



OSMANIA UNIVERSITY

POSE INVARIANT FACE RECOGNITION USING CNN PROJECT REPORT

Submitted

IN THE PARTIAL FULFILLMENT OF the award of the degree of

Bachelor of Engineering In

ELECTRONICS AND COMMUNICATION ENGINEERING BY

Shravan Kanuganti (2451-16-735-046)

Under the Guidance of

Dr. Sarala Beeram

Designation: Professor



MATURI VENKATA SUBBA RAO (MVSR) ENGINEERING COLLEGE

(Sponsored by Matrusri Education Society, Estd 1981)

(Approved by AICTE & Affiliated to OU)

(Accredited by NBA & NAAC)

JUNE, 2019-2020

DECLARATION

We hereby declare that the results embodied in this dissertation entitled “**Pose Invariant Face Recognition using CNN**” are carried out by us at M.V.S.R. Engineering College during the year 2019 – 2020 in partial fulfilment of the award of **B.E. (Bachelor of Engineering)** from “**M.V.S.R. ENGINEERING COLLEGE**”, affiliated to Osmania University (O.U). We have not submitted the same to any other university or organization for the award of any other degree.

No part of the thesis is copied from books/journals/internet and wherever referred, the same has been duly acknowledged in the text. The reported data are based on the project work done entirely by us and not copied from any other sources.

SHRAVAN KANUGANTI (2451-16-735-046)

CERTIFICATE

This is to certify that the Project report titled **“Pose Invariant Face Recognition using CNN”** submitted by **SHRAVAN KANUGANTI (2451-16-735-046)**, Student of Electronics & Communication Engineering, MVSR Engineering College, Hyderabad in partial fulfilment of the requirements for the award of Degree of Bachelor of Engineering is a record of the bonafide work carried out by them during the academic year 2019-2020.

Dr. S. P. VENU MADHAVA RAO
(Prof. & Head, ECE)

Dr. Sarala Beeram
(Prof., Internal Guide)

ACKNOWLEDGEMENT

The satisfaction and euphoria that accompany the successful completion of our project until now would be incomplete without mentioning of the people who made it possible, whose constant guidance and encouragement crowded our effort with success.

We are greatly thankful to **Dr. S. P. Venu Madhava Rao**, Professor and Head of the Department of Electronics and Communication Engineering, for his encouragement and advice during our work.

We gratefully acknowledge to our project coordinators **Mrs. T. Kavitha** and **S.V. R. Manimala** our internal guide Professor **Dr. B. Sarala**, panel members and lectures of Department of Electronics and Communication Engineering, for their excellent guidance on selecting the topics, constant inspiration and encouragement in pursuing the project work.

We would like to express our thanks to all the members of the project review committee for their valuable inputs.

We would like to express our gratitude to all people behind the screen who helped us to transform an idea into a real application. We would like to express our heart-felt gratitude to our parents and our siblings without whom we would not have been privileged to achieve and fulfill our dreams and all the faculty members, Staff of Department of Electronics and Communication Engineering of Maturi Venkata Subba Rao Engineering College and all other members who directly or indirectly rendered help in making this seminar a successful one.

SHRAVAN KANUGANTI (2451-16-735-046).

Contents

Abstract	i
List of figures	ii
List of tables	iv
Chapter One	3
INTRODUCTION	3
1.1 Context	3
1.2 Applications	4
1.3 Difficulties	6
1.3.1 Illumination	6
1.3.2 Pose	7
1.3.3 Facial Expressions	7
1.3.4 Partial Occlusions	8
1.3.5 Other types of variations	8
1.4 Project Significance	9
1.5 Project Approach	9
1.6 Project Outcomes	10
1.7 Conclusion	10
Chapter Two	11
LITERATURE REVIEW	11
2.1 Introduction	11
2.2. Emerging Trends in Face Recognition	11
2.3 Face Recognition Structure	14
2.4 Face Recognition Methods	16
2.4.1 Holistic Matching Methods	16
2.4.2 Feature Based(structural) Methods	17
2.4.3 Hybrid Methods	17
2.5 Conclusion	18
Chapter Three	19
POSE INVARIANT FACIAL RECOGNITION USING CONVOLUTIONAL NEURAL NETWORK (CNN)	19
3.1 Introduction	19
3.1.1 Neural Networks	20
3.1.2 Perceptron	20
3.1.3 Multi-Layer Perceptron	21

3.1.4 Training Neural Networks	21
3.1.5. Backpropagation Algorithm	22
3.1.6. Advantages of using Neural Networks	25
3.1.7. Disadvantages of using Neural Networks	25
3.2. Convolutional Neural Networks	26
3.2.1. Design	26
3.2.2. Layers	30
3.2.3. Classification	36
3.3. Pose-Invariant Face Recognition	40
3.4. CNN Block Diagram	41
3.4.1. Specification of the Architecture	42
3.4.2. Advantages of CNN	42
3.4.3. Limitations of CNN	42
3.4.4. Case Studies	43
3.5. Conclusion	44
Chapter Four	45
TOOLS & LIBRARIES USED	45
4.1. Introduction	45
4.2. Programming Languages	45
4.2.1. Python 3.7	45
4.3. Applications Used	46
4.3.1. Irfan View	46
4.4. Conclusion	46
Chapter Five	47
RESULTS & PROJECT ANALYSIS	47
5.1. Databases Used	47
5.2. Training Results	50
5.3. Testing Results	52
5.5. Hardware Comparison	53
5.6 Code	54
CHAPTER Six	55
CONCLUSION & FUTURE SCOPE	55
6.1 CONCLUSION	55
6.2 FUTURE SCOPE	57
CHAPTER Seven	58



ABSTRACT

This abstract aims at the development of a robust face recognition system which is suitable in identifying faces irrespective of various alignments or poses of the face. Most of the existing algorithms in face detection are only successful in controlled environments and fail drastically when non linear functions such as pose variation, illumination and occlusion are affecting the image. Face recognition under uncontrolled environment is most important aspect in enhancing HCI (Human Computer Interface). Convolutional Neural Networks is used in applying feature extraction to normalize the data causing the system to cope with faces subject to pose variation.

CNN is a class of deep neural networks which is most commonly applied in analysing and classification of visual imagery. It uses relatively little pre- processing compared to other image classification algorithms. The main benefit of using CNN over simple ANN (Artificial Neural Networks) is that CNN's are constrained to deal with image data exclusively. One of the main features of this algorithm is **weight sharing**, as a result it reduces the number of weights significantly. Furthermore processing of high quality images which consists of very high number of pixels is a tedious task when implemented in ANN, here we can understand the need of CNN and how they can scale to high resolution images also. Another advantage of it is that they are very good feature extractors. This means that you can extract useful attributes from an already trained CNN with its trained weights by feeding your data on each level and tune it a bit for the specific task. It is also used in video recognition, recommender systems, natural language processing, etc.

In this project, we implement an enhanced CNN architecture based on the superior **AlexNet**, developed by SuperVision group in 2012. It consists of 11x11, 5x5, 3x3 convolutions, max pooling and **ReLU** activations. Based on the features for extraction of various features of a face which will be helpful in detecting the face irrespective of non-linear functions as stated above, we compare the performance of existing algorithms and proposed algorithm.

KEYWORDS: Face Recognition, CNN, Image processing and Pose Variation.



LIST OF FIGURES

1.1	Block diagram of Face Recognition	3
1.2	An example of faces under a fixed view and varying illumination	6
1.3	An example of faces under varying pose	7
1.4	An example of faces under fixed illum. And pose but varying facial expression	7
1.5	An example of faces with occlusions	8
1.6	An example of faces with other type of variations	8
2.1	Block diagram of facial recognition	14
2.2	Eigen Faces	16
2.3	FR based on features	16
2.4	Hybrid Feature extraction using SVM	17
3.1	Representation of a CNN	18
3.2	Perceptron	19
3.3	Different types of activation functions	20
3.4	Typical CNN Architecture	25
3.5	Convolution	26
3.6	Pooling	27
3.7	Fully Connected	27
3.8	Neurons(blue) connected to their Receptive Field(red)	28
3.9	Weight sharing in CNNs	28
3.10	Convolution Layer	30
3.11	Pooling Layers	31
3.12	Linear Activation	32
3.13	Non-Linear Activation	32
3.14	Sigmoid Activation	33



3.15	Tanh Activation	33
3.16	ReLu Activation	34
3.17	Leaky ReLu	34



3.18	Classification Using Log. Reg.	36
3.19	CNN Block Diagram	40
5.1	ATT Database Images	46
5.2	FERET Database Images	47
5.3	Webcam Database Images	48
5.4	ATT Training Results	49
5.5	FERET Training Results	49
5.6	Webcam Training Results	50
5.7	Testing Comparison	51
5.8	Hardware Comparison	52



LIST OF TABLES

3.1	Different Types of Activation Layers	35
-----	---	----

Chapter One

INTRODUCTION

1.1 Context

Face recognition has been one of the most intensively studied topics in computer vision for more than four decades. Compared with other popular biometrics such as finger print, iris, and retina recognition, face recognition has the potential to recognize uncooperative subjects in a non-intrusive manner. Therefore, it can be applied to surveillance security, border control, forensics, digital entertainment, etc. Indeed, numerous works in face recognition have been completed and great progress has been achieved, from successfully identifying criminal suspects from surveillance cameras to approaching human level performance. These successful cases, however, may be unrealistically optimistic as they are limited to near-frontal face recognition (NFFR). Recent studies reveal that the best NFFR algorithms perform poorly in recognizing faces with large poses. In fact, the key ability of pose-invariant face recognition (PIFR) desired by real-world applications remains largely unsolved.

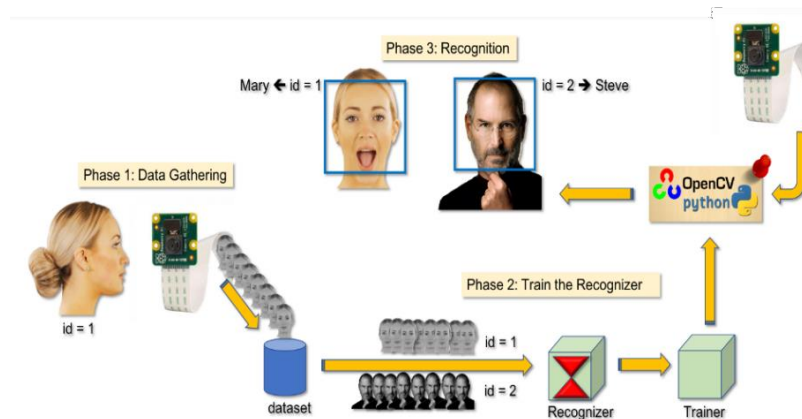


Fig. 1.1 Block diagram of Face Recognition

In the context of human-computer interaction (HCI), it might also be important to detect the position of specific facial characteristics or recognize facial expressions, in order to allow, for example, a more intuitive communication between the device and the user or to efficiently encode and transmit facial images coming from a camera. Thus, the automatic analysis of face images is crucial for many applications involving visual content retrieval or extraction.

The principal aim of facial analysis is to extract valuable information from face images, such as its position in the image, facial characteristics, facial expressions, the person's gender or identity. We will outline the most important existing approaches to facial image analysis and present novel methods based on Convolutional Neural Networks (CNN) to detect, normalize and recognize faces and facial features. CNNs show to be a powerful and flexible feature extraction and classification technique which has been successfully applied in other contexts, i.e. hand-written character recognition, and which is very appropriate for face analysis problems as we will experimentally show in this work. We will focus on the processing of two-dimensional gray-level images as this is the most widespread form of digital images and thus allows the proposed approaches to be applied in the most extensive and generic way. However, many techniques described in this work could also be extended to color images, 3D data or multi-modal data.

1.2 Applications

There are numerous possible applications for facial image processing algorithms. The most important of them concern face recognition. In this regard, one has to differentiate between closed world and open world settings. In a closed world application, the algorithm is dedicated to a limited group of persons, e.g. to recognize the members of a family. In an open world context, the algorithm should be able to deal with images from "unknown" persons, i.e. persons that have not been presented to the system during its design or training. For example, an application indexing large image databases like Google images or television programs should recognize learned persons and respond with "unknown" if the person is not in the database of registered persons.

Concerning face recognition, there further exist two types of problems: face identification and face verification (or authentication). The first problem, face

identification, is to determine the identity of a person on an image. The second one only deals with the question: “Is ‘X’ the identity of the person shown on the image?” or “Is the person shown on the image the one he claims to be?”. These questions only require “yes” or “no” as the answer. Possible applications for face authentication are mainly concerned with access control, e.g. restricting the physical access to a building, such as a corporate building, a secured zone of an airport, a house etc. Instead of opening a door by a key or a code, the respective person would communicate an identifier, e.g. his/her name, and present his/her face to a camera. The face authentication system would then verify the identity of the person and grant or refuse the access accordingly. This principle could equally be applied to the access to systems, automatic teller machines, biometric scanners, mobile phones, and Internet sites etc. where one would present his face to a camera instead of entering an identification number or password. Clearly, also face identification can be used for controlling access. In this case the person only has to present his/her face to the camera without claiming his/her identity. A system recognizing the identity of a person can further be employed to control more specifically the rights of the respective persons stored in its database. For instance, parents could allow their children to watch only certain television programs or web sites, while the television or computer would automatically recognize the persons in front of it.

Video surveillance is another application of face identification. The aim here is to recognize suspects or criminals using video cameras installed at public places, such as banks or airports, in order to increase the overall security of these places. In this context, the database of suspects to recognize is often very large and the images captured by the camera are of low quality, which makes the task rather difficult. With the vast propagation of digital cameras in the last years the number of digital images stored on servers and personal home computers is rapidly growing. Consequently, there is an increasing need of indexation systems that automatically categorize and annotate this huge number of images in order to allow effective searching and so-called content- based image retrieval. Here, face detection and recognition methods play a crucial role because a great part of photographs actually contain faces. A similar application is the temporal segmentation and indexation of video

sequences, such as TV programs, where different scenes are often characterized by different faces.

1.3 Difficulties

There are some inherent properties of faces as well as the way the images are captured which make the automatic processing of face images a rather difficult task. In the case of face recognition, this leads to the problem that the intra- class variance, i.e. variations of the face of the same person due to lighting, pose etc., is often higher than the inter-class variance, i.e. variations of facial appearance of different persons, and thus reduces the recognition rate. In many face analysis applications, the appearance variation resulting from these circumstances can also be considered as noise as it makes the desired information, i.e. the identity of the person, harder to extract and reduces the overall performance of the respective systems. In the following, we will outline the most important difficulties encountered in common real-world applications.

1.3.1 Illumination

Changes in illumination can entail considerable variations of the appearance of faces and thus face images. Two main types of light sources influence the overall illumination: ambient light and point light (or directed light). The former is somehow easier to handle because it only affects the overall brightness of the resulting image. The latter however is far more difficult to analyze, as face images taken under varying light source directions follow a highly non-linear function. Additionally, the face can cast shadows on itself..



Fig. 1.2 An example of faces under a fixed view and varying illumination

Many approaches have been proposed to deal with this problem. Some face detection or recognition methods try to be invariant to illumination changes by implicitly modeling them or extracting invariant features. Others propose a

separate processing step, a kind of normalization, in order to reduce the effect of illumination changes.

1.3.2 Pose

The variation of head poses or, in other words, the viewing angle from which the image of the face was taken is another difficulty and essentially impacts the performance of automatic face analysis methods. For this reason, many applications limit themselves to more or less frontal face images or otherwise perform a pose-specific processing that requires a preceding estimation of the pose, like in multi-view face recognition approaches.

If the rotation of the head coincides with the image plane the pose can be normalized by estimating the rotation angle and turning the image such that the face is in an upright position. This type of normalization is part of a procedure called face alignment or face registration.



Fig. 1.3 An example of faces under varying pose

1.3.3 Facial Expressions

The appearance of a face with different facial expressions varies considerably (see Fig. 1.4). In general, the mouth is subject to the largest variation. The respective person on an image can have an open or closed mouth, can be speaking, smiling, laughing or even making grimaces. Eyes and eyebrows are also changing subject to varying facial expressions, e.g. when the respective person blinks, sleeps or widely opens his/her eyes.



Fig. 1.4 An example of faces under fixed illum. and pose but varying facial expression

1.3.4 Partial Occlusions

Partial occlusions occur quite frequently in real-world face images. They can be caused by a hand occluding a part of the face, e.g. the mouth, by long hair, glasses, sun glasses or other objects or persons.



Fig. 1.5 An example of faces with occlusions

1.3.5 Other types of variations

Appearance variations are also caused by varying make-up, varying hair-cut and the presence of facial hair (beard, mustache etc.). Varying age is also an important factor influencing the performance of many face analysis methods. There are also variations across the subjects' identities, such as race, skin color or, more generally, ethnic origin. The respective differences



in the

appearance of the face images can cause difficulties in applications like face or facial feature detection or gender recognition.



Fig. 1.6 An example of faces with other type of variations

1.4 Project Significance

The basic idea behind this project is to identify faces under different pose variations and achieve a higher recognition rate using CNN architecture.

1.5 Project Approach

The goals pursued in this work principally concern the evaluation of Convolutional Neural Networks (CNN) in the context of facial analysis applications. More specifically, we will focus on the following objectives:

- Evaluate the performance of CNNs w.r.t. appearance-based facial analysis.
- Investigate the robustness of CNNs against classical datasets like FERET, ATT, etc.
- Propose different CNN architectures designed for specific facial analysis problems such as face alignment, facial feature detection and face recognition.
- Extend the project to appearance-based facial feature detection, face alignment as well as face recognition under real-world conditions.
- Extend the project to synthesize and train images for facial 3D pose estimation.



1.6 Project Outcomes

To identify images with higher recognition rate with minimum pre-processing. In terms of performance, CNNs outperform NNs on conventional image recognition tasks and many other tasks.

1.7 Conclusion

This chapter gives the brief introduction to project and describes the development in Face recognition.

Chapter Two

LITERATURE REVIEW

2.1 Introduction

Recognizing an individual is one thing that every human does effortlessly and without much conscious thought. But it remains a problem when a computer has to identify or recognize a person. Biometric is one of the areas where face recognition is used widely. Some of the biometric techniques are iris detection, finger print detection or voice detection. However, face recognition holds many advantages when compared with the above mentioned biometric technique. Below are some of the evolving techniques in FR field.

2.2. Emerging Trends in Face Recognition

- **Face Recognition by Support Vector Machines (2000):** Uses Cambridge ORL database. The experimental results show that this method is better than nearest center approach for face recognition.
- **Face Recognition with Support Vector Machines (2001):** the database includes the faces rotated in depth up to 40°. This method shows that using facial components as facial features simplifies the recognition task rather than using whole face.
- **Facial Component Extraction and Face Recognition with Support Vector Machines (2002):** this method is quick and robust where the algorithm is applied to the faces of different sizes.
- **Face Recognition Using Support Vector Machines with the Robust Feature (2003):** The ORL database is used for testing. This algorithm introduces the Gabor wavelet, KCPA and SVM. The SVM is used to classify the robust features and this obtains high accuracy.
- **Improving the Performance of Multi-Class SVMs in Face Recognition with Nearest Neighbor Rule (2003):** NNR method is introduced to mainly reduce training class levels. This process performs much faster than the available SVM methods and the classification process is much faster.
- **A SVM-Based Method for Face Recognition using a Wavelet PCA Representation of Faces (2004):** This method uses the concept of SVM. The proposed method includes a substantial reduction in error rate.

- **Bayesian Face Recognition Using Support Vector Machine and Face Clustering (2004):** We first develop a direct Bayesian based Support Vector Machine by combining the Bayesian analysis with the SVM. Then the experiment is carried out using the two traditional databases FERET and XM2VTS. This method also yields some acceptable results.
- **Face recognition in color image using PCA and FSVM (2005):** This method uses color images for recognition process. This method uses skin color as one clue to recognize the face. Experiments were conducted using the indoors photograph and the results demonstrated that the proposed method was efficient for the frontal face detection and recognition.
- **A SVM Face Recognition Method Based on Gabor-Featured Key Points (2005):** A novel faces recognition approach based on Support Vector Machine and Gabor-Featured Key Points. This experiment is conducted on FERET and AT&T databases.
- **Facial Feature Selection Based on SVMs by Regularized Risk Minimization (2006):** Proposed method is an effective solution for high dimensional facial feature selection. And this is based on SVM by regularized risk minimization.
- **Face Recognition using Multiple Classifiers (2006):** Here various classification techniques like support vector machine (SVM), linear discriminate analysis (LDA) and K nearest neighbor (KNN) are studied.
- **Face Recognition Using Total Margin-Based Adaptive Fuzzy Support Vector Machines (2007):** A new classifier called total margin-based adaptive fuzzy support vector machines (TAF-SVM) is presented to deal with the several problems that occur in SVM. Results show that the TAF- SVM is superior to SVM in terms of the face-recognition accuracy.
- **Face recognition using Multi Scale PCA and support vector machine (2008):** This approach has obtained a good recognition performance through reorganizing the Gabor feature, reducing the dimension with MS-PCA and classifies SVM. This method has two limitations one is the selection of kernel parameters and other is the number of SSM classifiers used here are more.
- **Face Recognition System using SVM and Feature Extraction by PVCA and LDA Combination (2009):** First PCA is used for dimension

reduction and LDA is used for feature extraction then finally SVM is used as a Classifier. The experiments results show that PCA+LDA+SVM method has higher recognition rate than the other two methods PCA+NCC and PCA+LDA+NCC (nearest neighbor classifier). ORL database is used here. The recognition rate of this system is 94.3%.

- **Face Recognition Based on Face Gabor Image and SVM (2009):** This is an effective algorithm for face recognition using face Gabor image and Support Vector Machine (SVM). Face Gabor image is firstly derived by down sampling and concatenating the Gabor wavelets representations which are the convolution of the face image with a family of Gabor kernels, and then the 2D Principle Component Analysis (2DPCA) method is applied to the face Gabor image to extract the feature space. Finally, Support Vector Machine (SVM) is used to classify. This method uses ORL data base.
- **ISVM for Face Recognition (2010):** Similarity of human faces, unpredictable variations and aging are the crucial obstacles in face recognition to handle this large set of training samples are required which increases the complexity of the system. Since both classification and feature information are necessary for a recognition system DCT is used to lower the computational complexity and SVM for classification. Since SVM is a popular classification tool but the main disadvantage of SVM is its large memory requirement and computation time to deal with large data set. The biggest advantage of using the proposed technique is that it not only decreases the training time and updating time but also improves the classification accuracy rate up to 100 %.
- **Face recognition based on principle component analysis and support vector machine (2011):** Face recognition is done using PCA as feature extractor and SVM is as classifier. Experiments are carried out on ORL data base. Compares proposed method with PCA and nearest neighbors (PCA and NN) methods of face recognition and support vector machine on recognition rate and recognition time respectively. This method is beneficial only for small sample of training data.
- **Sparse Representation based Classification (2011):** Sparse representation-based classification (SRC) has been widely used for face recognition (FR). SRC first codes a testing sample as a sparse linear



combination of all the training samples, and then classifies the testing sample by evaluating which class leads to the minimum representation error.

- **Collaborative Representation Based Face Recognition (2011):** This technique is the special case of sparse based classification.
- **Sparse Based Representation using k nearest subspace (2014):** (SRC) has been proven to be a robust face recognition method. For this modular method, face images are partitioned into a number of blocks first and then we propose an indicator to remove the contaminated blocks and choose the nearest subspaces. Finally, SRC is used to classify the occluded test sample in the new feature space.
- **Bilinear Convolutional Neural Networks (BCNN) (2015):** It is applied for facial recognition in large public data sets with pose variability. 4 sets of features are evaluated: traditional features (eyes, nose, mouth, and eyebrows), features correlated to accessories, features correlated with hair, and features correlated with background.
- **Deep Learning Face Recognition/ Deep Convolutional Neural Networks (DNN) (2016):** Deep learning does a better job than humans at figuring out which parts of a face are important to measure. Deep learning algorithm is employed in many real time applications for example face book uses DNN for face recognition.
- **Super Pixels (2017):** This is the newest technique for face recognition. Here the patterns or the similar pixels are grouped together for feature extraction and recognition is done. This method is still under development.

2.3 Face Recognition Structure

Let us consider a picture captured from a digital camera, and we would like to know who the person is in the picture. To achieve this, we perform face recognition in three following steps.

- Face Detection
- Feature Extraction
- Face Recognition

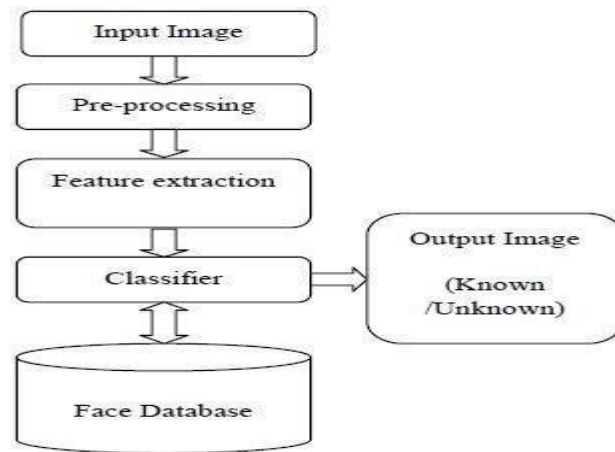


Fig. 2.1 Block Diagram of Facial Recognition

Face Detection:

The main aim of this step is to identify whether a face is present in the image. The expected outcomes are some patches containing the features of the face. Face alignment of the face is performed to know the scales and orientations of these patches. Face detection also helps in retargeting, video and image classification and also helps in concentrating the region of interest.

Face Extraction:

After the detection step patches are extracted by directly using the patches we cannot build a recognition system which is robust as image has some 1000 pixels, which is too large. Other disadvantages of directly using the patch is that each patch has different illumination, pose, camera alignments which may clutter. To overcome these problems features are extracted so that dimension reduction and noise cleaning can be done. After this step, face patches are transformed to a vector.

Face Recognition:

This is the last step of this process. In order to process this step a face database has to be built. For a person several images are taken and are stored in this database. After obtaining a face that has to be recognized face detection and feature extraction are done then the patches are compared with the faces available in the database. There are many algorithms proposed for classification.

2.4 Face Recognition Methods

Developing a very effective face recognition system is not an easy task. Several factors have to be taken in to consideration. Few of them are stated below:

- Overall speed from detection to recognition should be acceptable.
- The system should be accurate.
- It should be easy to increase the size of the subjects at any given point of time.

In the beginning people used to recognize faces using distances e.g. measurement between eyes or measurements between other nodal points. But with the changing patterns in technology face recognition can be classified into three methods:

- Holistic Matching Methods
- Feature Based (structural) Methods
- Hybrid Methods

2.4.1 Holistic Matching Methods

Here in holistic approach complete face is taken into consideration. One of the best examples of holistic approach is Eigen faces. The other methods of this approach are PCA, LDA etc.

Eigen faces:

At first data base is constructed and it is termed as the training set, next Eigen faces are made by extracting face features. Then the images are normalized and resized so that they have same size. By using the mathematical tool called Principal Component Analysis (PCA) Eigen faces are extracted. When Eigen faces are created, each image will be a vector with weights. Now, when an image is given as query, the weights of the query image is compared with the weights of the images in database. The image with the closest weight is the recognized face.



Fig. 2.2 Eigen Faces

2.4.2 Feature Based(structural) Methods

This is a local method where features such as eyes, nose, mouth are extracted and their locations and local statistics are stored into the classifier. This may be feature based or appearance based. The biggest disadvantage is the feature restoration i.e., it is difficult to retrieve an image with a large pose variation.



Fig. 2.3 FR based on features

The three different extraction methods are:

- Generic methods based on edges, lines, and curves
- Feature-template-based methods
- Structural matching methods that take into consideration geometrical Constraints on the features.

2.4.3 Hybrid Methods

Hybrid methods are those approaches that use both holistic and local feature-based methods. The main problem that affects the performance of a hybrid system is how to select the features that are to be combined and how to combine them so that their advantages are preserved and disadvantages are removed. For example, components of a hybrid system, either feature or

classifier, should be both accurate and diverse such that a complementary advantage can be feasible.

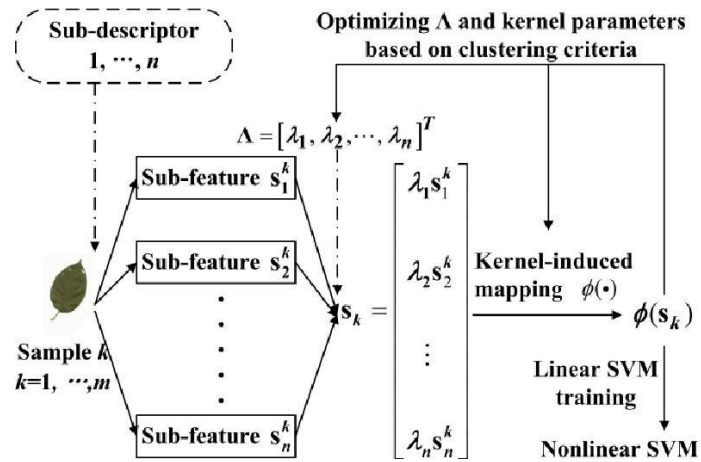


Fig 2.4 Hybrid Feature extraction using SVM

2.5 Conclusion

This chapter gives the brief description of face recognition technology techniques, applications and challenges.

Chapter Three

POSE INVARIANT FACIAL RECOGNITION USING CONVOLUTIONAL NEURAL NETWORK (CNN)

3.1 Introduction

In deeplearning, a **Convolutional Neural Network (CNN or ConvNet)** is a class of deep neural networks, most commonly applied to analyzing visual imagery. CNNs use a variation of multilayer perceptrons designed to require minimal preprocessing. They are also known as shift invariant or space invariant artificial neural networks (SIANN), based on their shared-weights architecture and translation invariance characteristics.

Convolutional networks were inspired by biological processes. In that, the connectivity pattern between neurons resembles the organization of the animal visual cortex. Individual cortical neurons respond to stimuli only in a restricted region of the visual field known as the receptive field. The receptive fields of different neurons partially overlap such that they cover the entire visual field. CNNs use relatively little pre-processing compared to other image classification algorithms. This means that the network learns the filters that in traditional algorithms were hand-engineered. This independence from prior knowledge and human effort in feature design is a major advantage. They have applications in image and video recognition, recommender systems, image classification, medical image analysis, and natural language processing.

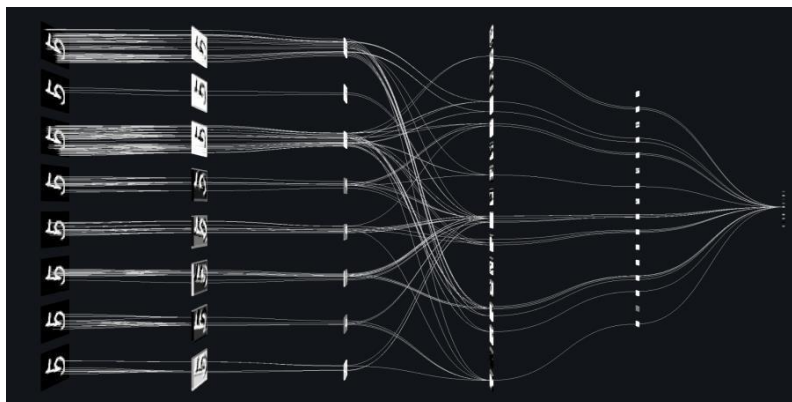


Fig. 3.1 Representation of a CNN

3.1.1 Neural Networks

Artificial Neural Networks (ANN), short Neural Networks (NN), denote a machine learning technique that has been inspired by the human brain and its capacity to perform complex tasks by means of inter-connected neurons performing each a very simple operation. Likewise, a NN is a trainable structure consisting of a set of inter-connected units, each implementing a very simple function, and together eventually performing a complex classification function or approximation task.

3.1.2 Perceptron

The most well-known type of neural unit is called Perceptron and has been introduced by Rosenblatt. Its basic structure is illustrated in Fig. 3.2. It has n inputs and one output where the output is a simple function of the sum of the input signals x weighted by w and an additional bias b . Thus,

$$y = \phi (x \cdot w + b) \quad (1)$$

Often, the bias is put inside the weight vector w such that $w_0 = b$ and the input vector x is extended correspondingly to have $x_0 = 1$. Equation (1) then becomes:

$$y = \phi (x \cdot w) \quad (2)$$

where ϕ is the Heaviside step function $\phi: \mathbb{R} \rightarrow \mathbb{R}$

$$\phi(x) = 1 \text{ if } x \geq 0; 0 \text{ else.}$$

The Perceptron thus implements a very simple two-class classifier where w is the separating hyper plane such that $w \cdot x \geq 0$ for examples from one class and $w \cdot x < 0$.

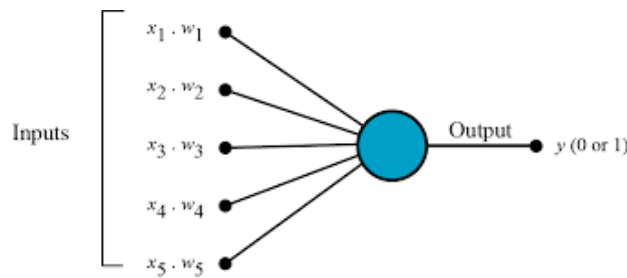


Fig. 3.2 Perceptron

3.1.3 Multi-Layer Perceptron

Multi-Layer Perceptrons (MLP) are capable of approximating arbitrarily complex decision functions. With the advent of a practicable training algorithm in the 1980's, the so-called Backpropagation algorithm, they became the most widely used form of NNs. There is an input layer, one or more hidden layer(s) and an output layer of neurons, where each neuron except the input neurons implements a perceptron as described in the previous section. Moreover, the neurons of one layer are only connected to the following layer. We call this type of network: feed-forward network, i.e. the activation of the neurons is propagated layer-wise from the input to the output layer. And there is a connection from each neuron to every neuron in the following layer. Further, the neurons' activation function has to be differentiable in order to adjust the weights by the Backpropagation algorithm. Commonly used activation functions are for example:

$$\varphi(x) = x$$

linear

$$\varphi(x) = 1/(1 + e^{-cx}) ; (c > 0)$$

sigmoid

$$\varphi(x) = (1 - e^{-x})/(1 + e^{-x})$$

hyperbolic tangent

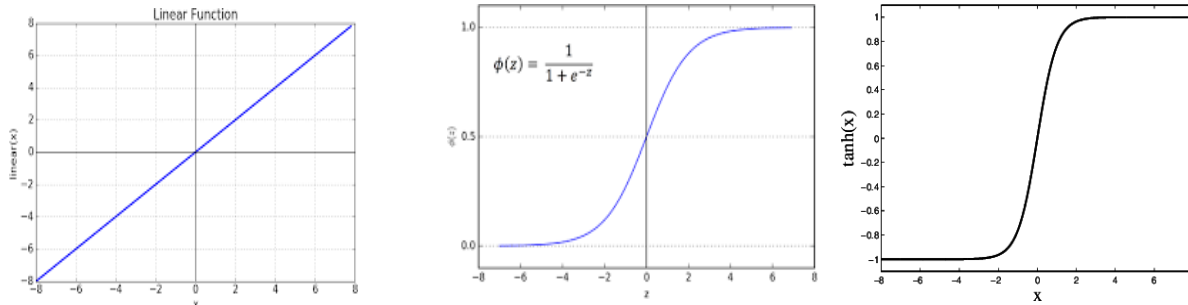


Fig. 3.3 Different types of activation functions

3.1.4 Training Neural Networks

In general, the parameters of a NN, i.e. the weights and biases, are learned using a training data set. However, as the space of possible weights can be very large and of high dimension the analytical determination of these weights might be very difficult or even infeasible. For this reason, an iterative approach is adopted in most cases. There are two principal training modes which determine the way the weights are updated:

Online training: After presentation of each training example, the error is calculated, and the weights are updated accordingly.

Offline training: The whole training set is propagated through the NN, and the respective errors are accumulated. Finally, the weights are updated using the accumulated error. This is also called batch training.

Many different NN training algorithms have been published in the literature. Some work only in online, some only in offline mode, and some can be executed in both ways. Which algorithm is best for a given problem depends on the NN architecture, the nature and cardinality of the training data set and the type of function to learn. Therefore, there is no basic rule for the choice of the training algorithm. In the following, we will focus on the Backpropagation algorithm since it is the most common and maybe most universal training algorithm for feed-forward NNs, especially for MLPs. Some alternative methods will also be presented at the end of this section.

3.1.5. Backpropagation Algorithm

In the context of NNs, the Backpropagation (BP) algorithm has initially been presented by Rumelhart et al. It is a supervised learning algorithm defining an error function E and applying the gradient descent technique in the weight space in order to minimize E . The combination of weights leading to a minimum of E is considered to be a solution of the learning problem. Note that the BP algorithm does not guarantee to find a global minimum which is an inherent problem of gradient descent optimization. However, we will discuss some approaches to overcome this problem in the following section. In order to calculate the gradient of E , at each iteration, the error function has to be continuous and differentiable. Thus, the activation function of each individual perceptron must also have this property. Mostly, a sigmoid or hyperbolic tangent activation function is employed, depending on the range of desired output values, i.e. $[0, 1]$ or $[-1, +1]$. Note that BP can be performed in online or offline mode, i.e. E represents either the error of one training example or the sum of errors produced by all training examples. In the following, we will explain the standard online BP algorithm, also known as Stochastic Gradient Descent, applied to MLPs. There are two phases of the algorithm:

- The forward pass, where a training example is presented to the network and the activations of the respective neurons is propagated layer by layer until the output neurons.
- The backward pass, where at each neuron the respective error is calculated starting from the output neurons and, layer by layer, propagating the error back until the input neurons.

Now, let us define the error function as:

$$E = 1/2 \sum_{p=1}^P ||op - tp||^2 \quad (6)$$

Where P is the number of training examples, op are the output values produced by the NN having presented example p, and tp are the respective target values. The goal is to minimize E by adjusting the weights of the NN. With online learning we calculate the error and try to minimize it after presenting each training example. Thus,

$$E_p = 1/2 ||op - tp||^2 = 1/2 \sum_{k=1}^K (opk - tpk)^2 \quad (7)$$

Where K is the number of output units. When minimizing this function by gradient descent, we calculate the steepest descent of the error surface in the weight space, i.e. the opposite direction of the gradient $\nabla E_p = (\partial E_p / \partial w_1, \dots, \partial E_p / \partial w_k)$. In order to ensure convergence, the weights are only updated by a proportion of the gradient. Thus,

$$\Delta w_k = -\lambda \partial E_p / \partial w_k \quad (8)$$

Now, let us define

$e_{pk} = o_{pk} - t_{pk}$	the error of pattern p at output neuron k
w_{kj}^o	the weight from hidden neuron j to output neuron k
w_{ji}^h	the weight from input neuron i to hidden neuron j
$z_{pj} = \phi \left(\sum_i w_{ji} x_{pi} \right)$	the output of hidden neuron j
$V_{pk} = \sum_j w_{kj}^o z_{pj}$	the weighted sum of all inputs z_{pj} of output neuron k
$V_{pj} = \sum_i w_{ji}^h x_{pi}$	the weighted sum of all inputs x_{pi} of hidden neuron j .

The subscript p always refers to pattern p, i refers to input neuron i, j to hidden neuron j and k to output neuron k. By applying the chain rule to equation (8), for one particular weight w_{kj}^o and training example p, we have:

$$\begin{aligned}
 \Delta w_{kj}^o &= -\lambda \frac{\partial E_p}{\partial w_{kj}^o} \\
 &= -\lambda \frac{\partial E_p}{\partial e_{pk}} \frac{\partial e_{pk}}{\partial o_{pk}} \frac{\partial o_{pk}}{\partial V_{pk}} \frac{\partial V_{pk}}{\partial w_{jk}^o} \\
 &= -\lambda e_{pk} \phi'(V_{pk}) z_{pj} \\
 &= -\lambda \delta_{pk} z_{pj} \quad ,
 \end{aligned}$$

$$\delta_{pk} = e_{pk} \phi'(V_{pk})$$

This holds for output neurons. For the hidden neurons the respective equations are slightly different:

$$\begin{aligned}
 \Delta w_{ji}^h &= -\lambda \frac{\partial E_p}{\partial w_{ji}^h} \\
 &= -\lambda \frac{\partial E_p}{\partial z_{pj}} \frac{\partial z_{pj}}{\partial V_{pj}} \frac{\partial V_{pj}}{\partial w_{ji}^h} \\
 &= -\lambda \left(\sum_{k=1}^K e_{pk} \frac{\partial e_{pk}}{\partial z_{pj}} \right) \phi'(V_{pj}) x_{pi} \\
 &= -\lambda \left(\sum_{k=1}^K e_{pk} \frac{\partial e_{pk}}{\partial o_{pk}} \frac{\partial o_{pk}}{\partial V_{pk}} \frac{\partial V_{pk}}{\partial z_{pj}} \right) \phi'(V_{pj}) x_{pi} \\
 &= -\lambda \left(\sum_{k=1}^K e_{pk} \phi'(V_{pk}) w_{kj}^o \right) \phi'(V_{pj}) x_{pi} \\
 &= -\lambda \delta_{pj} x_{pi}
 \end{aligned}$$

where

$$\delta_{pj} = \left(\sum_{k=1}^K \delta_{pk} w_{kj}^o \right) \phi'(V_{pj}) \quad (10)$$

The local gradient for hidden neuron j ($j = 1 \dots J$). Algorithm 3 summarizes the standard online Backpropagation algorithm. The respective variables are noted as functions of iteration n, e.g. $w_{kj}(n)$. φ is a small constant that determines the convergence criteria and max_iter is the maximum number of iterations.

3.1.6. Advantages of using Neural Networks

- Multi-layer feed-forward Neural Networks (NN) have shown to be a very powerful machine learning technique as they can be trained to approximate complex non-linear functions from high-dimensional input examples.
- Classically, standard Multi-Layer Perceptrons (MLP) has been utilized in pattern recognition systems to classify signatures coming from a separate feature extraction algorithm operating on the input data. However, the manual choice of the feature extraction algorithm and the features to classify is often empirical and therefore sub-optimal. Thus, a possible solution would be to directly apply the NN on the “raw” input data and let the training algorithm, e.g. Backpropagation, find the best feature extractors by adjusting the weights accordingly.

3.1.7. Disadvantages of using Neural Networks

- The problem with this approach is that when the input dimension is high, as in images, the number of connections, thus the number of free parameters are also high because each hidden unit would be fully connected to the input layer.
- This number may be in the order of several 10,000 or rather several 100,000 according to the application.
- The number of training examples, however, might be relatively small compared to the pattern dimension, which means that the NN would have a too high complexity and, thus, would tend to over fit the data.
- Its input layer has a fixed size and the input patterns have to be presented well aligned and/or normalized to this input window, which, in practice, is a rather complicated task. Thus, there is no built-in invariance w.r.t. small translations and local distortions.

- Finally, fully-connected NN architectures do not take into account correlations of neighboring input data. However, in pattern recognition problems there is generally a high amount of local correlation.

Thus, to avoid all the problems faces in Neural Networks as mentioned above, we go for a more sophisticated algorithm called Convolutional Neural Networks which we will discuss in the further sections. CNNs are an approach that tries to alleviate the above-mentioned problems. That is, they automatically learn local feature extractors, they are invariant to small translations and distortions in the input pattern, and they implement the principle of weight sharing which drastically reduces the number of free parameters and thus increases their generalization capacity compared to NN architectures without this property.

3.2. Convolutional Neural Networks

3.2.1. Design

A convolutional neural network consists of an input and an output layer, as well as multiple hidden layers. The hidden layers of a CNN typically consist of convolutional layers, RELU layer i.e. activation function, pooling layers, fully connected layers and normalization layers.

Description of the process as a convolution in neural networks is by convention. Mathematically it is a cross-correlation rather than a convolution (although cross-correlation is a related operation). This only has significance for the indices in the matrix, and thus which weights are placed at which index.

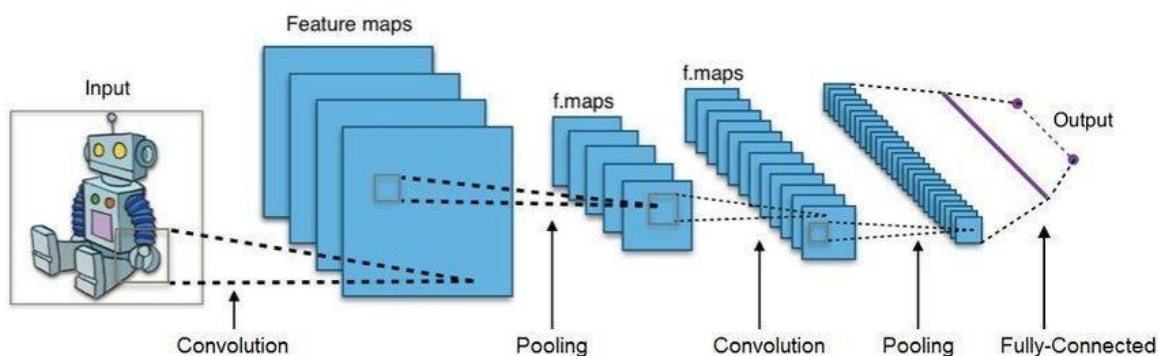


Fig. 3.4 Typical CNN Architecture

Convolutional:

Convolutional layers apply a convolution operation to the input, passing the result to the next layer. The convolution emulates the response of an individual neuron to visual stimuli.

Although fully connected feed forward neural networks can be used to learn features as well as classify data, it is not practical to apply this architecture to images. A very high number of neurons would be necessary, even in shallow (opposite of deep) architecture, due to the very large input sizes associated with images, where each pixel is a relevant variable.

For instance, a fully connected layer for a (small) image of size 100 x 100 has 10000 weights for *each* neuron in the second layer. The convolution operation brings a solution to this problem as it reduces the number of free parameters, allowing the network to be deeper with fewer parameters. For instance, regardless of image size, tiling regions of size 5 x 5, each with the same shared weights, requires only 25 learnable parameters. In this way, it resolves the vanishing or exploding gradients problem in training traditional multi-layer neural networks with many layers by using backpropagation.

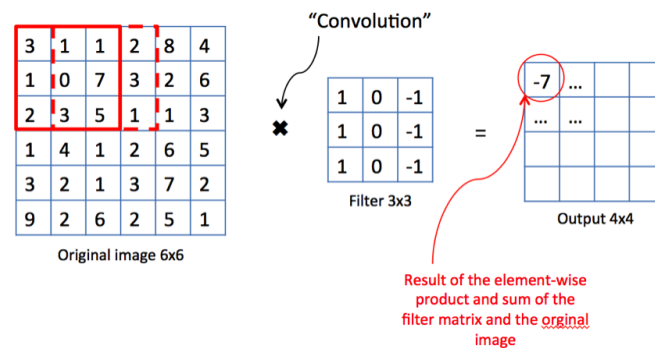


Fig. 3.5 Convolution

Pooling:

Convolution networks may include local or global pooling layers, which combine the outputs of neuron clusters at one layer into a single neuron in the next layer. For example, *max pooling* uses the maximum value from each of a cluster of neurons at the prior layer. Another example is *average pooling*, which uses the average value from each of a cluster of neurons at the prior layer.

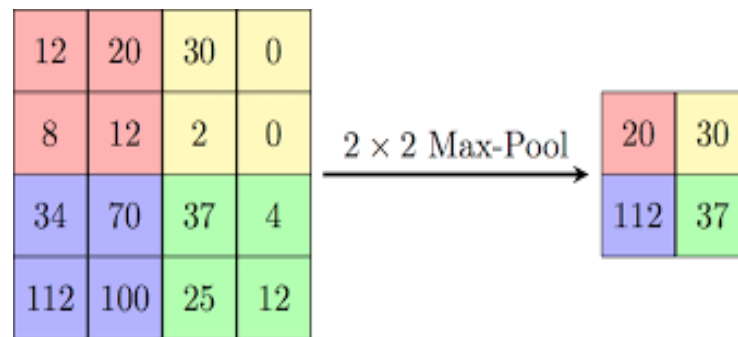


Fig. 3.6 Pooling

Fully Connected:

Fully connected layers connect every neuron in one layer to every neuron in another layer. It is in principle the same as the traditional perceptron neural network (MLP). The flattened matrix goes through a fully connected layer to classify the images.

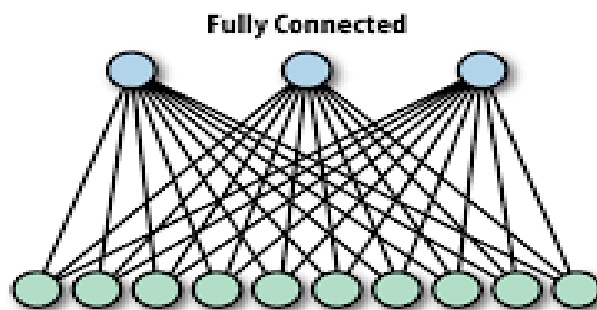


Fig. 3.7 Fully Connected

Receptive Field:

In neural networks, each neuron receives input from some number of locations in the previous layer. In a fully connected layer, each neuron receives input from *every* element of the previous layer. In a convolutional layer, neurons receive input from only a restricted subarea of the previous layer. Typically, the subarea is of a square shape (e.g., size 5 by 5). The input area of a



neuron is

called its *receptive field*. So, in a fully connected layer, the receptive field is the entire previous layer. In a convolutional layer, the receptive area is smaller than the entire previous layer.

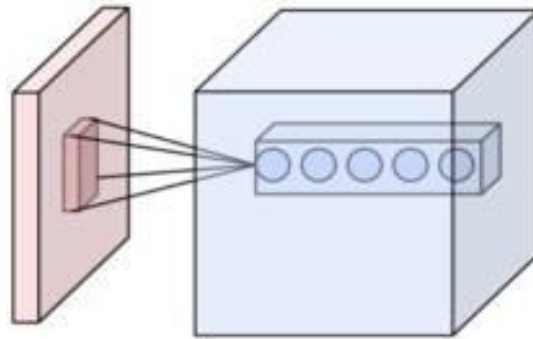


Fig. 3.8 Neurons(blue) connected to their Receptive Field(red)

Weights:

The function that is applied to the input values is specified by a vector of weights and a bias (typically real numbers). Learning in a neural network progresses by making incremental adjustments to the biases and weights. The vector of weights and the bias are called a *filter* and represent some feature of the input (e.g., a particular shape). A distinguishing feature of CNNs is that many neurons share the same filter. This reduces memory footprint because a single bias and a single vector of weights is used across all receptive fields sharing that filter, rather than each receptive field having its own bias and vector of weights.

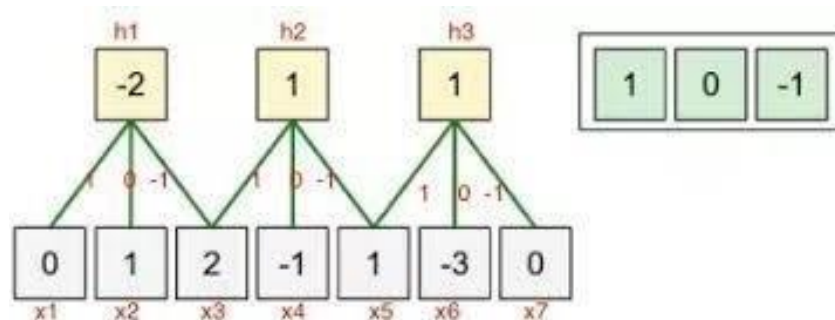


Fig. 3.9 Weight sharing in CNNs

3.2.2. Layers

3.2.2.1. Convolutional Layer

We now discuss the details of the neuron connectivity, their arrangement in space, and their parameter sharing scheme.

Local Connectivity:

When dealing with high-dimensional inputs such as images, as we saw above it is impractical to connect neurons to all neurons in the previous volume. Instead, we will connect each neuron to only a local region of the input volume. The spatial extent of this connectivity is a hyper parameter called the receptive field of the neuron (equivalently this is the filter size).

Example: suppose that the input volume has size $[32 \times 32 \times 3]$, (e.g. an RGB CIFAR-10 image). If the receptive field (or the filter size) is 5×5 , then each neuron in the Conv Layer will have weights to a $[5 \times 5 \times 3]$ region in the input volume, for a total of $5 \times 5 \times 3 = 75$ weights (and +1 bias parameter).

Spatial Arrangement:

Three hyperparameters control the size of the output volume: the **depth**, **stride** and **zero-padding**.

- **Depth:** It corresponds to the number of filters we would like to use, each learning to look for something different in the input.
- **Stride:** We must specify the **stride** with which we slide the filter. When the stride is 1 then we move the filters one pixel at a time. When the stride is 2 then the filters jump 2 pixels at a time as we slide them around. This will produce smaller output volumes spatially.
- **Zero-padding:** Sometimes it will be convenient to pad the input volume with zeros around the border. The size of this **zero-padding** is a hyperparameter. The nice feature of zero padding is that it will allow us to control the spatial size of the output volumes.

Parameter Sharing:

A parameter sharing scheme is used in convolutional layers to control the number of free parameters. It relies on one reasonable assumption: if a patch feature is useful to compute at some spatial position, then it should also be useful to compute at other positions. In other words, denoting a single 2- dimensional slice of depth as a **depth slice**, we constrain the neurons in each depth slice to use the same weights and bias.

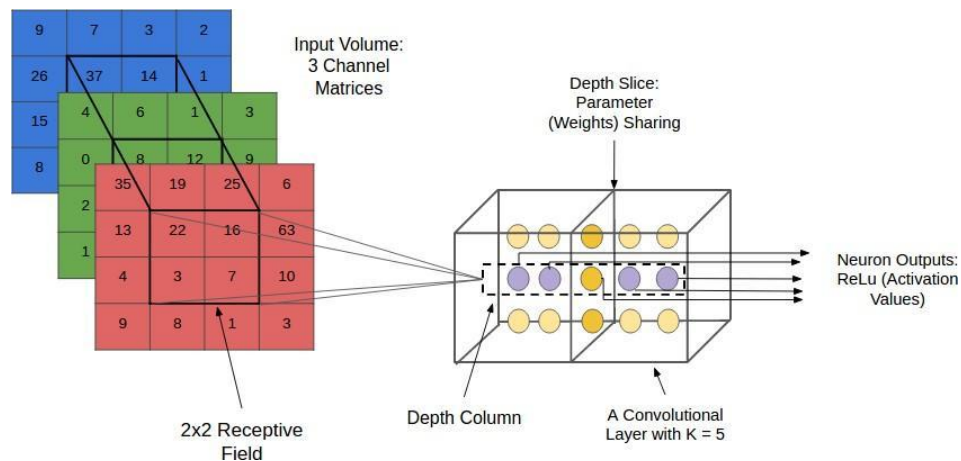


Fig. 3.10 Convolutional Layer

3.2.2.2. Pooling Layer

It is common to periodically insert a Pooling layer in-between successive Conv layers in a ConvNet architecture. Its function is to progressively reduce the spatial size of the representation to reduce the number of parameters and computation in the network, and hence to also control overfitting.

The Pooling Layer operates independently on every depth slice of the input and resizes it spatially, using the MAX operation. The most common form is a pooling layer with filters of size 2x2 applied with a stride of 2 down samples every depth slice in the input by 2 along both width and height, discarding 75% of the activations. Every MAX operation would in this case be taking a max over 4 numbers.

The depth dimension remains unchanged. More generally, the pooling layer:

- Accepts a volume of size $W1 \times H1 \times D1$
- Requires two hyperparameters:

- o their spatial extent FF,
- o the stride SS,
- Produces a volume of size $W2 \times H2 \times D2$ where:
 - o $W2 = (W1 - F) / S + 1$
 - o $H2 = (H1 - F) / S + 1$
 - o $D2 = D1$
- Introduces zero parameters since it computes a fixed function of the input,
- For pooling layers, it is not common to pad the input using zero-padding.

General Pooling: In addition to max pooling, the pooling units can also perform other functions, such as average pooling or even L2-norm pooling. Average pooling was often used historically but has recently fallen out of favor compared to the max pooling operation, which has been shown to work better in practice.

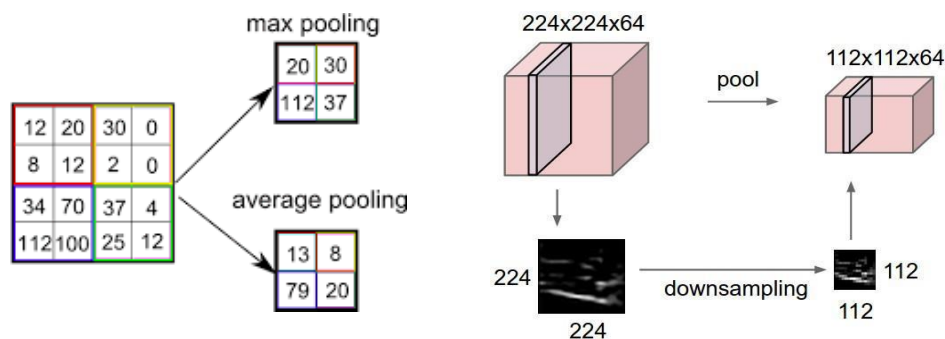


Fig. 3.11 Pooling Layers

In our project, we have used the max pooling layer.

3.2.2.3. Activation Functions

It's just a function that you use to get the output of node. It is also known as **Transfer Function**. It maps the resulting values in between 0 to 1 or -1 to 1 etc. (depending upon the function).

The Activation Functions can be basically divided into 2 types-

- Linear Activation Function

- Non-linear Activation Functions

Linear or Identity Activation Function:

As you can see the function is a line or linear. Therefore, the output of the functions will not be confined between any range.

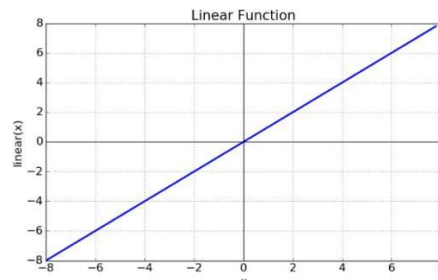


Fig. 3.12 Linear Activation Function

Non-Linear Activation Function:

The Nonlinear Activation Functions are the most used activation functions. Nonlinearity helps to makes the graph look something like this;

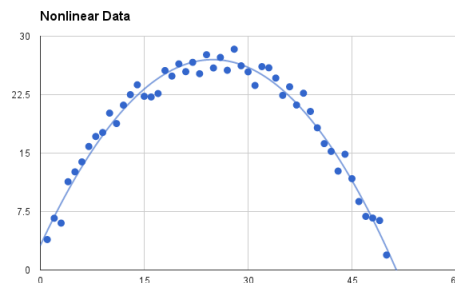


Fig. 3.13 Non-Linear Activation Function

The Nonlinear Activation Functions are mainly divided on the basis of their range or curves:

1) Sigmoid or Logistic Activation Function

The Sigmoid Function curve looks like a S-shape. The main reason why we use sigmoid function is because it exists between (0 to 1). Therefore, it is especially used for models where we have to predict the probability as an output. Since probability of anything exists only between the range of 0 and 1, sigmoid is the right choice. The function is differentiable. That means, we can find the slope of the sigmoid curve at any two points. The function is monotonic but function's derivative is not. The logistic sigmoid



function can cause a neural network to get stuck at the training time.

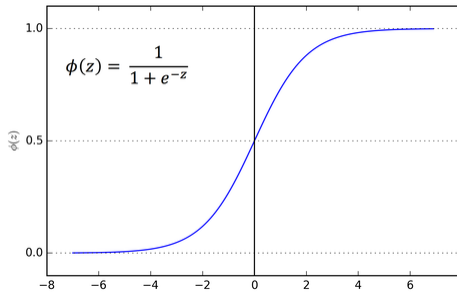


Fig. 3.14 Sigmoid Activation Function

2) Tanh or hyperbolic tangent Activation Function

Tanh is also like logistic sigmoid but better. The range of the tanh function is from (-1 to 1). tanh is also sigmoidal (s - shaped). The advantage is that the negative inputs will be mapped strongly negative and the zero inputs will be mapped near zero in the tanh graph. The function is differentiable. The function is monotonic while its derivative is not monotonic. The tanh function is mainly used classification between two classes.

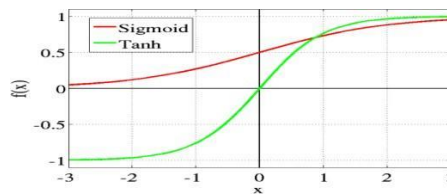


Fig. 3.15 Tanh Activation Function

3) ReLU (Rectified Linear Unit) Activation Function

The ReLU is the most used activation function in the world right now since it is used in almost all the convolutional neural networks which applies the non-saturating activation function $\mathbf{f(x)=max(0,x)}$. It effectively removes negative values from an activation map by setting them to zero.

The function and its derivative **both are monotonic**.

But the issue is that all the negative values become zero immediately which decreases the ability of the model to fit or train from the data properly. That means any negative input given to the ReLU activation function turns the value into zero immediately in the graph,



which in turns affects the resulting

graph by not mapping the negative values appropriately. We have used the ReLU activation function in our project.

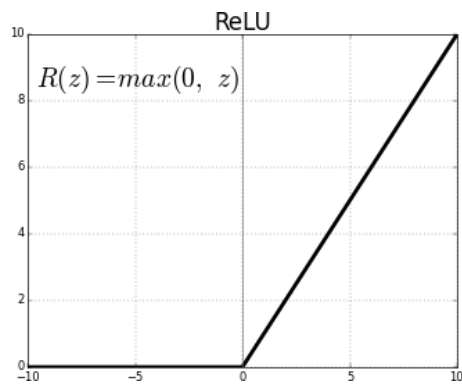


Fig. 3.16 ReLU Activation Function

4) Leaky ReLU:

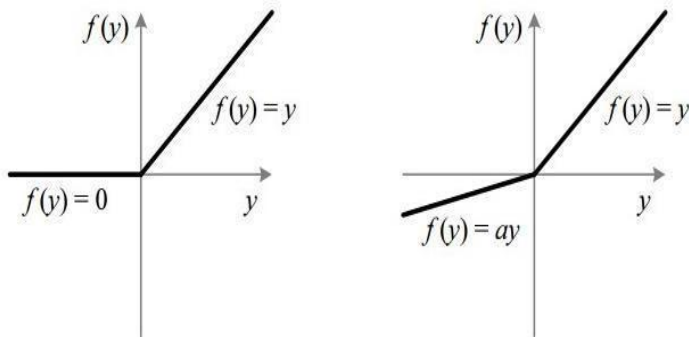











Fig. 3.17 ReLU v/s Leaky ReLU

The leak helps to increase the range of the ReLU function. Usually, the value of **a** is 0.01 or so, when **a is not 0.01** then it is called **Randomized ReLU**. Therefore, the **range** of the Leaky ReLU is (-infinity to infinity). Both Leaky and Randomized ReLU functions are monotonic in nature. Also, their derivatives also monotonic in nature.

Table 3.1 Different Activation Functions

Name	Plot	Equation	Derivative
Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
Rectified Linear Unit (ReLU)		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Parametric Rectified Linear Unit (PReLU) [2]		$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Exponential Linear Unit (ELU) [3]		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
SoftPlus		$f(x) = \log_e(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

3.2.3. Classification

In **machine learning** and statistics, **classification** is the problem of identifying to which of a set of categories (sub-populations) a new observation belongs, on the basis of a training set of data containing observations (or instances) whose category membership is known.

3.2.3.1. Flatten

Now that we have converted our input image into a suitable form for our Multi- Level Perceptron, we shall flatten the image into a column vector. This phenomenon is called 'Flattening'. The flattened output is fed to a feed-forward neural network and backpropagation applied to every iteration of training. Over a series of epochs, the model is able to distinguish between dominating and certain low-level features in images and classify them using the **Logistic Regression Classification** technique.

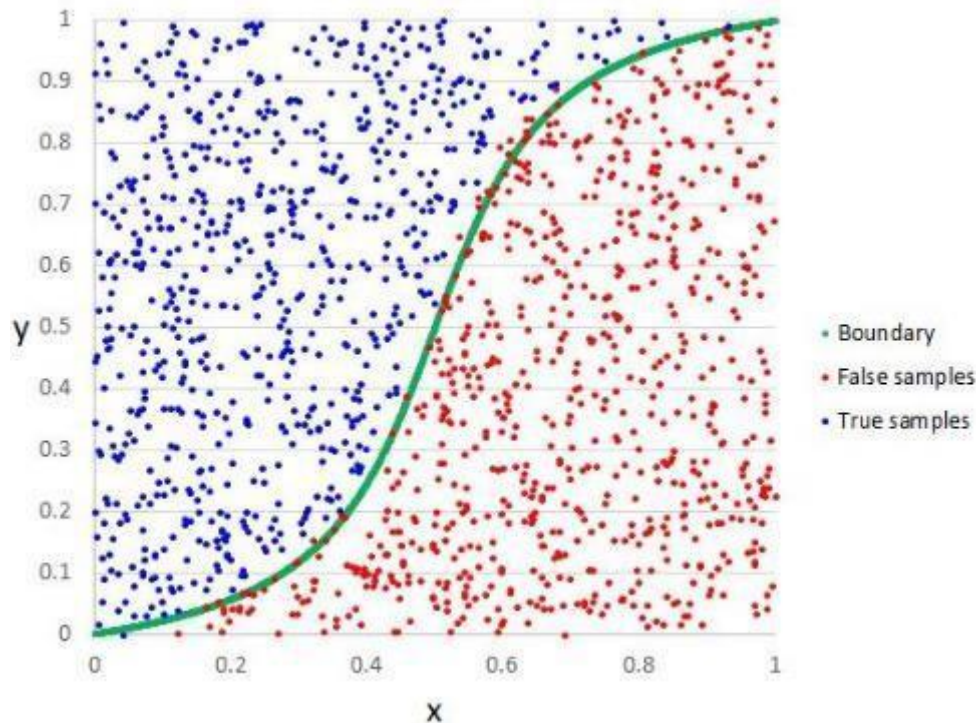


Fig. 3.18 Classification Using Log. Reg. Function

3.2.3.2. *Fully-Connected*

Finally, after several convolutional and max pooling layers, the high-level reasoning in the neural network is done via fully connected layers. Neurons in a fully connected layer have connections to all activations in the previous layer, as seen in regular (non-convolutional) artificial neural networks as shown in Fig. 3.7 Their activations can thus be computed as an affine transformation, with matrix multiplication followed by a bias offset (vector addition of a learned or fixed bias term).

3.2.3.3. *Logistic Regression Classifier*

Logistic regression is named for the function used at the core of the method, the logistic function.

The logistic function, also called the sigmoid function was developed by statisticians to describe properties of population growth in ecology, rising quickly and maxing out at the carrying capacity of the environment. It's an S- shaped curve that can take any real-valued number and map it into a value between 0 and 1, but never exactly at those limits.

Logistic regression uses an equation as the representation, very much like linear regression.



Input values (x) are combined linearly using weights or coefficient values (referred to as the Greek capital letter β) to predict an output value (y). A key difference from linear regression is that the output value being modeled is a binary value (0 or 1) rather than a numeric value.

Below is an example logistic regression equation:

$$y = e^{(b_0 + b_1 * x)} / (1 + e^{(b_0 + b_1 * x)})$$

Where y is the predicted output, b_0 is the bias or intercept term and b_1 is the coefficient for the single input value (x). Each column in your input data has an associated b coefficient (a constant real value) that must be learned from your training data.

The actual representation of the model that you would store in memory or in a file is the coefficients in the equation.

Logistic regression models the probability of the default class (e.g. the first class).

For example, if we are modeling people's sex as male or female from their height, then the first class could be male and the logistic regression model could be written as the probability of male given a person's height, or more formally:

$$P(\text{sex}=\text{male}|\text{height})$$

Written another way, we are modeling the probability that an input (X) belongs to the default class ($Y=1$), we can write this formally as:

$$P(X) = P(Y=1|X)$$

Note that the probability prediction must be transformed into a binary value (0 or 1) in order to actually make a probability prediction. More on this later when we talk about making predictions.

Logistic regression is a linear method, but the predictions are transformed using the logistic function. The impact of this is that we can no longer understand the predictions as a linear combination of the inputs as we can with linear regression, for example, continuing on from above, the model can be stated as:

$$p(X) = e^{(b_0 + b_1 * X)} / (1 + e^{(b_0 + b_1 * X)})$$

we can turn around the above equation as follows (remember we can remove the e from one side by adding a natural logarithm (ln) to the other):

$$\ln(p(X) / 1 - p(X)) = b_0 + b_1 * X$$

This is useful because we can see that the calculation of the output on the right is linear again (just like linear regression), and the input on the left is a log of the probability of the default class.

This ratio on the left is called the odds of the default class (it's historical that we use odds, for example, odds are used in horse racing rather than probabilities). Odds are calculated as a ratio of the probability of the event divided by the probability of not the event, e.g. $0.8/(1-0.8)$ which has the odds of 4. So, we could instead write:

$$\ln(\text{odds}) = b_0 + b_1 * X$$

Because the odds are log transformed, we call this left-hand side the log-odds or the probit. It is possible to use other types of functions for the transform (which is out of scope, but as such it is common to refer to the transform that relates the linear regression equation to the probabilities as the link function, e.g. the probit link function.

We can move the exponent back to the right and write it as:

$$\text{odds} = e^{(b_0 + b_1 * X)}$$

All of this helps us understand that indeed the model is still a linear combination of the inputs, but that this linear combination relates to the log- odds of the default class.

In our project we have used the Linear + Log-Regression Classifier.

3.2.3.4. *Softmax Classifier*

The softmax activation is normally applied to the very last layer in a neural net, instead of using ReLU, sigmoid, tanh, or another activation function. The reason why softmax is useful is because it converts the output of the last layer in your neural network into what is essentially a probability distribution. If you look at the origins of the cross-entropy loss function in information theory, you

will know that it "expects" two probability distributions as input. That's why softmax output with cross entropy loss is very common.

Just to reiterate; softmax is typically viewed as an activation function, like sigmoid or ReLU. Softmax is NOT a loss function, but is used to make the output of a neural net more "compatible" with the cross entropy or negative log likelihood loss functions.

The Softmax classifier is given by:

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$$

Which allows us to interpret the outputs as probabilities, while cross-entropy loss is what we use to measure the error at a softmax layer, and is given by

$$L(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N H(p_n, q_n) = -\frac{1}{N} \sum_{n=1}^N \left[y_n \log \hat{y}_n + (1 - y_n) \log(1 - \hat{y}_n) \right],$$

3.3. Pose-Invariant Face Recognition:

According to C. Ding and D. Tao-A comprehensive survey on pose-invariant face recognition, existing PIFR methods can be classified into four categories including: 1) multi-view subspace learning, 2) pose-invariant feature extraction, face synthesis, and 3) a hybrid approach of the above three. Our work belongs to the second category of extracting pose-invariant features. Some previous work in this category treat each pose separately by learning different models for face images with different poses.

Since pose variation is the most challenging one among other non-identity variations, and the proposed m-CNN already classifies all images into different pose groups, we propose to apply divide-and-conquer to CNN learning.

Specifically, we develop a novel pose-directed multi-task CNN (pCNN) where the pose labels can categorize the training data into three different pose groups, direct them through different routes in the network to learn pose-specific identity features in addition to the generic identity features.

Similarly, the loss weights for extracting these two types of features are learnt dynamically in the CNN framework. During the testing stage, we propose a stochastic routing scheme to fuse the generic identity features and the pose- specific identity features for face recognition that is more robust to pose estimation errors.

This work utilizes all data in FERET, ATT, Webcam and a self-curated dataset, i.e., faces with the full range of pose variations, as the main experimental dataset — ideal for studying MTL for Pose-invariant face recognition. Since the ground truth label of the side task is unavailable, we use the estimated poses as labels for training. In summary, we make four contributions:

- We explore how and why pose estimations can help face recognition;
- We propose a dynamic-weighting scheme to learn the loss weights for different tasks automatically in the CNN framework.
- We develop a pose-directed multi-task CNN to handle pose variation.
- We also developed an algorithm in which we directly map faces using PCA and show the frontal face of the subject irrespective of the pose.

3.4. CNN Block Diagram

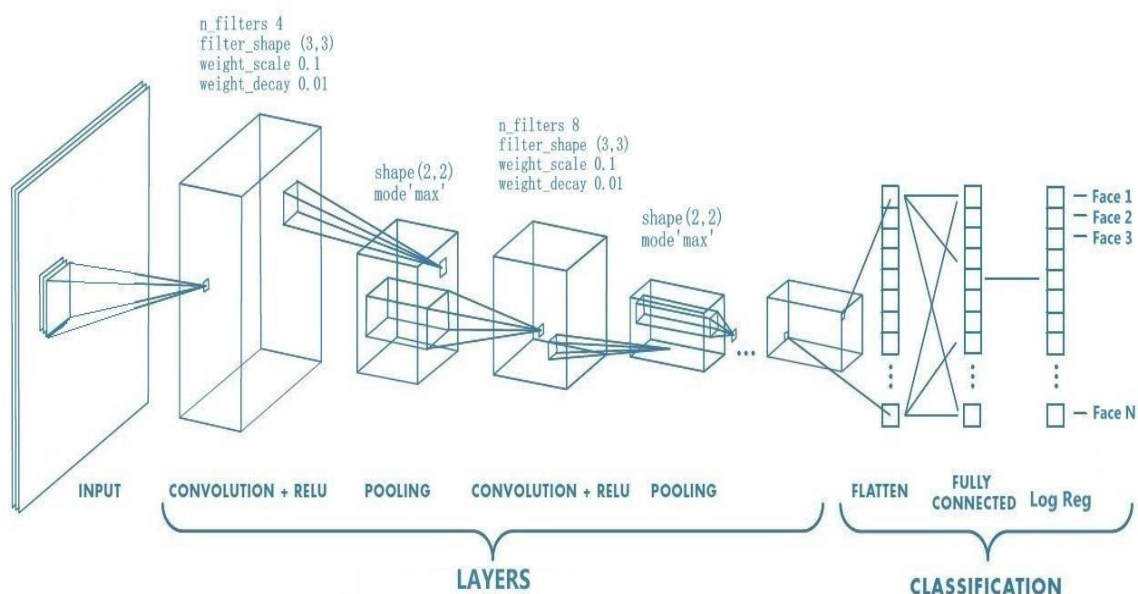


Fig. 3.19 CNN Block Diagram

3.4.1. Specification of the Architecture

- In the first stage of architecture we used convolution layer with 4 filters and shape 3 x 3 so that we can grab minute details in face which will be helpful for the classifier to segment different faces.
- In RELU we generally suppress the negative values which are generated in convolution layer because practically we don't have negative pixel values.
- In the second stage we used 2 x 2 pooling layer with max mode so that we reduce the image size and only obtain high frequency components.
- Same layers can be added again and again with change in parameters for better performance.
- In flattening, we convert 2D images to vectors which represent the features of a person.
- Finally, we use Logistic Regression Layer for classification.

3.4.2. Advantages of CNN

- CNN (Convolutional Neural Networks) is more efficient in extracting the minute features from picture which is very much important in classifying human faces.
- We reduce the size of intermediate images required for processing using pooling layers.
- Convolution simplifies computation to a great extent without losing the essence of the data.
- Minimize computation compared to a regular neural network.
- We can use pre trained models using transfer learning options thus reducing time for training the model.
- We can add multiple layers easily with little modification for better classification of images.

3.4.3. Limitations of CNN

- Large number of training data and annotations are needed, which may not be practical in some problems.
- If you don't have a good GPU they are quite slow to train (for complex tasks).

- A convolution is a significantly slower operation than, say maxpool, both forward and backward. If the network is pretty deep, each training step is going to take much longer.
- As the data increases, time complexity and memory footprint of the architecture increases.
- Reconstruction of intermediate images is not possible.

3.4.4. Case Studies

There are several architectures in the field of Convolutional Networks that have a name. The most common are:

- **LeNet.** The first successful applications of Convolutional Networks were developed by Yann LeCun in 1990's. Of these, the best known is the LeNet architecture that was used to read zip codes, digits, etc.
- **AlexNet.** The first work that popularized Convolutional Networks in Computer Vision was the AlexNet, developed by Alex Krizhevsky, Ilya Sutskever and Geoff Hinton. The AlexNet was submitted to the ImageNet ILSVRC challenge in 2012 and significantly outperformed the second runner-up (top 5 error of 16% compared to runner-up with 26% error). The Network had a very similar architecture to LeNet, but was deeper, bigger, and featured Convolutional Layers stacked on top of each other (previously it was common to only have a single CONV layer always immediately followed by a POOL layer).
- **ZF Net.** The ILSVRC 2013 winner was a Convolutional Network from Matthew Zeiler and Rob Fergus. It came to be known as the ZFNet (short for Zeiler & Fergus Net). It was an improvement on AlexNet by tweaking the architecture hyperparameters, in particular by expanding the size of the middle convolutional layers and making the stride and filter size on the first layer smaller.
- **GoogLeNet.** The ILSVRC 2014 winner was a Convolutional Network from Szegedy et al. from Google. Its main contribution was the development of an *Inception Module* that dramatically reduced the number of parameters in the network (4M, compared to AlexNet with 60M). Additionally, this paper uses Average Pooling instead of Fully Connected layers at the top of the



ConvNet, eliminating a large number of parameters that do not seem to matter much. There are also several follow-up versions to the GoogLeNet, most recently Inception-v4.

3.5. Conclusion

In this chapter we study about the evolution of neural networks, the architecture of Convolutional Neural Networks and the block diagram.

Chapter Four

TOOLS & LIBRARIES USED

4.1. Introduction

In this project we have used a wide range of tools, applications and libraries. Most of these tools are open source tools and contains in-built functions for easy development of programs and applications.

4.2. Programming Languages

We mostly used two languages in the development of architecture and testing various methods related to pose. They are:

4.2.1. Python 3.7

Python is a popular platform used for research and development of production systems. It is a vast language with number of modules, packages and libraries that provide multiple ways of achieving a task.

Python and its libraries like NumPy, SciPy, Scikit-Learn, Matplotlib are used in data science and data analysis. They are also extensively used for creating scalable machine learning algorithms. Python implements popular machine learning techniques such as Classification, Regression, Recommendation, and Clustering.

To understand machine learning, you need to have basic knowledge of Python programming. In addition, there are a number of libraries and packages generally used in performing various machine learning tasks as listed below:

numpy- is used for its N-dimensional array objects.

pandas – is a data analysis library that includes data frames.

matplotlib – is a 2D plotting library for creating graphs and plots.

scikit-learn - the algorithms used for data analysis and data mining tasks.

Python offers ready-made framework for performing data mining tasks on large volumes of



data effectively in lesser time. It includes several implementations



achieved through algorithms such as linear regression, logistic regression, Naïve Bayes, k-means, K nearest neighbor, and Random Forest.

4.3. Applications Used

We used two applications during the course of our project:

4.3.1. Irfan View

This application is one stop solution for all problems related to format of the images and basic processing of the images. Features of this application are:

32 and 64 bit version

Many supported file formats (.jpg, .png, .pgm, raw, gif etc)

Thumbnail/preview option

Paint option - to draw lines, circles, arrows, straighten image etc.

Lossless JPG rotation, crop and EXIF date change (also in batch mode)

Slideshow (save slideshow as EXE/SCR or burn it to CD)

Support for Adobe Photoshop Filters

Batch conversion (with advanced image processing of all files)

Change color depth

Scan (batch scan) support

4.4. Conclusion

In this chapter we discussed about the various tools and their functions used in the developments of our architecture.



Chapter Five

RESULTS & PROJECT ANALYSIS

5.1. Databases Used

CNN technique is tested on the ATT, FERET and WEBCAM.

5.1.1. *ATT Database*

ORL face database is composed of 400 images of size 112 x 92. There are 40 persons, 10 images per each person but we only used 10 people's data for simulation. The faces are in an upright position in frontal view, with a slight left-right rotation



Fig. 5.1 ATT Database Images

5.1.2. FERET Database

This database has 10 subjects and each subject has 10 images with various poses from -60° to 60° approximately. Size of the images is 112×92 .



Fig. 5.2 FERET Database Images

5.1.3. *WEBCAM Database*

This database has 10 subjects and each subject has 8 images with various poses from -60° to 60° approximately and these pictures are affected by irregular illumination. Size of the images is 112 x 92.

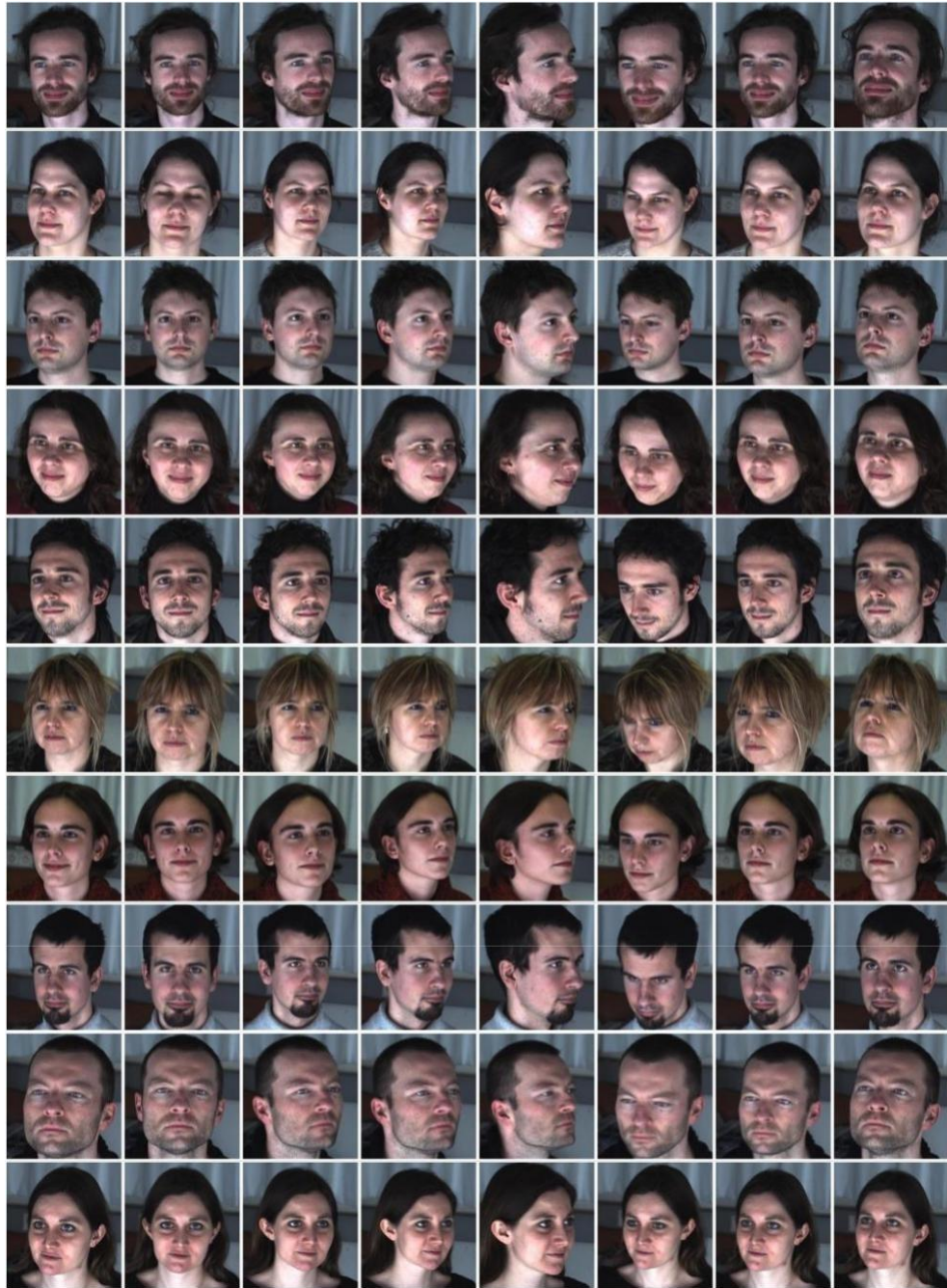


Fig. 5.3 Webcam Database Image

5.2. Training Results

Training is the most important part in CNN where all the layers, forward propagation and backward propagation come into picture. So, we took training recognition and training error at every iteration and plotted line graph for easy analysis.

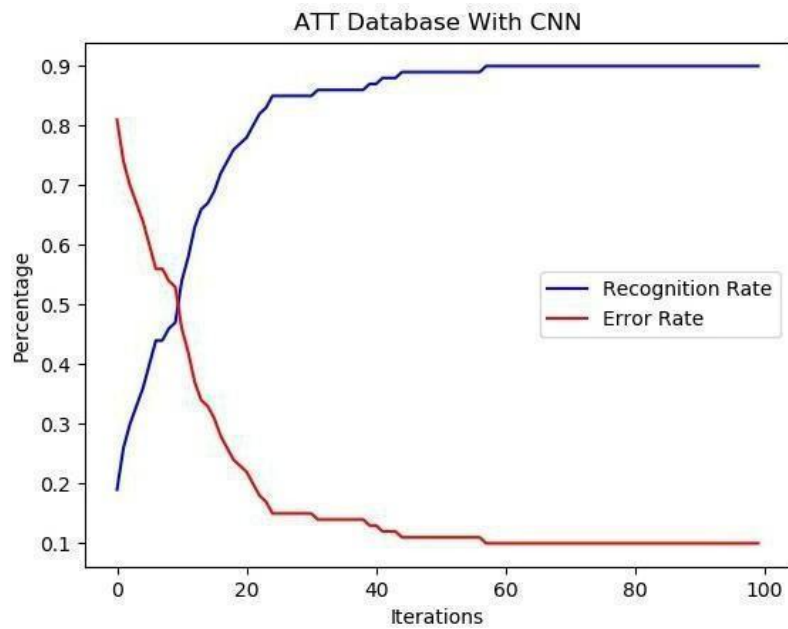


Fig. 5.4 ATT Training Result

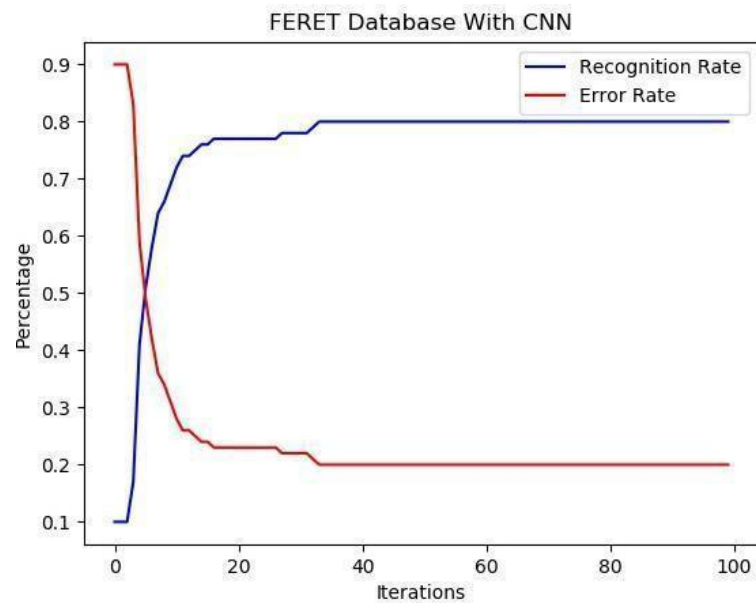


Fig. 5.5 FERET Training Result

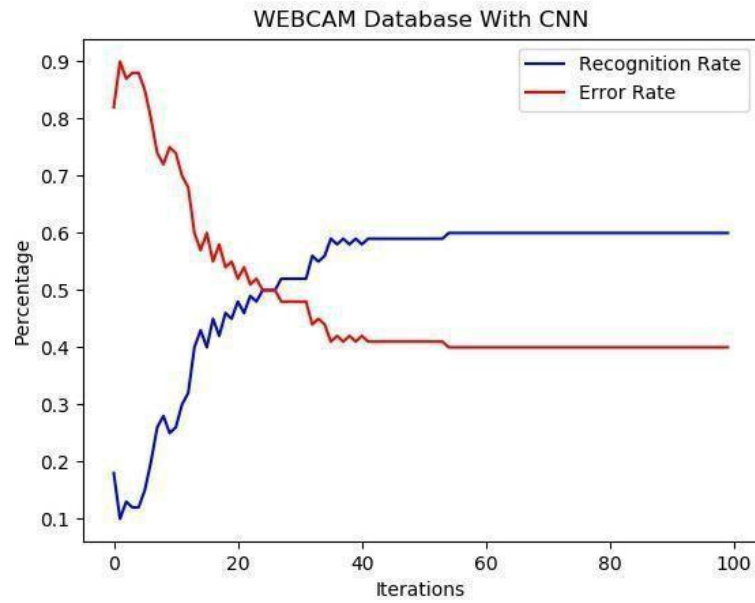


Fig. 5.6 WEBCAM Database Training Result

It is clearly visible that training of Webcam Database is low compared to other databases because we designed our architecture specifically for poses. But, webcam database is highly affected by illumination variation as you can see in the above database images. For FERET database there is a rapid growth in the training because there are only pose variations.

5.3. Testing Results

After training the CNN model we test the model using images of the same database to test its efficiency of recognition. Here, we compare the testing results of various databases.

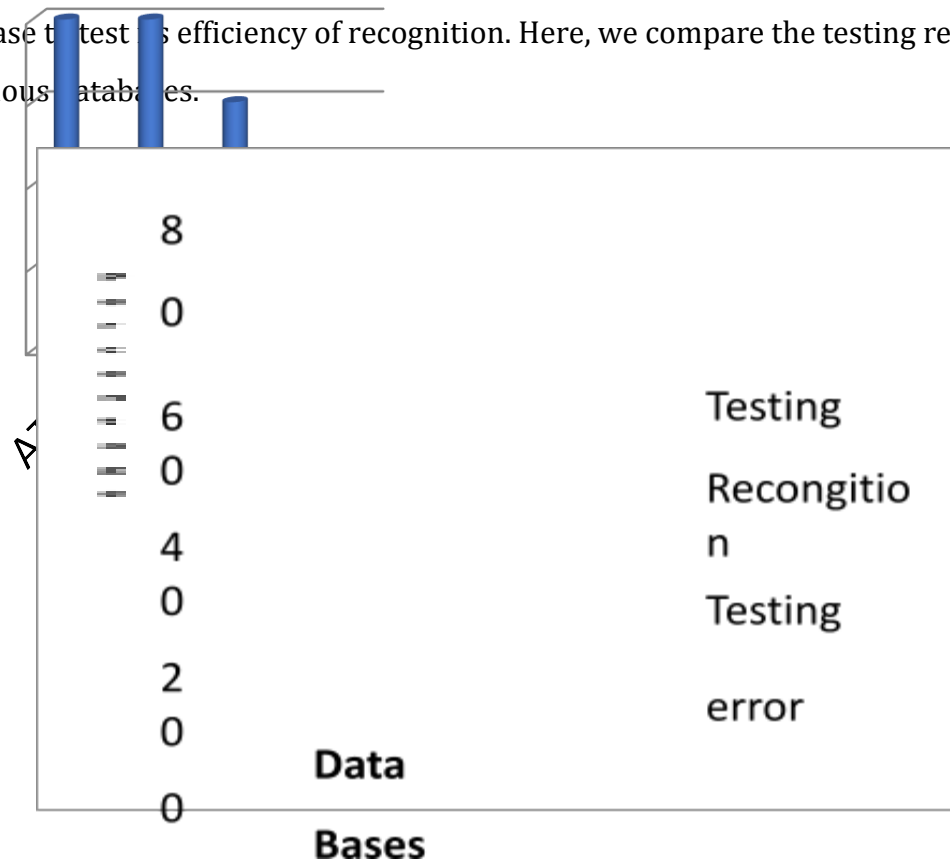


Fig. 5.7 Testing Comparison

5.5. Hardware Comparison

Time complexity of algorithm is the most important part in execution. Unfortunately many neural network and machine learning algorithms have more time complexity. So, we generally go for high end processors and GPU.

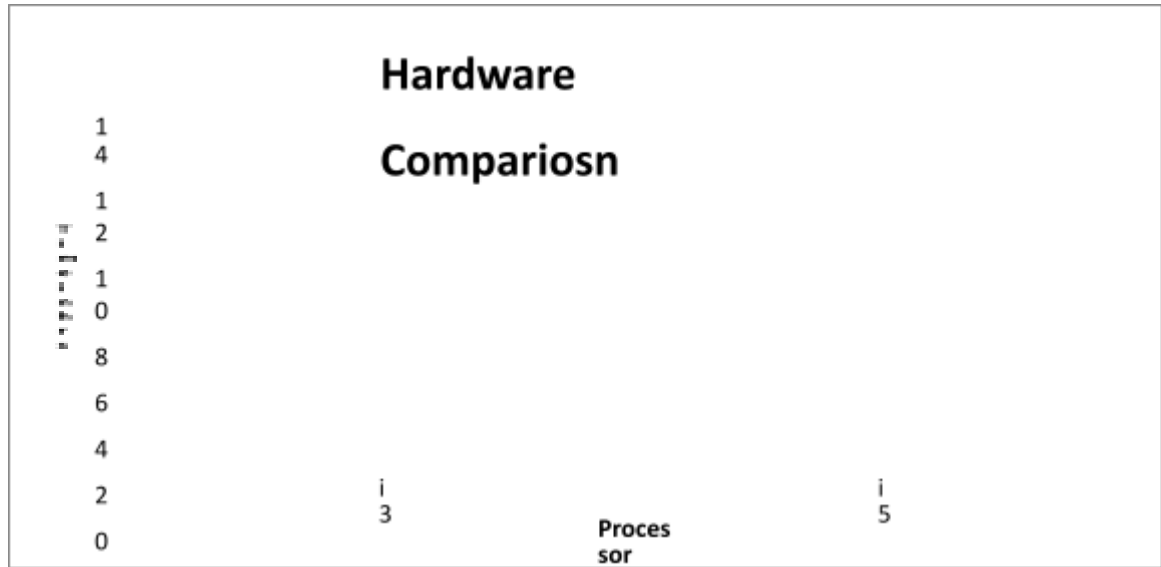


Fig 5.8 Hardware Comparison

5.6 Code

```
1: from getdataset import GetDataSet
2: from neuralnetwork import NeuralNetwork
3: from layers import *
4: import numpy
5:
6: class cnn:
7:     def run(self):
8:         getdataob = GetDataSet()
9:         X_train,Y_train,X_test,Y_test = getdataob.createDataSet(9)
10:        self.n_classes = len(numpy.unique(Y_train))
11:        # Create CNN Feature Extractor
12:        nn= NeuralNetwork(
13:            layers=[
14:                #CNN Feature Extractor
15:                Conv(
16:                    n_filters=4,
17:                    filter_shape=(3, 3),
18:                    weight_scale=0.1,
19:                    weight_decay = 0.01
20:                ),
21:                Activation('relu'),
22:                Pool(
23:                    pool_shape=(2, 2),
24:                    mode='max'
25:                ),
26:                Conv(
27:                    n_filters=8,
28:                    filter_shape=(3, 3),
29:                    weight_scale=0.1,
30:                    weight_decay=0.01
31:                ),
32:                Activation('relu'),
33:                Pool(
34:                    pool_shape=(2, 2),
35:                    mode='max'
36:                ),
37:                Flatten(), # Gives Feature Vectors
38:                # Classifier
39:                Linear(
40:                    n_out= self.n_classes,
41:                    weight_scale=0.1,
42:                    weight_decay=0.002
43:                ),
44:                LogRegression(),
45:            ],
46:        )
47:        # Initialize(Setup) The Layers of CNN
48:        nn._setup(X_train,Y_train)
49:        #Fit the Training Set to CNN to learn the task specific filters for Feature Extraction
50:        nn.fit(X_train,Y_train,learning_rate=0.01,max_iter=100,batch_size=40)
51:        print("\nModel Trained\n\n")
52:        print("\nTesting Prediction : \n")
53:        Y_pred = nn.predict(X_test)
54:        print("Actual : \n",Y_test)
55:        print("Predicted : \n",Y_pred)
56:        error = nn._error(Y_pred,Y_test)
57:        print("Testing Error : ",error)
58:        file = open("weights.txt", "r+")
59:        file.write("\nTesting Error : "+str(error))
60:        ob = cnn()
61:        ob.run()
```

CHAPTER Six

CONCLUSION & FUTURE SCOPE

6.1 CONCLUSION

In this thesis, various approaches to classify the poses of different face images are proposed and investigated. The experimental results show that CNN has the capability to capture the semantic contents from the image spectra. The proposed approaches provide a novel framework for pre-processing steps prior to image classification. The proposed techniques are evaluated through extensive simulations on large standard databases of different categories of pose invariant facial images. The enormous scope of spectra of pose invariant facial images is analyzed.

Although great progress in PIFR has been achieved, there is still much room for improvement, and the performance of existing approaches needs to be further evaluated on real-world databases. To meet the requirement of practical face recognition applications, we propose the following design criteria as a guide for future development.

- **Fully automatic:** The PIFR algorithms should work autonomously, i.e., require no manual facial landmark annotations or pose estimation, etc.
- **Full range of pose variation:** The PIFR algorithms should cover the full range of pose variations that might appear in the face image, including the yaw, pitch, and combined yaw and pitch. In particular, recognition of profile faces is very difficult and largely under-investigated. For pose normalization-based methods, the difficulty lies in the larger error of shape and pose estimation for profile faces.
- **Recognition from a single image:** The PIFR algorithms should be able to recognize a single non frontal face utilizing a single gallery

image per person. This is the most challenging but also the most common setting for real-world applications.

- **Robust to combined facial variations:** As explained in Chapter 1, the pose variation is often combined with illumination, expression, and image quality variations. A practical PIFR algorithm should also be robust to combined facial appearance variations.
- **Matching between two arbitrary poses:** The most common setting for existing PIFR algorithms is to identify non-frontal probe faces from frontal gallery images. However, it is desirable to be able to match two face images with arbitrarily different poses, for both identification and verification tasks. One extreme example is to match a left-profile face to a right-profile face. However, facial symmetry does not exactly hold true for high-resolution images where fine facial textures are clear.
- **Reasonable requirement of labeled training data:** Although a large amount of labeled multi pose training data helps to promote the performance of pose-robust feature extraction based PIFR algorithms, it is not necessarily available because labeled multi- pose data are difficult to collect in many practical applications. Possible solutions may incorporate making use of 3D shape priors of the human head and combining unsupervised learning algorithms.
- **Efficient:** The PIFR algorithms should be efficient enough to realize the requirement of practical applications, e.g., video surveillance and digital entertainment. Therefore, approaches that are free from complicated optimization operations in the testing time are preferable.

6.2 FUTURE SCOPE

- Besides obtaining a robust technique like CNN to classify pose variations in facial images, we are also trying to incorporate volumetric regression so that we can generate different poses with single image and then train them using the proposed architecture.
- Since we know that the problems and challenges faced in facial recognition are inevitable, we were specifically dealing only with poses in this project. But in the near future we can expect many more advancements in this field considering the rapid growth in the field of Machine Learning.
- CNN can also be used for large-scale video classification.
- Besides CNN, Deep learning techniques like Recurrent Neural Networks can be used to implement phase-sensitive and recognition-boosted speech separation.



CHAPTER

Seven

REFERENCES

1. G. Arulampalam and A. Bouzerdoum. A generalized feedforward neural network architecture for classification and regression. *Neural Networks*, 16:561–568, 2003.
2. S. O. Ba and J. M. Odobez. A probabilistic framework for joint head tracking and pose estimation. In *Proceedings of the 17th International Conference on Pattern Recognition*, volume 4, pages 264–267, August 2004.
3. Face Image Analysis with Convolutional Neural Networks Stefan Duffner Master of Science Thesis in Electrical Engineering Face Recognition with Preprocessing and Neural Networks David Habrman 2007.
4. Pose-invariant Face Recognition by Matching on Multi- resolution MRFs linked by Supercoupling Transform Shervin Rahimzadeh Arashloo, Josef Kittler and William J. Christmas.
5. A Comprehensive Survey on Pose-Invariant Face Recognition Changxing Ding Dacheng Tao Centre for Quantum Computation and Intelligent Systems, Faculty of Engineering and Information Technology University of Technology, 2016.
6. Disentangled Representation Learning GAN for Pose- Invariant Face Recognition Luan Tran, Xi Yin, Xiaoming Liu Department of Computer Science.
7. Face Frontalization Using Appearance Flow based Convolutional Neural Network Zhihong Zhang, Xu Chen, Beizhan, Edwin R. Hancock, *Fellow, IEEE*.
8. Illumination and Pose Invariant Face Recognition: A Technical Review K. R. Singh.



9. Logistic Regression and Convolutional Neural Networks Performance Analysis based on Size of Dataset Kartik Chopra, C. Srimathi VIT University.
10. Pose-Invariant Face Alignment via CNN-Based Dense 3D Model Fitting Amin Jourabloo New York 2017.
11. Supplementary Material for Large Pose 3D Face Reconstruction from a Single Image via Direct Volumetric CNN Regression Aaron S. Jackson.