



ARTIFICIAL INTELLIGENCE BASED FACIAL EXPRESSION AND RECOGNITION

Final Project Report

Submitted by

Javed Ali

*in partial fulfillment for the award of the degree
of*

Bachelor of Technology

IN

Computer Science and Engineering

SCHOOL OF COMPUTING SCIENCE AND ENGINEERING

**Under the supervision of
Dr. Tapas Kumar**

MAY-2020

TABLE OF CONTENTS

CHAPTER NO.	TITLE	PAGE NO.
1.	Abstract	1
2.	Introduction	2
3.	Existing System	4
4.	Proposed System	8
5.	Implementation or architecture diagrams	14
6.	Conclusion/Future Enhancement	23
7.	References	24

Facial Expression Recognition is the part of **Facial recognition** which is an important system and the need for it increases tremendously. Facial Expression Recognition is the part of Facial recognition which is an important system and the need for it increases tremendously. Facial expression recognition enables us to understand human emotions and it is the basis as well as essential for quantitative analysis of emotions. Facial expression recognition is an emerging field. In this paper, the current research states are analyzed mostly from the latest facial expression extraction algorithm and the Facial Expression Recognition algorithm based on AI/ML, deep learning a comparison is made of these methods. The paper also focuses on a system of recognizing human emotion from a human's face. The analyzed information is conveyed to the computer by the regions of the eye and the mouth into a merged new image in various facial expressions. Age estimation and face recognition are the most robust techniques to maintain Authenticity. In today's world fraud and scam are on the rise and to curb all of this we have implemented this project. Our field of study is computer vision, **Computer vision** is concerned with the automatic extraction, analysis and understanding of useful information from a single image or a sequence of images. Computer vision uses techniques from machine learning and, in turn, some machine learning techniques are developed especially for computer vision.

2.1 Purpose

There are increasing incidents of fraud related to age and identity below are the types related to the same.

All of these crimes are related to wrong identification of age of people. In order to curb this issue our system will be efficient and play an important role in combating this increasing age fraud cases. It is important to validate the authenticity of a person regarding age as date of birth is required in almost every legal official record in government or non-government organization for issue of important documents like citizen card, financial claims etc. Some of the cases which took headlines due to age fraud:

Credit Card fraud: The crime of credit card fraud begins when someone either fake credit card details or fraudulently obtains the card number and other account information like age is necessary for the card to be used successfully.

Employment and Tax Related: criminal frauds in producing false documents with fake age of a professional to get inside margin line of employment misleading the organization and taking job over deserved candidates. This also make them eligible for rebate on taxation policy.

Phone and Utilities: fraudsters using false identification to use mobile service and utilities offered by service provider as cable accounts are actually the general type of utility fraud we see, followed by the opening of fraudulent household electricity and gas accounts.

Bank Fraud and grant of Loan: opening a bank account with fake documents hiding actual identification like name, age, income etc. to fulfil the eligibility is kind of forgery banks are dealing with nowadays. Loan is granted on the basis of credit score which comprises of personal details and by producing forged documents will cost banks a huge stack of money.

Government benefits: to enjoy or utilize the benefits offered by government then an individual need to match desired criteria or standards set by governing body but scams and fraud take places due to submitting fake proofs and certificates by citizens.

2.2 A recent incident-

According to the article Bangalore Mirror age fraud in Tennis came into limelight when parents of former India No. 1 among the Under-12 girls, Siddhi Khandelwal, deliberately changed her year of birth in official records to play an extra season in a lower age group in the Indian tennis circuit. That's what her date of birth certificate clearly indicates. Siddhi, the national U-12 champion this year, is over the age of 12.

According to the report of The Guardian Kwankwo Kanu's official age is 33 but his real age is 42. Obafemi Martins is not 25 but 32. Jay-Jay Okocha was 10 years older than his "official" age throughout his career. And Taribo West, whose playing career ended only two years ago, is in his

late fifties. Who says so? A stream of bloggers on some of Nigeria's most popular websites, in response to comments made after the country's timid effort in last month's Africa Cup of Nation

2.3 Motivation and Scope

This project has its applications in various fields like content filtering and marketing according to user by detecting age and suggesting various products by age. Some of its applications are as follows:

Age based Access Control: Preventing certain goods like alcohol drinks, cigars, cigarettes by under age in most cases of age-based restrictions control is enforced by judgement of human an alternative approach for this automatic facial estimation.

Human Machine Interaction: people belonging to different age group should interact with machines differently for eg. The way a child uses a computer is very different and the interface of the system should be user friendly and for old people the icons should be larger.

Data Mining: age-based retrieval of images from the internet for various purposes.

Unlock Phones: A variety of phones including the latest iPhone are now using face recognition to unlock phones. This technology is a powerful way to protect personal data and ensure that, if a phone is stolen, sensitive data remains inaccessible by the perpetrator.

Smarter Advertising: Face recognition has the ability to make advertising more targeted by making educated guesses at people's age and gender. Companies like Tesco are already planning on installing screens at gas stations with face recognition built in. It's only a matter of time before face-recognition becomes an omni-present advertising technology.

Find Missing Persons: Face recognition can be used to find missing children and victims of human trafficking. As long as missing individuals are added to a database, law enforcement can become alerted as soon as they are recognized by face recognition—be it an airport, retail store or other public space.

CASINOS: Face recognition can help casinos recognize the moment that a cheater or advantage gambler enters a casino. In addition, face recognition can recognize members of voluntary exclusion lists, who can cost casinos hefty fines if they're caught gambling.

FIND LOST PETS: Finding Rover is an app that tries to help owners reunite with lost pets. The app uses face recognition (albeit it's the face of an animal in this case) to match photos that pet owners upload to a database of photos of pets in shelters. The app can then instantly alert owners if their pets are found.

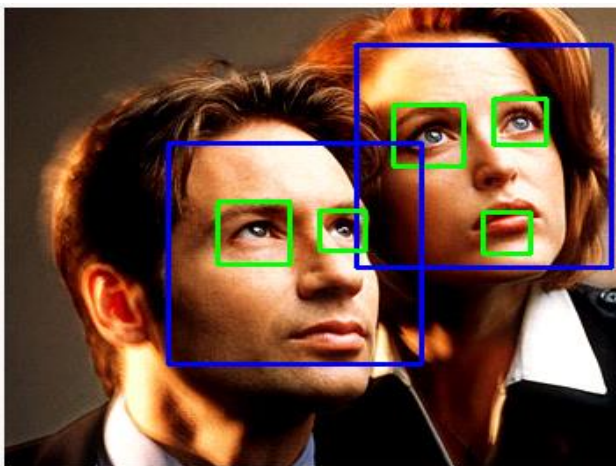
Healthcare: Machine learning is being combined with computer vision to more accurately track patient medication consumption and support pain management procedures.

Face Detection

Over the years there has been a lot of research in the field of Face detection and recognition, and various methodologies that can be used regarding it. We will list a majority of them and will explain about each one in detail. There are other algorithms also but as we have implemented the code in OpenCV and python so we listed only those which are present in OpenCV

- Haar Cascades
- Local Binary Pattern Histograms (LBP Cascades)
- Deep learning

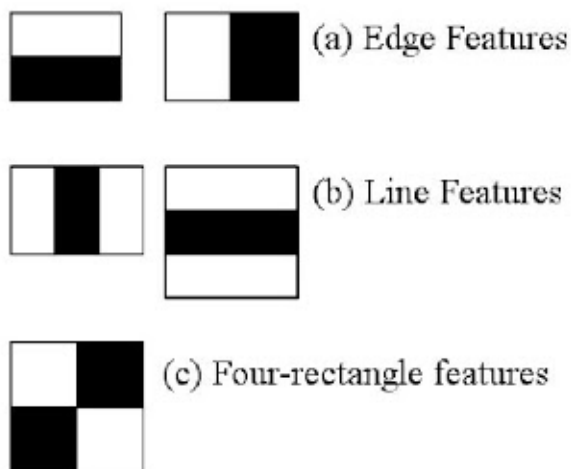
Haar Cascades: The first algorithm that was developed for face detection. It was researched by Paul Viola and Micheal Jones [6]. Haar features are digital image features used in object detection, the name is based on Haar wavelets. Viola and Jones used the idea of Haar wavelets and developed the so-called Haar-like features. A Haar-like feature considers adjacent rectangular regions at a specific location in a detection window, sums up the pixel intensities in each region and calculates the difference between these sums. This difference is then used to categorize subsections of an image. For example, let us say we have an image database with human faces. It is a common observation that among all faces the region of the eyes is darker than the region of the cheeks. Therefore, a common Haar feature for face detection is a set of two adjacent rectangles that lie above the eye and the cheek region. The position of these rectangles is defined relative to a detection window that acts like a bounding box to the target object (the face



in this case).

(Figure 3.1)

Initially, the algorithm needs a lot of positive images (images of faces) and negative images (images without faces) to train the classifier. Then we need to extract features from it. For this, Haar features shown in the below image are used. Each feature is a single value obtained by subtracting sum of pixels under the white rectangle from sum of pixels under the black rectangle.



(Figure 3.2)

Haar Cascade Detection in OpenCv:

OpenCV comes with a trainer as well as detector. If you want to train your own classifier for any object like car, planes etc. you can use OpenCV to create one. First, we need to load the required XML classifiers. Then load our input image (or video) in grayscale mode.

Now we find the faces in the image. If faces are found, it returns the positions of detected faces as Rect(x,y,w,h). Once we get these locations, we can create a ROI for the face and apply eye detection on this ROI , since eyes are always on the face .

LBP Cascade classifier

It is also an object detection in digital images, but the working of LBP is different that the Haar features. Each training image is divided into some blocks as shown in the picture below.

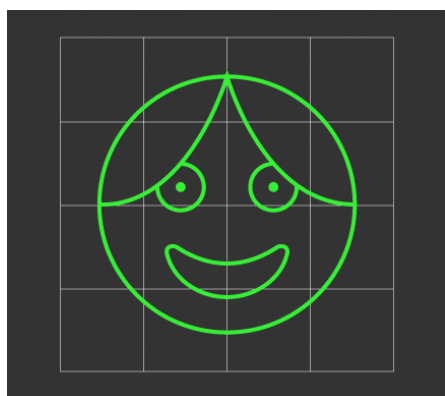


Figure (2.3)

For each block, LBP looks at 9 pixels (3×3 window) at a time, and with a particular interest in the pixel located in the center of the window. Then, it compares the central pixel value with every

neighbor's pixel value under the 3×3 window. For each neighbor pixel that is greater than or equal to the center pixel, it sets its value to 1, and for the others, it sets them to 0.

After that, it reads the updated pixel values (which can be either 0 or 1) in a clockwise order and forms a binary number. Next, it converts the binary number into a decimal number, and that decimal number is the new value of the center pixel. We do this for every pixel in a block.

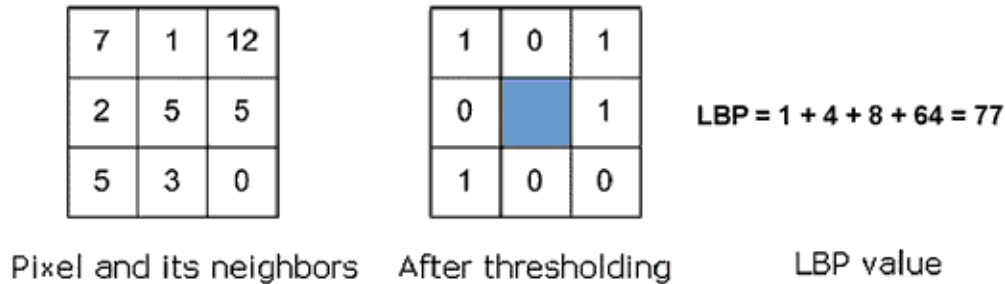
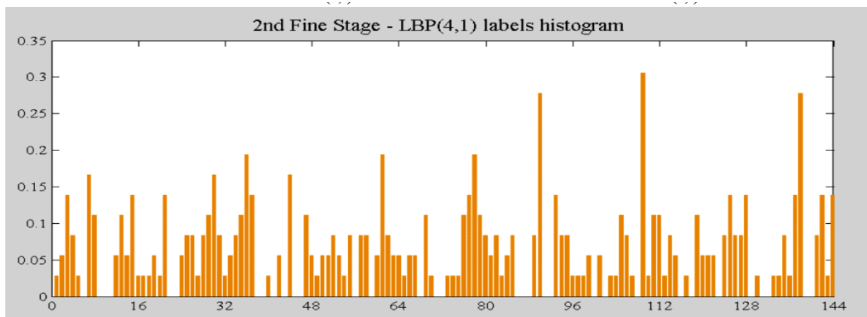


Figure (3.3)

It then converts each block values into a histogram, so now we have gotten one histogram for each block in an image, like this:



(Figure 3.4)

Finally, it then concatenates these block diagrams into one to form a feature vector for one image which contains all the features we are interested in.

Implementing LBP in OpenCv:

We just need to change a .xml classifier file, in the previous code as

```
#load cascade classifier training file for lbpcascade
lbp_face_cascade = cv2.CascadeClassifier('data/lbpcascade_frontalface.xml')
```

Deep learning for face detection

Deep learning module was incorporated later in the OpenCV 3.1, Before these the above two methods only could be used for the purpose. This module now supports a number of deep learning frameworks, including Caffe, TensorFlow, and Torch/PyTorch. With OpenCV 3.3, we can utilize pre-trained networks that support various popular deep learning frameworks. The fact that they

are pre-trained implies that we don't need to spend much time training the network, rather we can complete a forward pass and utilize the output to make a decision within our application.

Popular architectures that are supported by this module are:

- GoogleLeNet
- AlexNet
- SqueeZNet
- VGG
- ResNet

The Caffee module: Caffee is a deep learning framework made with expression, speed, and modularity in mind. It is developed by Berkeley AI Research (BAIR) and by community contributors. Yangqing Jia created the project during his PhD at UC Berkeley.

TensorFlow: TensorFlow is an open-source software library for dataflow programming across a range of tasks. It is a symbolic math library, and is also used for machine learning applications such as neural networks.

Pytorch: PyTorch is an open-source machine learning library for Python, based on Torch, used for applications such as natural language processing. It is primarily developed by Facebook's artificial-intelligence research group, and Uber's "Pyro" software for probabilistic programming is built on it.

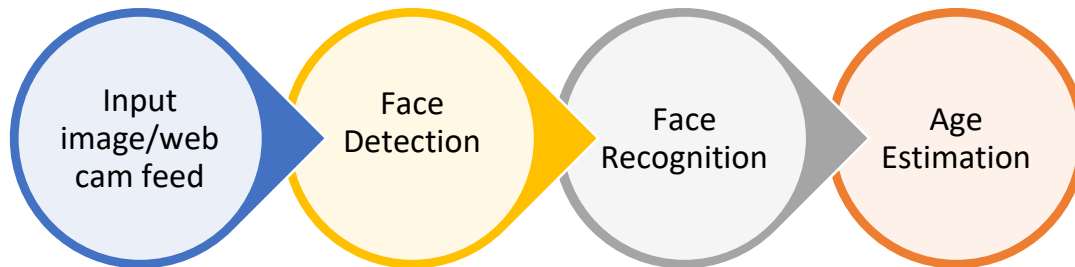
The GoogLeNet architecture (now known as "Inception" after the novel micro-architecture) was introduced by Szegedy et al. in their 2014 paper,[7] Going deeper with convolutions. It is important to note at this step that we aren't training a CNN, rather, we are making use of a pre-trained network. Therefore, we are just passing the blob through the network (i.e., forward propagation) to obtain the result (no back-propagation).

This is a direct implementation of deep learning modules, which have been trained by the processes of feed forward and back propagation as explained earlier in the report. Plus, they are heavily improved by the architectures like **GoogleLeNet**.

Face Recognition

There are many methods for face recognition which have been developed, but since we are using OpenCV face recognition, so we will restrict ourselves to those methods. All the face recognizers in work by first preparing the image or camera feed for face recognition and then the face recogniser is trained to recognize the faces.

There were mainly three types of face recognizers, but we will explain about the fourth one later in the report that we have used.



Figure(4.1)

Now, since the parts of the project are presented by the above steps, we would like to explain a variety of things on these parts. But before just moving on to these parts directly we will explain the background behind all of them –

- **Face detection-** The analogy behind it, different methods of implementing it and the technologies used.
- **Face recognition-** The analogy behind it, different methods of implementing it and the technologies used.
- **Age recognition-** This we have only studied in brief so we will just explain how and which method we have used.

Artificial Intelligence, Machine learning and Neural networks

AI is a term that was coined by scientists in Dartmouth conference in 1956. Over the past years the term has become more popular and has been accepted as a key to a bright future for our civilization. The main aim of the researchers in the initial years was to construct complex machines that could think like humans. It should have had all the senses that we humans have and maybe even more, Like the ones we have seen in sci-fi movies “The Terminator”. The term AI is a very broad term in itself, which includes several other terms that we are going to explain further.

The complex machines that the scientists and researchers wanted to build at that time, that could think like humans had to have some intelligence like humans also. Naturally this intelligence can't be put into them magically. We have to create some approaches that could enable the machines to think. **This “intelligence” is what is called “Machine Learning”.** Machine learning is the process in which we study the data, learn from it and then make some future predictions. The machine learns the ability to do a task by learning from its experiences again and again until it gets the idea. The most basic example is a new born baby, who doesn't know anything. But the child starts to learn from his/her experiences what is good and what is bad. If the child accidentally puts his finger on fire, he feels pain and then immediately removes his finger from fire. Now the next time he knows not to do something like that. Machine learning uses a lot of algorithms that helps it to learn from data and predict the outcomes. We have listed some algorithms below-

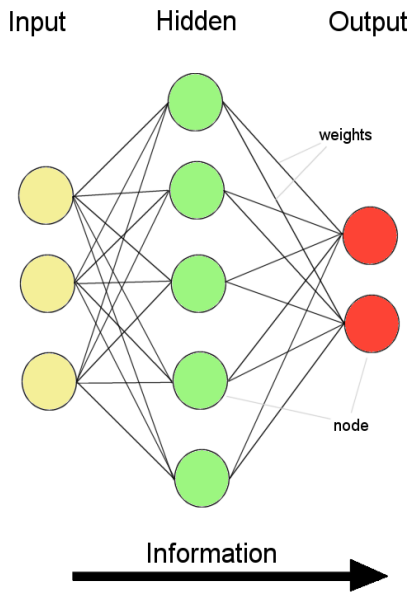
- Decision Tree
- Inductive logic programming
- Clustering
- Reinforcement learning
- Bayesian networks

Our project is based on **Neural networks, CNN and Deep learning** only so we will restrict our explanation up to that only, but for the sake of understanding we have given the above timeline with all the algorithms.

Neural Networks

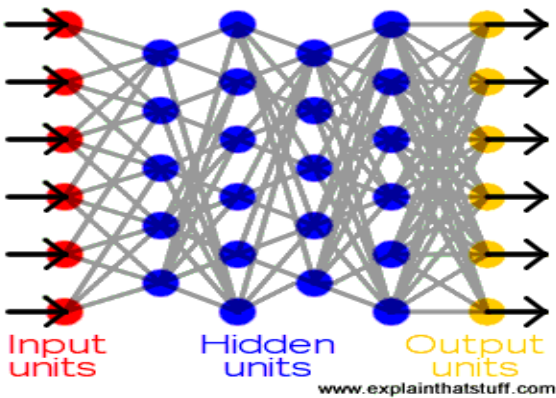
The Artificial Neural Nets are based on the biological neurons as mentioned earlier, so there two basic components of a Neural net

- Neurons(nodes)
- Synapses(weights)



(Figure 4.2)

A typical neural network consists of a lot of neurons (units). There are input units as well as output units, the former are designed to receive various types of information from the outside world and then it attempts to learn from that information. The output units respond to the information it has learned. There are one more type of units in between these units, called the hidden units which make up the majority of the neural nets. The connection between these units are called weights. The neural nets are fully connected.



Figure(4.3)

Learning in Neural Networks

The information in the network flows in two ways:

- **Learning phase:** Patterns of information is fed into the network via the input units, which activates the hidden layers and then the result arrives at the output layer. This design is called The **Feedforward network**. Each unit receives inputs from the units to its left, and the inputs are multiplied by the weights of the connections they travel along. Every unit adds up all the inputs it receives in this way and (in the simplest type of network) if the sum is more than a certain threshold value, the unit "fires" and triggers the units it's connected to (those on its right).
- **Back Propagation:** A feedback process is very important in a network similarly as we humans take feedback about our progress all the time. Hence a neural network also tends to do the same. Once the result is reached at the output layer, it is compared with the result it was supposed to produce. Then the difference between the two is used to adjust the weights of the connections between the units in the network that is going backwards. Hence it is called **Back Propagation**.

Once the network has been trained with enough learning examples, it reaches a point where you can present it with an entirely new set of inputs it's never seen before and see how it responds. For example, suppose you've been teaching a network by showing it lots of pictures of chairs and tables, represented in some appropriate way it can understand, and telling it whether each one is a chair or a table. After showing it, let's say, 25 different chairs and 25 different tables, you feed it a picture of some new design it's not encountered before, let's say a chaise longue, and see what happens. Depending on how you've trained it, it'll attempt to categorize the new example as either a chair or a table, generalizing on the basis of its past experience, just like a human. Hey presto, you've taught a computer how to recognize furniture.

Speech recognition using neural network learning:

Learning in the neural nets using feedforward and backpropagation can be explained using the example of speech recognition. Suppose there are two persons named "Steve" and "David". They both say the word "Hello", there are two frequency bins for each-

David= 1-0

Steve=0-1

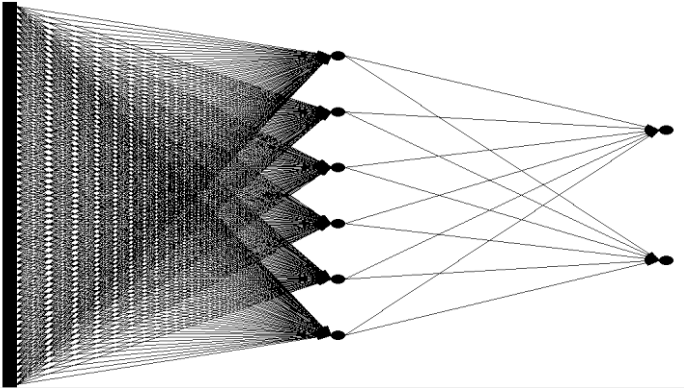
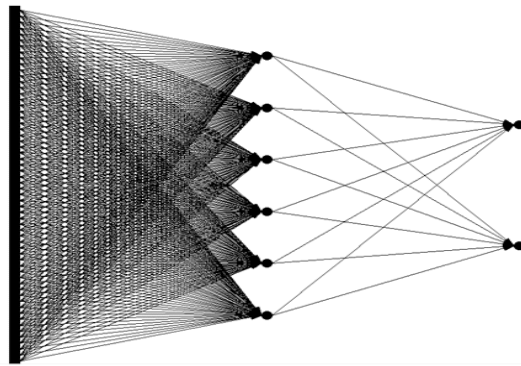
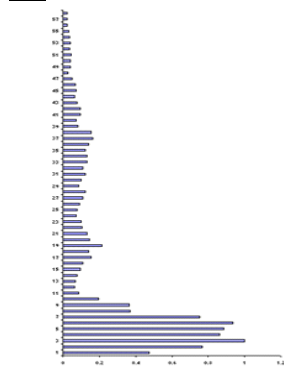


Figure (4.4)

There are 6 hidden layers in the network and 2 output layers. Now in the first phase, the untrained network is given the inputs for which it produces the output as shown. The initial values, output by the network obviously has errors.

Steve

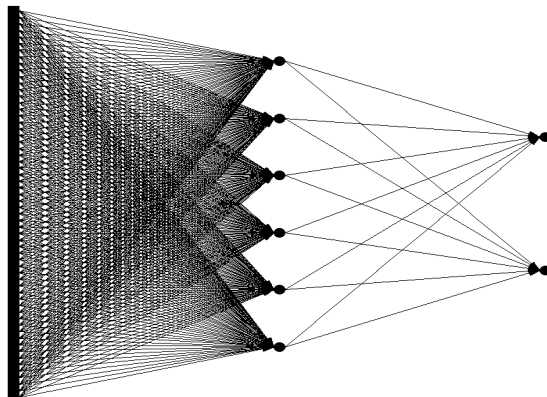
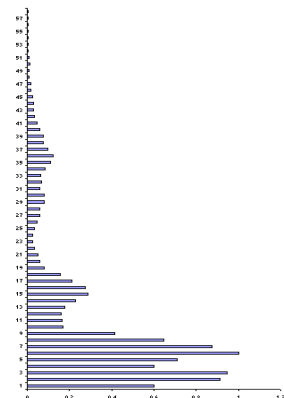


0.43

0.26

Figure(4.5)

David



0.73

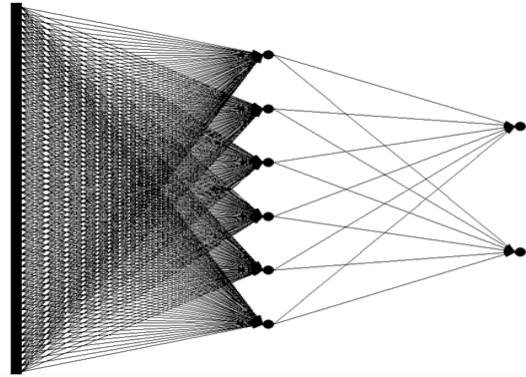
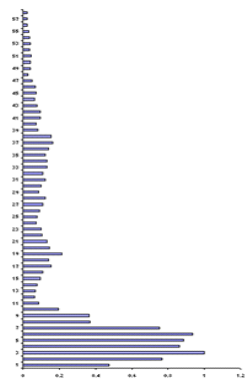
0.55

Figure (4.6)

Now in the next phase we will calculate the error in the output and then as we discussed earlier the neural network will make some changes to the weights so that correct output can be produced.

Calculating error

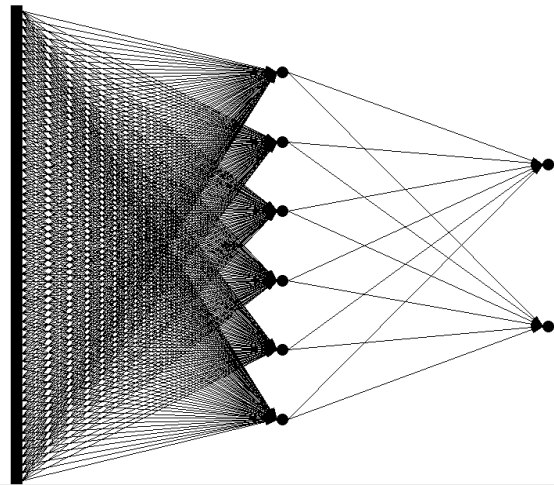
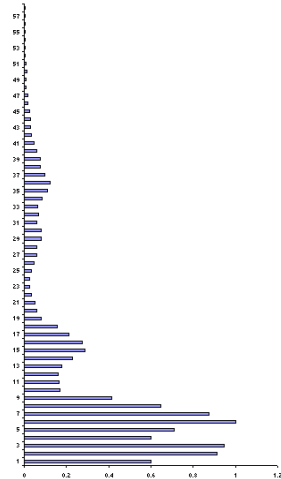
Steve



$$0.43 - 0 = 0.43$$

$$0.26 - 1 = 0.74$$

David



$$0.73 - 1 = 0.27$$

$$0.55 - 0 = 0.55$$

The total error in Steve's network is $0.43 + 0.74 = 1.17$

The total error in David's network is $0.27 + 0.55 = 0.82$

After this calculation the network will adjust its weights until its done with correct output.

As Neural network is a class of machine learning algorithms, there are different variations of neural networks. The class of Neural networks contains various architectures like **Convolutional neural networks (CNN)**, **Recurrent neural networks (RNN)** and **Deep belief networks**. The number of (layers of) units, their types, and the way they are connected to each other is called the **network architecture**.

Now, for our project we have used **CNN Architecture**. we will explain this in detail. A CNN consists of a **convolutional layer**, a **pooling layer** and **fully connected layer**.

- **The convolutional layer**: It is the first layer to extract features from the input image, it creates a relationship between pixels by learning image features using small squares of input data. It is a mathematical operation.

Suppose there is a 3 x 3 matrix with image pixel values 0,1 and a filter matrix as shown below

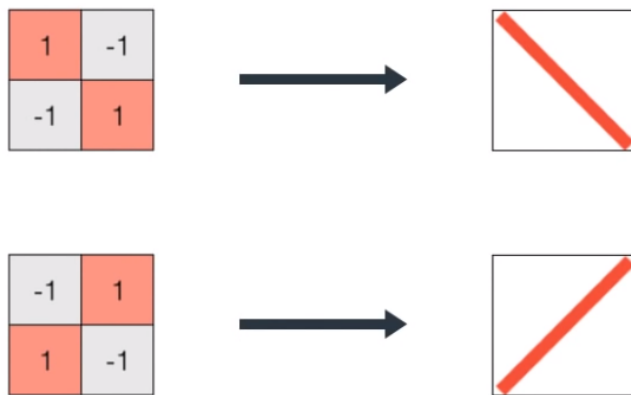
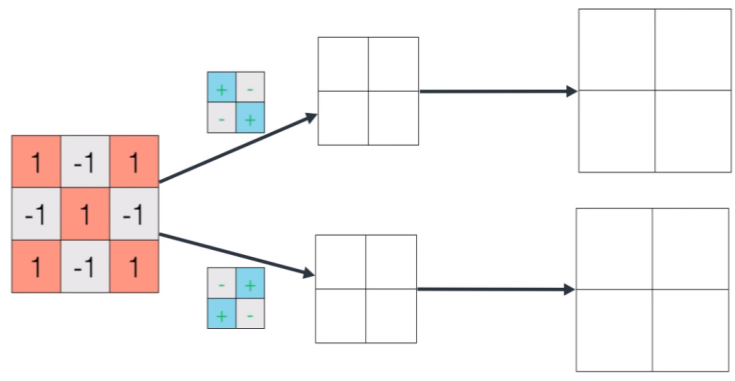


Figure (5.1)

This is the filter matrix that we know creates a forward slash and backward slash, this is the previous knowledge we have for a 2x2 matrix. Now this matrix will be superimposed on the 3x3 matrix from left to right and top to bottom to get smaller outputs from the larger input that is 3x3.



(Figure 5.2)

Filter one

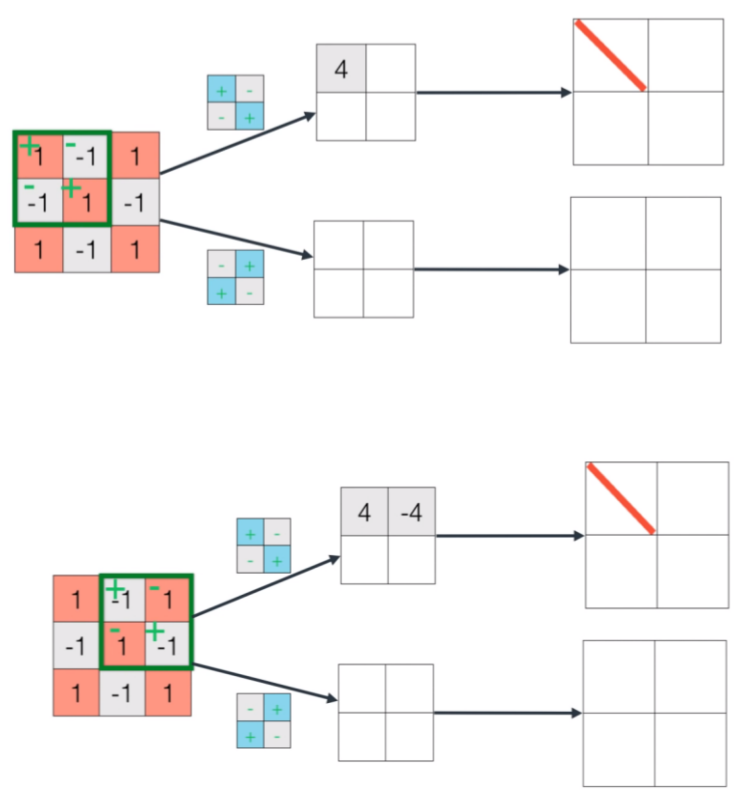
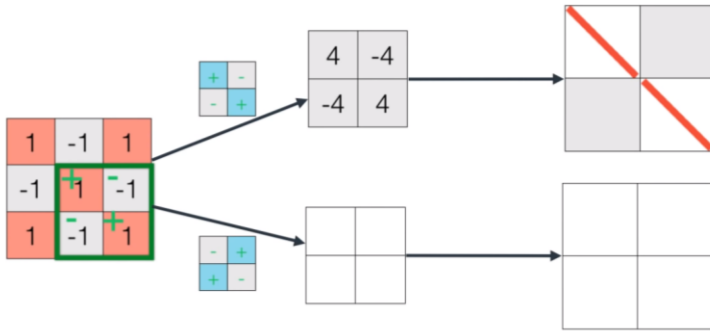
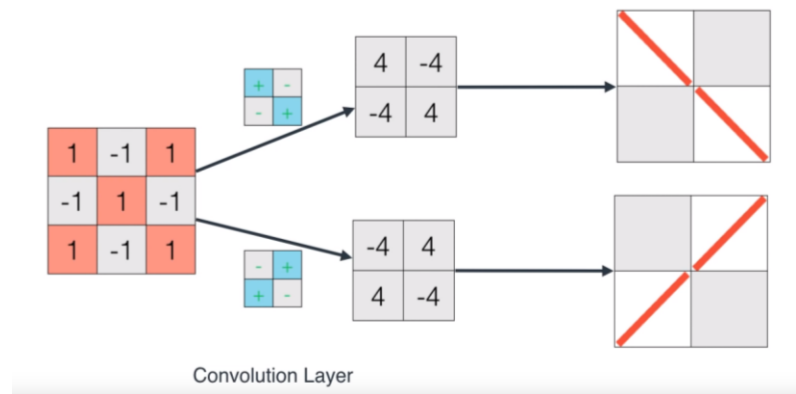


Figure (5.3)



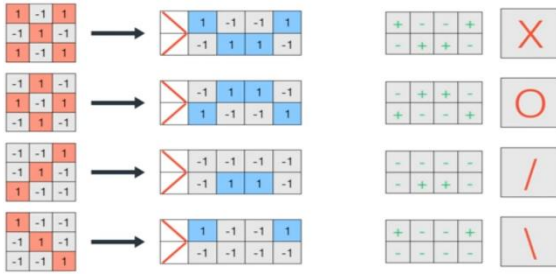
(Figure 5.4)

Filter 2



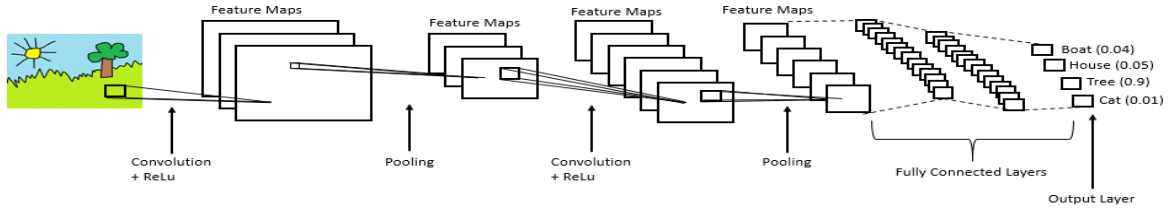
Figure(5.5)

- **Pooling layer:** Pooling layers section would reduce the number of parameters when the images are too large. In the above image the areas that are grey do not satisfy the criterion, so it not useful. This is the work of the pooling layer.
- **Fully Connected Layer:** This layer finds logic and tries to figure out which image that might be.



Figure(5.6)

The process can be visualised by the below image. These are the series of steps that are taken up in a Convolutional Neural Network Architecture for image recognition .



Figure(5.7)

Deep Learning Face recognition

The most accurate implementation of face recognition is using deep learning. Well with that being said this could only be possible with OpenCV 3 and above. Otherwise implementing this algorithm would have been a lot tougher.

Deep learning combined with face recognition is called **deep metric learning**. This means instead of output as a single label of an image, this algorithm outputs a 128 -d feature vector, i.e a list of 128 real valued numbers that are used to quantify a face.

The network quantifies the faces, constructing the 128-d embedding for each. From there, the general idea is that it'll tweak the weights of the neural network so that the 128-d measurements of the two Will Ferrel will be closer to each other and farther from the measurements for Chad Smith. This network architecture for face recognition is based on [5]ResNet-34 from the Deep Residual Learning for Image Recognition paper by Kaiming He., but with fewer layers and the number of filters reduced by half. This network was trained on the dataset LFW (Labelled faces in the wild), by Davis King (creator of Dlib library) and He claims the accuracy to be 99.38%, This accuracy is supported by the dlib documentation.

Deep learning and Convolutional neural networks

The confusion arises to a lot of people about what is what. So we would like to explain a little about both. Everything that is done in CNN is also done in Deep learning that is described in the below picture-

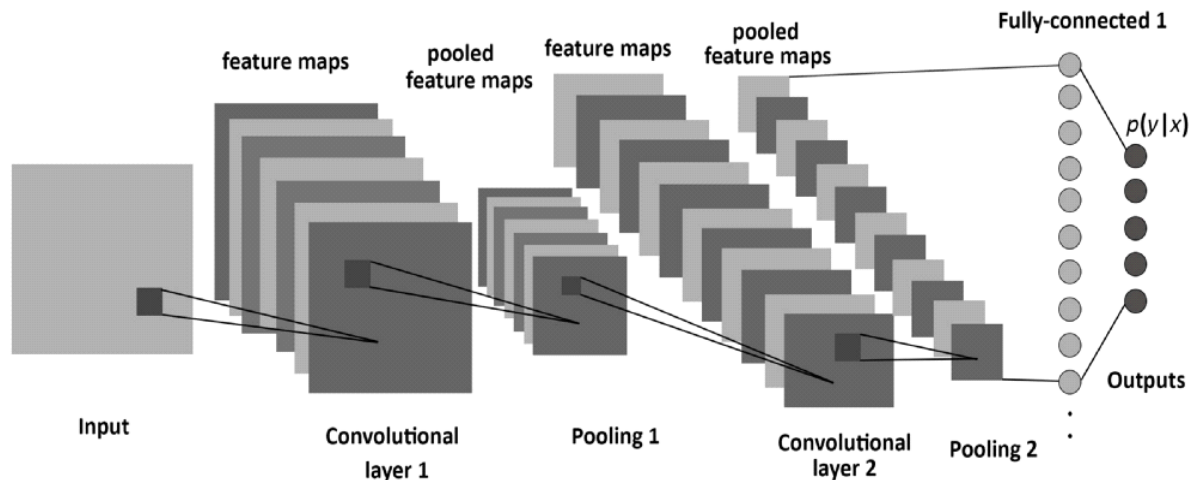
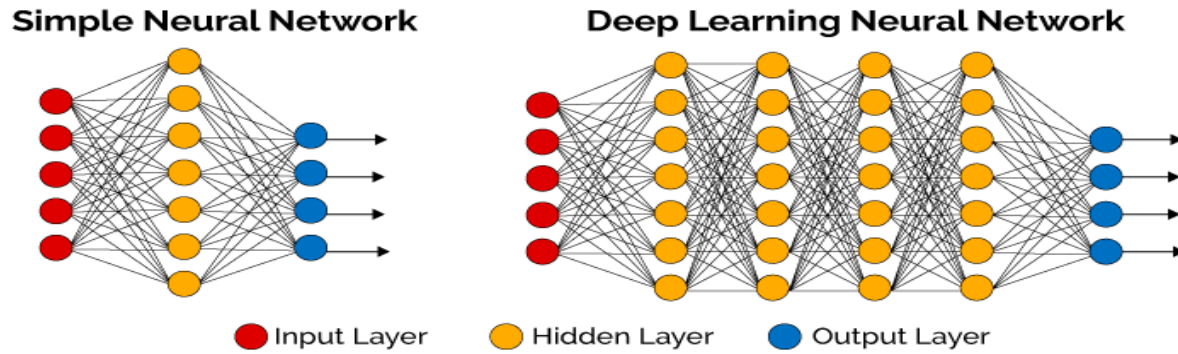


Figure (5.8)

The Deep learning uses CNN , and the clear distinction between both is any CNN that has more than 1 hidden layer is called a Deep learning network. This also has an effect on the training time. Only computers with a GPU can use Deep learning networks to train. The accuracy between both is also different, mainly deep neural networks have a lot of layers and hence more the layers, more is the accuracy.



(Figure 5.9)

Age Recognition using Deep learning

There have been various attempts on Age recognition, like the early ones include one facial feature (eyes, nose, mouth, chin), are localized and their sizes and distances are measured, ratios between them are calculated and then used for classifying images. Another model similar to this was age progression of subjects under 18 years. All of these methods require an accurate localization of features which is a tough job in itself.

As it is, we know that Deep learning is possible with today’s technology and it gives better results than any other methodology, so we used Deep learning for this model taken from [1]. The proposed model was used by the researchers throughout their experiments and is shown below.

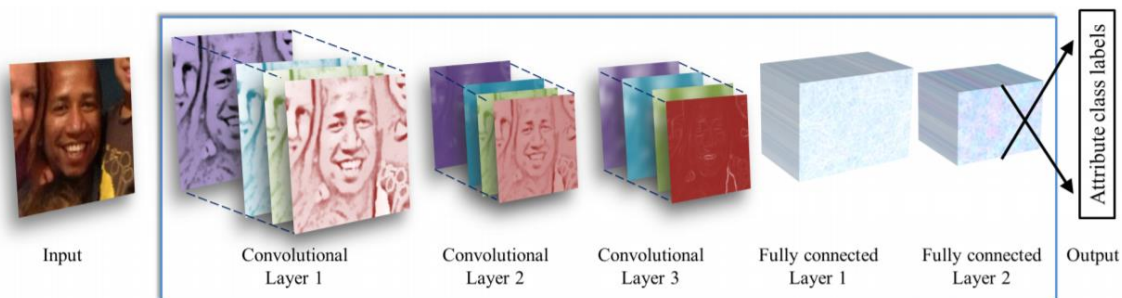


Figure 2. **Illustration of our CNN architecture.** The network contains three convolutional layers, each followed by a rectified linear operation and pooling layer. The first two layers also follow normalization using local response normalization [28]. The first Convolutional Layer contains 96 filters of 7×7 pixels, the second Convolutional Layer contains 256 filters of 5×5 pixels, The third and final Convolutional Layer contains 384 filters of 3×3 pixels. Finally, two fully-connected layers are added, each containing 512 neurons. See Figure 3 for a detailed schematic view and the text for more information.

(Figure 5.10)

The network comprises of only three convolutional layers and two fully-connected layers with a small number of neurons, all three color channels are processed directly by the network. Images are first rescaled to 256×256 and a crop of 227×227 is fed to the network. The three subsequent convolutional layers are then defined as follows.

- 96 filters of size $3 \times 7 \times 7$ pixels are applied to the input in the first convolutional layer, followed by a rectified linear operator (ReLU), a max pooling layer is taking the maximal value of 3×3 regions with two-pixel strides and a local response normalization layer.
- The $96 \times 28 \times 28$ output of the previous layer is then processed by the second convolutional layer, containing 256 filters of size $96 \times 5 \times 5$ pixels. Again, this is followed by ReLU, a max pooling layer and a local response normalization layer with same hyper parameters as before.
- Finally, the third and the last convolution layer operates on the $256 \times 14 \times 14$ blob by applying a set of 384 filters of size $256 \times 3 \times 3$ pixels, followed by ReLU and a max pooling layer.
- A first fully connected layer that receives the output of the third convolutional layer and contains 512 neurons, followed by a ReLU and a dropout layer.
- A second fully connected layer that receives the 512- dimensional output of the first fully connected layer and again contains 512 neurons, followed by a ReLU and a dropout layer.
- A third, fully connected layer which maps to the final classes for age.
- Finally, the output of the last fully connected layer is fed to a soft-max layer that assigns a probability for each class. The prediction itself is made by taking the class with the maximal probability for the given test image.

Training the network

Initialization. The weights in all layers are initialized with random values from a zero mean Gaussian with standard deviation of 0.01. To stress this, they do not use pre-trained models for initializing the network; the network is trained, from scratch, without using any data outside of the images and the labels available by the benchmark. This, again, should be compared with CNN implementations used for face recognition, where hundreds of thousands of images are used for training. Target values for training are represented as sparse, binary vectors corresponding to the ground truth classes. For each training image, the target, label vector is in the length of the number of classes (two for gender, eight for the eight age classes of the age classification task), containing 1 in the index of the ground truth and 0 elsewhere.

A possible future application for facial recognition systems lies in retailing. A retail store (for example, a grocery store) may have cash registers equipped with cameras; the cameras would be aimed at the faces of customers, so pictures of customers could be obtained. The camera would be the primary means of identifying the customer, and if visual identification failed, the customer could complete the purchase by using a PIN (personal identification number). After the cash register had calculated the total sale, the face recognition system would verify the identity of the customer and the total amount of the sale would be deducted from the customer's bank account. Hence, face-based retailing would provide convenience for retail customers, since they could go shopping simply by showing their faces, and there would be no need to bring debit cards, or other financial media. Wide-reaching applications of face-based retailing are possible, including retail stores, restaurants, movie theaters, car rental companies, hotels, etc. e.g. Swiss European surveillance: facial recognition and vehicle make, model, color and license plate reader.

Some other possible applications that can be developed are:

- 1.** In order to prevent the frauds of ATM, it is recommended to prepare the database of all ATM customers with the banks in India & deployment of high-resolution camera and face recognition software at all ATMs. So, whenever user will enter in ATM his photograph will be taken to permit the access after it is being matched with stored photo from the database.
- 2.** Duplicate voter are being reported in India. To prevent this, a database of all voters, of course, of all constituencies, is recommended to be prepared. Then at the time of voting the resolution camera and face recognition equipped of voting site will accept a subject face 100% and generates the recognition for voting if match is found.
- 3.** Passport and visa verification can also be done using face recognition technology as explained above.
- 4.** Driving license verification can also be exercised face recognition technology as mentioned earlier.
- 5.** To identify and verify terrorists at airports, railway stations and malls the face recognition technology will be the best choice in India as compared with other biometric technologies since other technologies cannot be helpful in crowded places.

- [1] M. El Ayaadi, F. Karrae and M. S. Kamal “Survey on expression recognition: Features, classification scheme, and database,” *Pattern Recognit.*, vol. 44, no. 3, pp. 572–587, 2011.
- [2] Y. Lecan, Y. Bengio, and G. Hintin, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [3] J. Schmidhuber, “Deep Learning in neural networks: An overview,” *Neural Networks*, vol. 61, pp. 85–117, 2015.
- [4] J. Ngiem, A. Khoslaa, M. Kam, J. Nim, H. Lei, and A. Y. Nig, “Multimodal Deep Learning,” *Proc. 28th Int. Conf. Mach. Learn.*, pp. 689–696, 2011.
- [5] S. Lugivic M. Harva and I. Dander “Techniques and applications of emotion recognition,” 2016 39th Int. Conv. Inf. Commun. Technol. Electron. Microelectron. MIPRO 2016 - Proc., no. November 2017, pp. 1278–1283, 2016.
- [6] B. Schaller, G. Rigull, and M. Ling, “Emotion recognition combining acoustic features and information in a neural network - belief network architecture,” *Acoust. Speech, Signal Process.*, vol. 1, pp. 577–580, 2004.
- [7] J. Rieng, G. Leie, and Y. P. P. Chen, “Acoustic feature selection for automatic emotion recognition from speech,” *Inf. Process. Manag.*, vol. 45, no. 3, pp. 315–328, 2009.
- [8] F. Noruzi, G. Anbarjafaari and N. Akraami “Expression-based emotion recognition and next reaction prediction,” 2017 25th Signal Process. Commun. Appl. Conf. SIU 2017, 2017.
- [9] G. Hintin ,Greves and A. Mohamed, “Emotion Recognition with Deep Recurrent Neural Networks,” in 2013 IEEE International Conference on Acoustics, Speech and Signal Processing, 2013, pp. 6645–6649.
- [10] K. Weint and C.-W. Huang, “Characterizing Types of Convolution in Deep Convolutional Recurrent Neural Networks for Robust Speech Emotion Recognition,” 2017.
- [11] Face Recognition Data, University of Essex, UK, Face 94, http://cswww.essex.ac.uk/mv/all_faces/faces94.html.
- [12] Face Recognition Data, University of Essex, UK, Face 95, http://cswww.essex.ac.uk/mv/all_faces/faces95.html.
- [13] Face Recognition Data, University of Essex, UK, Face 96, http://cswww.essex.ac.uk/mv/all_faces/faces96.html.
- [14] Face Recognition Data, University of Essex, UK, Grimace, http://cswww.essex.ac.uk/mv/all_faces/grimace.html.