visualizes the concepts discussed.

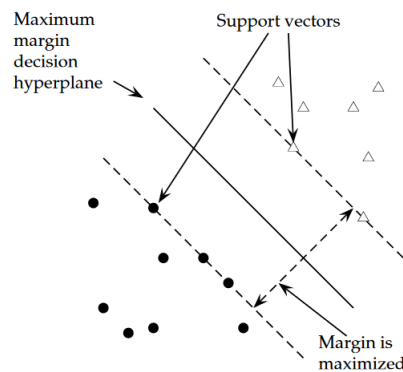


Figure 4.4: Key Functions of a SVM [12]

**Cross-Validation-** The feature dataset compiled after completing the

feature extraction is used as input for the SVM. Where the dataset is divided into a training and testing dataset. The percentage of the split is dependent on how well the model performs with each split. A 60% training and 40% testing split was used for the SVM. Both these datasets should contain an equal distribution of the emotion labels in each dataset. This step ensures that all labels appear equally in both datasets. The Cross-validation score using a stratified K-fold of 3 is approximately 84%. [13]Cross-validation measures the overall performance of the SVM model on different training and testing data splits. This tests the independence of the model to the dataset and helps to prevent overfitting in our model. K-fold cross-validation is done by splitting the dataset into K subsets of equal length. Each subset is then tested on the SVM model of the remaining K-1 subsets. Stratified K-fold cross-validation ensures that the classes are distributed equally in each subset. The training dataset is used as input for training our SVM and creating the SVM model that will be used in our system.

**Grid-Search-** After optimizing the SVM model with grid-search a Linear Kernel with a C of 1 was used. [13]Grid search helps to find the optimal parameters for the SVM model. Where C determines the extent to which the SVM model should avoid misclassification in the training. Larger values of C decrease the margin of the hyperplane which aims to increase the accuracy of the training. A smaller value for C increases the margin of the hyperplane, but carries the downside of more misclassified training data points. In [13]the optimal range for  $C = 2^{-5}, 2^{-3}, \dots, 2^{15}$ .

- **Emotion Classification** - The testing dataset is used to assess the performance of the the trained SVM model on unseen data. Where the SVM model is given the testing dataset features without the corresponding labels. The results of the SVM model classification are then compared to the original labels to test the accuracy of the SVM model. The SVM model trained for our system has an overall accuracy score of 88.2%.

**Confusion Matrix-**See Figure 4.5,a confusion matrix summarises the

outcomes of the classification based on the the actual labels of the testing data and those obtain from testing. The values obtained from the confusion matrix help with analyzing the SVM mode. Figure 4.6 shows the confusion matrix for our SVM model. The diagonal starting at the top-left index till the bottom-right index contains all the true positives/negatives for our SVM model.

		Actual Result/ Classification	
		Positive	Negative
Predicted Result/ Classification	Positive	TP (true positive)	FP (false positive) Type I Error
	Negative	FN (false negative) Type II Error	TN (true negative)

Figure 4.5: Confusion Matrix

	Angry	Disgust	Fear	Happy	Neutral	Sad	Surprise
Angry	14	0	0	0	3	0	0
Disgust	1	22	0	1	0	0	0
Fear	1	0	6	0	1	2	0
Happy	0	0	0	27	0	0	0
Neutral	2	0	1	1	10	0	0
Sad	1	0	0	0	1	9	0
Surprise	0	0	0	0	0	1	32

Figure 4.6: Confusion Matrix of SVM Model Classification

Table 4.4 gives an overview of the formulas and terms used in Figure 4.5 and Table 4.6. Where Table 4.6 contains the 'SVM Model Classification Report'. The report indicates the performance of each individual class and the overall estimated performance of the SVM model with regards to the precision, recall and f1-score. Classes that had more data available performed better overall as compared to those that had less. There was more testing data available for these classes. Working with an uneven dataset makes it harder to judge the performance of each class in comparison to the other classes.

SVM Model Evaluation		
Term	Formula	Description
Type I Error	FP	False Positive
Type II Error	FN	False Negative
Accuracy	$\frac{TP+TN}{TP+TN+FP+FN}$	Evaluates the degree of correctness for the predictions
Precision	$\frac{TP}{TP+FP}$	The Positive predictive value
Recall	$\frac{TP}{TP+FN}$	True positive rate
F1-Score	$2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$	Evaluates the accuracy of predictions

Table 4.4: Terminology and formulas used for evaluating the SVM Model[14]



SVM Model Classification Report				
Label	precision	recall	f1-score	support
Angry	0.74	0.82	0.78	17
Disgust	1.00	0.92	0.96	24
Fear	0.86	0.60	0.71	10
Happy	0.93	1.00	0.96	27
Neutral	0.67	0.71	0.69	14
Sad	0.75	0.82	0.78	11
Surprise	1.00	0.97	0.98	33
avg total	0.89	0.88	0.88	136

Table 4.5: SVM Classification Report

### 4.3.3 Low-Level View of Final System

The third low-level view shows how the process of Automatic Human Emotion Detection is streamlined for user interaction. Where classifications are performed live as the user changes their facial expressions. At this point we use the 'Trained SVM Model' with the HOG parameters from Table 4.3 above.

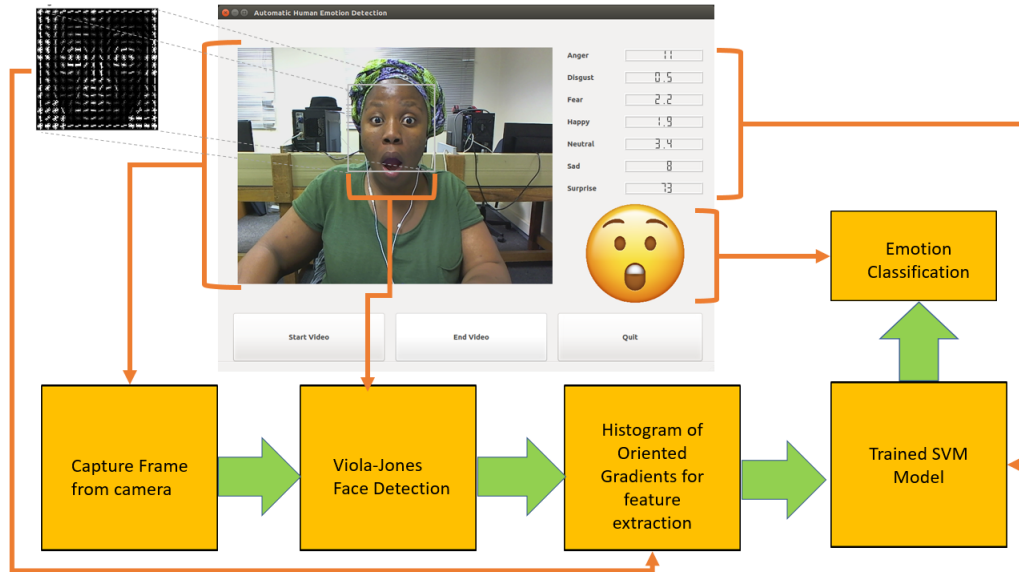


Figure 4.7: Low-Level View of Final System

### 4.3.4 Optimizing HOG features

This section covers the different combinations for the HOG parameters that were considered before the values in Table 4.3 were chosen. The HOG implemented from scratch by us is compared to the OpenCV HOG using the Python Sklearn SVM library. All the other parameters remained the same at

this point, the dataset had a 60% training and 40% testing dataset. The SVM model optimized to a Linear Kernel with a C of 1 for each iteration of the HOG features optimization. Both HOGs maintained a bin size of 9 and 50% overlap for Block(2, 2), with a 66.66% overlap for Block(3, 3).

<b>Accuracy of HOG Parameters</b>		
-	<b>Cell(8, 8) Block(3, 3)</b>	<b>Cell(4, 4) Block(3, 3)</b>
OpenCV HOG	86%	88.2%
AHED HOG	77.9%	75.7%
-	<b>Cell(8, 8) Block(2, 2)</b>	<b>Cell(4, 4) Block(2, 2)</b>
OpenCV HOG	83.8%	83.8%
AHED HOG	77.9%	82.3%

Table 4.6: HOG Optimization

## 4.4 Code Documentation

All the code used for the implementation of the Automatic Human Emotion Detection system can be found at :<https://github.com/tsweedie/Final>.

## 4.5 Conclusion

The Automatic Human Emotion Detection system was implemented entirely using Python. OpenCV is used for the Viola-Jones face detection, Sklearn was used to implement the Histograms of Oriented gradients and the Support Vector Machines. The HOG implementation done from scratch used python numpy arrays and followed the implementation method discussed in Chapter 3. The highest accuracy achieved for the HOG implemented with OpenCV was 88.2% and 82.3% with the HOG implemented in this project.

# Bibliography

- [1] A. Mehrabian. *Nonverbal communication*. CURRENT CONTENT, CA, USA, 1984.
- [2] P. Ekman. *NONNERBAL BEHAVIOR AND COMMUNICATION*. Lawrence Erlbaum Association, New Jersey, 1977.
- [3] R. Goyal and T. Mittal. Facial expression recognition using artificial neural network. *HCTL Open Int. J. of Technology Innovations and Research*, 10:1–10, July 2014.
- [4] Ch. Satyananda Reddy and T. Srinivas. Improving the classification accuracy of emotion recognition using facial expressions. *International Journal of Applied Engineering Research*, 11(1):650–655, 2016.
- [5] H. Boubenna and D. Lee. Feature selection for facial emotion recognition based on genetic algorithm. *12th International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery*, 12:511–517, August 2016.
- [6] P. Viola and M. J. Jones. Robust real-time face detection. *International Journal of Computer Vision*, 2(57):137–154, 2004.
- [7] A. Bosch, A. Zisserman, and X. Munoz. Representing shape with a spatial pyramid kernel. *CVIR*, pages 401–408, July 2007.
- [8] A. Abraham. *Handbook of Measuring System Design*. John Wiley & Sons, Ltd., OK, USA, 2005.
- [9] K. Bose and S. Bandyopadhyay. Crack detection and classification in concrete structure. *Journal for Research*, 2(4):29–38, June 2016.
- [10] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference*, 1:886–893, 2005.

- [11] Patrick Lucey, Jeffrey F. Cohn, Takeo Kanade, Jason Saragih, Zara Ambadar, and Iain Matthews. The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression.
- [12] P. Raghavan D. Manning and H. Schtze. *Introduction to Information Retrieval*. Cambridge University Press, 2009.
- [13] Chih-Wei Hsu, Chih-Chung Chang, Chih-Jen Lin, et al. A practical guide to support vector classification. 2003.
- [14] Claude Sammut and Geoffrey I Webb. *Encyclopedia of machine learning*. Springer Science & Business Media, 2011.