

ABSTRACT

Tuberculosis is one of the most ancient diseases and still it is one of the top 10 causes of death across the world. Most people who get infected with tuberculosis can be saved with proper treatment and their Life can be saved but due to lack of medical support to detect tuberculosis in most parts of World still mortality rate due to tuberculosis is high. This project helps to detect tuberculosis by using image processing techniques over chest x-rays.

Our objective is to prepare a model which classifies the chest x-ray. This model contains two classes normal and abnormal (infected with TB) we need to classify between these two classes and also to achieve high accuracy while classifying.

In our proposed approach we will be improving the accuracy of model by using deep neural networks to train the model and that model helps us in classifying new chest x-ray given as input to the model thus meeting our objective.

To achieve good accuracy, we need to pre-process the images first we pre-processed the images we took images from both datasets Shenzhen and Montgomery and together there are 800 chest x rays we did augmentation over these images and normalized them followed by giving these pre-processed images as inputs to our models.

In this project we used two models baseline CNN model and pretrained VGG16 model and gave pre-processed images as inputs to these both models and evaluated the models to see which performed better comparing using different performance metrics like accuracy, specificity, sensitivity, precision and f1-score and depicted them using graphs and tables From the above results we made a classification reports for both models.

CONTENT

| DESCRIPTION | PAGE NO. |
|--|-----------------|
| MINOR PROJECT REPORT | 1 |
| MINOR PROJECT REPORT | 2 |
| DECLARATION | 3 |
| DECLARATION | 4 |
| CERTIFICATE | 5 |
| CERTIFICATE | 6 |
| ACKNOWLEDGEMENT | 7 |
| ABSTRACT | 8 |
| CONTENT | 9 |
| LIST OF FIGURES | 12 |
| LIST OF TABLES | 13 |
| CHAPTER-1 | 14 |
| INTRODUCTION | 14 |
| 1.1 OVERVIEW | 15 |
| 1.2 Importance of Project and objectives | 16 |
| 1.3 Scope of Project | 16 |
| 1.4 Motivation | 16 |
| 1.5 Organisation of Project Report | 17 |
| CHAPTER-2 | 18 |
| LITERATURE REVIEW | 18 |
| 2.1 OVERVIEW | 19 |
| 2.2 SUMMARY | 20 |
| CHAPTER-3 | 22 |
| METHODOLOGY | 22 |
| 3.1 Proposed Method | 23 |
| 3.2 Datasets used | 23 |

| | |
|---|-----------|
| 3.3 Tools and technologies used | 25 |
| 3.3.1 SciKit Learn | 25 |
| 3.3.2 Pandas | 25 |
| 3.3.3 NumPy | 26 |
| 3.3.4 Keras | 26 |
| 3.3.5 Kaggle Kernel | 26 |
| 3.4 Pre-processing | 27 |
| 3.5 About CNN | 28 |
| 3.5.1 Convolutional layers: | 29 |
| 3.5.2 Pooling layers: | 29 |
| 3.5.3 Fully connected layers: | 30 |
| 3.5.4 Feature Extraction | 30 |
| 3.6 Architecture of models | 30 |
| 3.6.1 Architecture of baseline CNN model | 30 |
| 3.6.2 Architecture of VGG16 model | 34 |
| 3.7 Flow diagram for models and evaluation of models | 35 |
| 3.7.1 Flow diagram for baseline CNN model | 35 |
| 3.7.2 Flow diagram for VGG16 model | 36 |
| 3.7.3 Training and evaluation of models | 38 |
| CHAPTER-4 | 40 |
| RESULT AND PERFORMANCE EVALUATION | 40 |
| 4.1 preprocessing over images | 41 |
| 4.2 Performance evaluation over validation set | 41 |
| 4.2.1 Baseline CNN | 42 |
| 4.2.1 Pretrained VGG16 | 43 |
| 4.3 Training and validation accuracy and loss versus epochs | 43 |
| 4.3.1 Baseline CNN | 43 |
| 4.3.2 Pretrained VGG16 | 45 |
| 4.4 Comparison between Models for different performance metrics | 46 |
| 4.4.1 Accuracy | 46 |

| | |
|---|-----------|
| 4.4.2 Specificity | 47 |
| 4.4.3 Sensitivity | 48 |
| 4.4.4 Precision and F1-score | 49 |
| 4.5 Confusion Matrix | 49 |
| 4.5.1 CM for baseline CNN model | 49 |
| 4.5.2 Confusion matrix for VGG16 model | 50 |
| 4.6 Overall classification report | 51 |
| 4.6.1 Classification report for baseline CNN model | 51 |
| 4.6.2 Classification report for VGG16 model | 52 |
| 4.7 Why VGG16 performed better compared to baseline CNN | 53 |
| CHAPTER-5 | 55 |
| CONCLUSION AND FUTURE SCOPE | 55 |
| 5.1 Conclusion | 56 |
| 5.2 Future Scope | 57 |
| References and useful links | 58 |

LIST OF FIGURES

| DESCRIPTION | PAGE NO. |
|---|-----------------|
| Figure 1:proposed flow of work | 23 |
| Figure 2: Shenzhen dataset images | 24 |
| Figure 3 Montgomery dataset images | 25 |
| Figure 4 : Images after augmentation | 28 |
| Figure 5 CNN architecture (ref: research gate) | 29 |
| Figure 6:Baseline CNN architecture | 31 |
| Figure 7:layers in baseline CNN model | 32 |
| Figure 8:Summary of baseline CNN model | 33 |
| Figure 9 Architecture of VGG16 | 34 |
| Figure 10:flow diagram for baseline CNN model | 35 |
| Figure 11:flow of VGG16 model | 37 |
| Figure 12:Neural networks depicting dropout | 38 |
| Figure 13:Epochs versus loss for baseline CNN | 44 |
| Figure 14:epochs versus accuracy for baseline CNN | 44 |
| Figure 15:epochs versus loss for VGG16 | 45 |
| Figure 16:epochs versus accuracy for VGG 16 | 45 |

LIST OF TABLES

| DESCRIPTION | PAGE NO. |
|--|-----------------|
| Table 1: Accuracy and loss over validation set | 42 |
| Table 2: Accuracy obtained by Models | 47 |
| Table 3: Specificity achieved by models | 48 |
| Table 4: Sensitivity obtained by models | 48 |
| Table 5: Precision and f1-score comparison | 49 |
| Table 6: Confusion matrix for Baseline CNN model | 50 |
| Table 7: Confusion matrix for VGG16 model | 51 |
| Table 8: Classification Report for Baseline CNN mode | 52 |
| Table 9: Classification report for VGG16 model | 53 |

CHAPTER-1
INTRODUCTIO
N

1.1 OVERVIEW

According to WHO TB (tuberculosis) [1] is one of the top 10 causes of death across the world in 2018 around 10.4 million people fell ill from TB i.e. around 28,500 people per day out of them 1.8 million people have died i.e. around 4500 per day. With proper treatment TB can be cured but most of the people have no access for treatment there aren't many specialists to check whether people are infected or not. Manually checking the chest x-rays and detecting the infection needs radiology experts and there are not many experts compared to number of people getting infected.

Image processing is that the wide used modern field in automating utterly totally different techniques and procedures. Image processing has nice impact in medical sector [2] as a result of as patients unit increasing day by day so treating them manually is very powerful, so this approach is utilized for quick automation. we tend to tend to automatize utterly totally different medical procedures previously and it helped at intervals the designation and treatment of various diseases. Variety of the designed automatic systems unit CT scan processed axial imaging digital analysis associate disease that's transmitted through a medium that's a bacterium referred to as mycobacteria. The chest x rays samples of T.B patients unit taken as a medium for screening. The paper collectively includes discussion and analyzation of the techniques that unit previously designed for T.B detection mistreatment analysis photos and chest x rays.

Major issues involved in tuberculosis control programs [2]are the management of pulmonary tuberculosis and Case-finding. The countries which are affected by epidemics of human immunodeficiency virus infection are facing tuberculosis as a more serious problem. The diagnosis of tuberculosis using accurate methods is one of the crucial steps involved to control the occurrence and prevalence of TB. However, the manual diagnosis of tuberculosis is quite complex these days, so there is no standard method at present.

Hence, we need to develop a solution detects whether the person is infected or not with high accuracy using image processing.

1.2 Importance of Project and objectives

Objective of this project is to detect whether the person is infected with TB or not using the chest x-ray of the person.

- To classify across various chest x-rays with high accuracy.
- With this it is easy to detect whether the person is infected or not without involvement of a specialist.
- Thus, it helps many people to know infection at an early stage and take necessary steps to cure it.
- Deaths caused by this epidemic will be drastically reduced as detection becomes easy.

1.3 Scope of Project

- This project has a lot of scope in medical sector as there are many TB infected patients compared to radiology specialists.
- Requirement to achieve this project is dataset containing CXR of normal and infected images.

1.4 Motivation

- Many people are dying across the world because of TB and most of these can be saved by just giving the proper treatment at right time i.e. at early stages.
- First step to do that would be detecting whether a person is infected or not.
- Thus, many lives across the world can be saved and most of the developing countries are facing the problem of TB thus helps them.

1.5 Organisation of Project Report

Organisation of project is done in the following way, initially we started with introduction part where we explained the overview of project including objective, scope of project and motivation.

In the preceding chapters we explained all the research papers we used as a reference to this project. we went through all the papers and studied various advantages and limitations of these papers and explained the summary of all papers in the below sections. In the next chapters we explained the technologies being used, proposed methodology, results and performance evaluation along with conclusion and future scope followed by reference of the project.

SciKit Learn, Pandas library, NumPy, Keras, Kaggle kernel are the technologies being used in this project.

In the proposed methodology we explained about the datasets being used along with pre- processing techniques and baseline CNN architecture and pretrained VGG16 model architecture and followed by the steps involved in training the models.

In chapter 4 evaluation of project is explained and performance results are depicted using tables and graphs followed by conclusion and future scope of project in chapter 5.

2.1 OVERVIEW

- In a paper [3] published in 2017 at “IEEE International Conference on Signal and Image method Applications (ICSIPA)” detailed study is done and written about the potential technique for communicable disease classification victimization CNN is used which classifies the images into 2 categories, throughout this paper they have used a CNN style with 7 convolutional layers and 3 all connected layers and performance of assorted optimizers are a validation accuracy of 82.09% has been obtained.
- In a paper [4] Deep Learning at Chest radiography published by Paras Lakhani and Baskaran Sundaram both AlexNet pre-trained, AlexNet untrained and GoogleNet pre-trained and GoogleNet untrained CNN networks were used and compared pretrained models have provided a better accuracy compared to untrained models and image augmentation has also improved the accuracy of the models. Untrained models have shown an AUC of 0.88 and with augmentation it has improved slightly. Advantage with deep learning is it provides a better accuracy with higher dimensional datasets like images and pre-trained models have performed better compared to untrained model as starting layers of neural networks are the same for all images like edges and blobs.
- In a paper [5] published in 2018 in NCBI (National centre for biotechnology information) focusing on different technologies and methods that can be used for pre-processing segmentation and processing of images paper reviews on common methods of computer aided detection of chest radiographs based on AI. Usually rib structure is not removed but removing rib structure and clavicle the accuracy can be improved. Traditional algorithms like SVM (support vector machine) and random forest may be better, but deep learning methods are providing better performance compared thus they deep learning methods are becoming mainstream in terms of image processing.
- In paper [6] “Comparing deep learning models for population screening using chest radiography” compared across many models and published the results in this