

PLANT LEAF DISEASE DETECTION USING MACHINE LEARNING ALGORITHMS

*A Project Report submitted in partial fulfilment of the requirements for
the award of the degree of*

**BACHELOR OF TECHNOLOGY
IN
ELECTRONICS AND COMMUNICATION ENGINEERING**

Submitted by

T.LALITHYA RAMA(319126512058)

V.PRAVEEN(319126512061)

G.JAYANTH(319126512019)

S.JAYA PRAKASH(319126512051)

Under the guidance of

Ms. P. DEVI M.Tech,(Ph.D)

Assistant professor, ECE



DEPARTMENT OF ELECTRONICS AND COMMUNICATION ENGINEERING

ANIL NEERUKONDA INSTITUTE OF TECHNOLOGY AND SCIENCES

(UGC AUTONOMOUS)

(Permanently Affiliated to AU, Approved by AICTE and Accredited by NBA & NAAC)

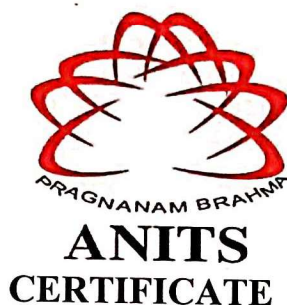
Sangivalasa, Bheemili mandal, Visakhapatnam dist.(A.P)

2022-2023

DEPARTMENT OF ELECTRONICS AND COMMUNICATION ENGINEERING
ANIL NEERUKONDA INSTITUTE OF TECHNOLOGY AND SCIENCES
(UGC AUTONOMOUS)

(Permanently Affiliated to AU, Approved by AICTE and Accredited by NBA & NAAC)

Sangivalasa, bheemili mandal, Visakhapatnam dist.(A.P)



*This is to certify that the project report entitled “PLANT LEAF DISEASE DETECTION USING MACHINE LEARNING ALGORITHMS” submitted by T.Lalithya Rama(319126512058), V.Praveen(319126512061), G.Jayanth(319126512019), S.Jaya Prakash(319126512051) in partial fulfillment of the requirements for the award of the degree of **Bachelor of Technology in Electronics & Communication Engineering** of Anil Neerukonda Institute of Technology and Sciences(A), Visakhapatnam is a record of bonafide work carried out under my guidance and supervision.*


Project Guide

Ms. P. DEVI

M.Tech,(Ph.D)

Assistant Professor

Department of E.C.E

ANITS

Assistant Professor
Department of E.C.E.

Anil Neerukonda

Institute of Technology & Sciences
Sangivalasa, Visakhapatnam-531 162


Head of the Department

Dr. B. Jagadeesh

B.E., M.E., Ph.D

Professor&HOD

Department of E.C.E

ANITS

Head of the Department
Department of E C E

Anil Neerukonda Institute of Technology & Sciences
Sangivalasa - 531 162

ACKNOWLEDGEMENT

We would like to express our deep gratitude to our project **Ms. P. DEVI, Assistant Professor**, Department of Electronics and Communication Engineering, ANITS, for her guidance with unsurpassed knowledge and immense encouragement. We are grateful to, **Dr.B.Jagadeesh**, Head of the Department, Electronics and Communication Engineering, for providing us with the required facilities for the completion of the project work.

We are very much thankful to the **Principal and Management, ANITS, Sangivalasa**, for their encouragement and cooperation to carry out this work.

We express our thanks to all **teaching faculty** of Department of ECE, whose suggestions during reviews helped us in accomplishment of our project. We would like to thank **all non- teaching staff** of the Department of ECE, ANITS for providing great assistance in accomplishment of our project.

We would like to thank our parents, friends, and classmates for their encouragement throughout our project period. At last but not the least, we thank everyone for supporting us directly or indirectly in completing this project successfully.

PROJECT STUDENTS

T. Lalithya Rama (319126512058)

V.Praveen(319126512061)

G. Jayanth Naga Sri Sai(319126512019)

S. Jaya Prakash (319126512051)

ABSTRACT

Abstract. Agriculture is essential for many nations' incomes. Diseases in plants caused by pathogens such as viruses, fungi, and bacteria cause global financial losses in agriculture. Effective disease management ensures crop quality and yield. Disease symptoms are often visible on various plant parts, with leaves being the most affected. Researchers have utilized computer vision, deep learning, few-shot learning, and soft computing techniques to automatically identify plant diseases using leaf images. These technologies help farmers act promptly to protect crops. By automating disease detection, these methods resolve limitations of traditional methods, enhancing both research speed and technology effectiveness. Additionally, molecular techniques have been developed to mitigate pathogenic threats. This review examines the use of machine learning, deep learning, and few-shot learning for automated plant disease detection, reviews diagnostic techniques and future advancements. The integration of these advanced techniques into agricultural practices not only aids in timely disease detection but also supports sustainable farming by reducing reliance on chemical treatments. By leveraging machine learning and molecular diagnostics, farmers can implement targeted interventions, minimizing environmental impact and improving resource efficiency. These innovations are vital in addressing the growing challenges of food security and climate change, ensuring the resilience of agricultural systems worldwide. The CNN model that we built achieves an accuracy of 81.83 %.

KEYWORDS: Deep learning, diagnosis, image processing, machine learning, and plant disease.

CONTENTS

| | |
|---|------------|
| LIST OF FIGURES | ix |
| LIST OF TABLES | xi |
| LIST OF ABBREVIATIONS | xii |
| LIST OF SYMBOLS | xiv |
| CHAPTER 1: INTRODUCTION | |
| 1.1 Plant objective | 02 |
| 1.2 Plant leaf Disease | 02 |
| 1.3 Importance of Plant Disease Detection | 07 |
| 1.4 Techniques for disease detection | 07 |
| 1.4.1 ML methods | 07 |
| CHAPTER 2: LITERATURE SURVEY | 11 |
| CHAPTER 3: IMAGE PROCESSING | 21 |
| 3.1 Introduction to Image Processing | 19 |
| 3.1.1 Purpose of Image Processing | 22 |
| 3.1.2 Types of Images | 23 |
| 3.1.2.1 Binary Image | 23 |
| 3.1.2.2 Gray Scale Images | 24 |
| 3.1.3 Types of Image Processing | 24 |
| 3.1.3.1 Analog Image Processing | 24 |
| 3.1.3.2 Digital Image Processing | 25 |
| 3.2 Methodology for image processing | 25 |
| 3.2.1 Confusion Matrix | 27 |
| 3.2.2 Evaluation Parameters | 28 |
| 3.2.3 Introduction to Machine Learning | 29 |
| 3.2.4 Algorithms used | 30 |
| 3.2.4.1 Random Forest classifier | 30 |

| | |
|---|----|
| 3.2.4.2 Logistic Regression | 31 |
| 3.2.4.3 Support vector machine (SVM) | 32 |
| 3.2.4.3 K-Nearest Neighbour | 32 |
| 3.2.4.4 Naive Bayes | 32 |
| 3.2.5 Technologies used | 32 |
| 3.2.5.1 Python | 32 |
| 3.2.5.2 Libraries | 33 |
| 3.2.5.3 Pandas | 33 |
| 3.2.5.4 Scikit-learn | 34 |
| 3.2.5.5 NumPy | 34 |
| 3.2.5.6 Seaborn | 35 |
| 3.3 Results | 36 |
| 3.3.1 Confusion matrix Result | 36 |
| 3.3.2 Comparison Table for Different Machine Learning Algorithms | 36 |
| 3.3.3 Boxplot Comparison for Different Algorithms | 36 |
| 3.3.4 Results of Plant Leaf Images | 37 |
| CHAPTER 4: METHODOLOGY | 39 |
| 4.1 Introduction | 40 |
| 4.1.1 Support Vector Machine | 41 |
| 4.1.2 K-Nearest Neighbours | 41 |
| 4.1.3 Naive Bayes | 41 |
| 4.1.4 K Means Clustering | 41 |
| 4.1.5 Random Forest | 42 |

| | |
|---|----|
| 4.4.2 Image augmentation in Keras | 55 |
| 4.4.2.1 Techniques present in ImageDataGenerator class | 55 |
| 4.4.2.2 ImageDataGenerator methods | 57 |
| 4.5 Features extraction using Resnet50 | 58 |
| 4.5.1 Need of Feature extraction | 58 |
| 4.5.2 Building the Resnet50 model | 58 |
| 4.6 Training the model | 59 |
| 4.7 Results | 60 |
| CHAPTER 5: CONCLUSION AND FUTURE WORK | 63 |
| REFERENCES | 65 |

LIST OF FIGURES

| Figure No | Title | Page no |
|------------------|------------------------------------|----------------|
| Fig.3.1 | Binary Image | 23 |
| Fig.3.2 | Gray scale Image | 24 |
| Fig.3.3 | Flow Chart | 25 |
| Fig.3.4 | Working of Random Forest Algorithm | 31 |
| Fig.3.5 | Plot of Confusion Matrix | 36 |
| Fig.3.6 | Boxplot Comparison | 37 |
| Fig.3.7 | Healthy Leaf | 37 |
| Fig.3.8 | Diseased (Low) | 37 |
| Fig.3.9 | Diseased (Medium) | 38 |
| Fig.3.10 | Diseased (High) | 38 |
| Fig.4.1 | Transfer Learning | 43 |
| Fig.4.2 | Residual block | 45 |
| Fig.4.3 | Resnet50 Architecture | 46 |
| Fig.4.4 | Identity Block | 47 |
| Fig.4.5 | Convolution Block | 47 |
| Fig.4.6 | Resnet50 Architecture | 47 |
| Fig.4.7 | Methodology | 47 |
| Fig.4.8 | Training accuracy | 60 |
| Fig.4.9 | Validation accuracy | 60 |
| Fig.4.10 | Diseased leaf | 61 |
| Fig.4.11 | Healthy leaf | 61 |
| Fig.4.12 | Diseased with Tomato early blight | 61 |

| | | |
|----------|-----------------------------------|----|
| Fig.4.13 | Diseased with Tomato spider mites | 62 |
| Fig.4.14 | Diseased with Potato early blight | 62 |
| Fig.4.15 | Diseased with Tomato late blight | 62 |

LIST OF TABLES

| Figure No | Title | Page no |
|-----------|--|---------|
| Table 3.1 | Confusion Matrix | 27 |
| Table 3.2 | Comparison Table for Different Machine Learning Algorithms | 36 |

LIST OF ABBREVIATIONS

| Abbreviation | Full form | Page no |
|---------------------|---|----------------|
| KNN | K-Nearest Neighbour | 07 |
| SVM | Support Vector Machine | 08 |
| CNN | Convolutional Neural Network | 12 |
| GPDCNN | Global pooling dilated CNN | 12 |
| MCNN | Multilayer convolutional neural network | 13 |
| PSO | Particle Swarm Optimization | 11 |
| DCNN | Deep Convolution Neural Network | 14 |
| CCM | Color Co-occurrence Matrix | 16 |
| GLCM | Gray Level Co-occurrence Matrix | 16 |
| MER | Minimum Enclosing Rectangle | 16 |
| DWT | Discrete Wavelet Transform | 16 |
| SIFT | Scale Invariant Feature Transform | 16 |
| HOG | Histogram of Oriented Gradients | 19 |
| STL-10 | Self-Taught Learning 10 | 19 |
| ResNet | Residual Network | 20 |
| HSV | Hue Saturation Value | 26 |
| RGB | Red, Green, Blue | 26 |
| HDF | Hierarchical Data Format | 27 |
| AI | Artificial intelligence | 29 |
| MDP | Markov decision process | 30 |
| ML | Machine Learning | 31 |
| HTTP | HyperText Transfer Protocol | 33 |

| | | |
|--------|---|----|
| WSGI | Web Server Gateway Interface | 33 |
| PEP | Python Enhancement Proposal | 33 |
| BSD | Berkeley Source Distribution | 33 |
| AQR | Applied Quantitative Research | 33 |
| DBSCAN | Density-based spatial clustering of applications with noise | 34 |
| API | Application programming interface | 35 |
| RF | Random Forest | 36 |
| LR | Logistic Regression | 36 |
| VGG | Visual Geometry Group | 36 |
| JPEG | Joint Photographic Experts Group | 53 |
| BMP | Bitmap | 53 |

LIST OF SYMBOLS

| Symbol | Description | Page no |
|---------------|-----------------------|----------------|
| pH | Potential of hydrogen | 03 |
| TP | True Positive | 28 |
| TN | True Negative | 28 |
| FP | False Positive | 28 |
| FN | False Negative | 28 |

CHAPTER 1

INTRODUCTION

1.1 Project Objective:

The main goal of plant disease detection is to accurately and promptly identify diseases. This is essential for preventing their spread and reducing crop losses. Early detection plays a vital role in limiting damage and controlling outbreaks. When using machine learning algorithms for disease detection, the results must be reliable and consistent.

1.2 Plant Leaf Disease:

A disruption in the supply of herbs can result from diseases that damage or alter essential qualities. All herbs—whether wild or cultivated—are vulnerable to disease. Each plant family is susceptible to certain illnesses, though these are generally uncommon. The occurrence and spread of plant diseases vary from season to season, depending on the presence of pathogens, environmental conditions, and the types and varieties of herbs being cultivated. Some herb species are especially prone to widespread disease outbreaks, while others show greater resistance.

- **Definitions of plant disease**

Generally, when plants are consistently disturbed, they might catch illness. pathogens that result in abnormal physiological processes that disrupt the normal structure, growth, function, or other activities of the plant. The essential physiological or biochemical processes of the plant are disrupted, which results in the typical diseased states or symptoms. Depending on whether the primary cause of the disease is infectious or non- infectious, plant diseases may be broadly categorised. Infectious plant diseases are caused by pathogens such as fungi, bacteria, mycoplasma, viruses, viroids, nematodes, or parasitic flowering plants. Within or on a host, infectious organisms can grow and spread from one vulnerable host to another. Unfavorable growth circumstances, including as excessive temperatures, unfavorable moisture-oxygen ratios, soil and air pollutants, and an abundance or shortage of vital minerals, are the root causes of non-infectious plant illnesses.

- **Temperature**

All pathogens have an optimal temperature range for growth. This ideal temperature can vary slightly depending on different stages of fungal development, such as spore production, spore germination, and the growth of mycelium (the thread-like structure of fungi). Storage conditions for certain produce like berries, leafy greens, and root vegetables are carefully controlled to inhibit fungal and bacterial spoilage—provided that temperature adjustments do not compromise the quality of the produce.

While there is limited control over temperature in open fields (aside from protective dips), greenhouse environments offer greater flexibility. By managing temperature and humidity, it is possible to reduce the risk of disease outbreaks. High humidity combined with favorable temperatures can encourage diseases such as downy mildew on vines (*Pseudoperonospora cubensis*), lima bean blight (*Phytophthora phaseoli*), and late blight on potatoes and tomatoes (*Phytophthora infestans*), as well as leaf spot in sugar beets.

However, temperature can also mask the symptoms of certain aggressive or mycoplasma-related diseases, making them more difficult to detect.

Relative humidity

Relative humidity plays a crucial role in the germination of fungal spores and the development of storage rot. *Rhizopus stolonifer*, a common storage disease affecting sweet potatoes, does not progress when the relative humidity of the tuber's surface is maintained between 85–90%, even if the storage temperature favors pathogen growth. Under these conditions, sweet potato roots develop a protective corky tissue that helps resist *Rhizopus* infection.

High humidity is generally essential for spore germination, bacterial proliferation, and infection spread. Most molds thrive at relative humidity levels of 90–95%. In greenhouse environments, diseases such as those caused by *Botrytis* species—which rot herb leaves, stems, and seedlings—can be managed by lowering humidity levels and avoiding overhead watering.

Soil pH

Soil pH plays a crucial role in the occurrence of certain plant diseases, such as clubroot in cruciferous vegetables and potato scab (*Streptomyces scabies*). A soil pH of 5.2 or lower inhibits the development of potato scab. (Note: a pH of 7 is neutral; values below 7 are acidic, and those above 7 are alkaline.) In naturally acidic soils with a pH around 5.2, potato scab is generally not a problem. To maintain this acidity, some growers apply sulfur to keep the soil pH near 5.0.

Conversely, for cruciferous vegetables affected by clubroot, increasing the soil pH to 7.2 or above—typically by thoroughly incorporating lime—can help suppress the disease and reduce its impact

.Requirements for disease development

An infectious plant disease cannot occur unless all three of the following key elements are present: (1) suitable environmental conditions—particularly factors like rainfall or heavy dew, relative humidity, and air and soil temperatures; (2) a virulent pathogen; and (3) a susceptible host plant. Effective disease prevention strategies focus on disrupting this "disease triangle" of environment, pathogen, and host.

For example, increasing the host plant's resistance or immunity through plant breeding or genetic engineering can reduce disease impact. Similarly, modifying environmental conditions to favor the growth of the plant over the pathogen can help prevent infection. Lastly, targeting the pathogen directly—by eliminating it or stopping its spread—can also break the cycle of disease development.

Classification of plant diseases by causal agent

Plant diseases are often classified based on their physiological effects or visible symptoms. However, many diseases display similar symptoms and signs despite being caused by entirely different pathogens or agents, each requiring distinct control measures. Symptom-based classification is also limited because a single causal agent can produce a variety of

symptoms—even on the same plant organ—many of which may appear simultaneously.

Another approach is to classify diseases by the plant species they affect. Host indexes, which list known diseases for specific hosts across regions, countries, or continents, are valuable diagnostic tools. When a seemingly new disease appears on a familiar host, consulting the relevant host index can often help identify the causal agent. Diseases may also be categorized using other criteria, depending on the context and purpose.

- **Noninfectious disease-causing agents**

Non-infectious plant diseases can arise suddenly and are caused by a range of environmental and physical factors. These include poor soil moisture and oxygen levels, extreme soil pH (acidity or alkalinity), high or low temperatures, pesticide damage, toxic chemicals in the air or soil, changes in soil grade, root girdling, mechanical or electrical injury, and soil compaction. Preharvest and storage conditions that are unsuitable for fruits, vegetables, or nursery plants often contribute to significant losses.

Non-infectious diseases can affect many plant species within a particular area or ecosystem. Although they can lead to substantial damage, these conditions are often difficult to prevent or treat, as they stem from broader environmental factors beyond human control. Symptoms may not appear until weeks or even months after the initial stress event.

- **Infectious disease-causing agents**

Thousands of species from incredibly varied families of creatures can infect plants. Few are macroscopic, while the majority are tiny. The infectious agents, also known as pathogens, include bacteria, fungi, nematodes, generally known as mycoplasma-like organisms (MLOs), and parasitic seed herbs.

Diseases caused by viruses and viroids

- **General characteristics**

Viruses and viroids are the smallest of all infectious agents. A virion is an infectious particle that has achieved structural maturity. The size and shape of virions can range from about 20 nanometers (0.0000008 inch) to 250–400 nanometers. Unlike viruses, viroids lack structural proteins, such as the protein coat (capsid) found in viruses. Both viruses and viroids are obligate parasites, meaning they can only replicate or multiply within a living host cell. A single plant species can be vulnerable to multiple viruses or viroids. Viral infections often cause serious diseases in important crops.

- **Diseases Caused by Fungi**

Approximately two-thirds of infectious plant diseases are caused by fungi. All commercially important plants are susceptible to one or more fungal infections, and in many cases, multiple fungal species can affect a single plant species.

- **General Characteristics**

Fungi are a diverse and extensive group of eukaryotic microorganisms. They lack chlorophyll and possess rigid cell walls. Their cells also contain a membrane-bound nucleus. Many fungi have a vegetative body composed of microscopic branching filaments called hyphae (plural:

hypha), some of which extend into the air while others penetrate the substrate where they grow. The network of hyphae is known as the mycelium, and the "cottony" or "fuzzy" appearance of fungal growth is due to the bulk of the mycelium. Fungi reproduce in various ways, including both asexual and sexual reproduction, producing vast quantities of spores. For example, the color of moldy bread is due to the fungal spores.

- **Diseases Caused by Nematodes**

Nematodes, or roundworms, are unsegmented, active, parasitic roundworms that attack plants (also called nemas or eelworms). Most of them are too small to be seen with the naked eye, with adult nematodes typically measuring between 0.25 and 2 mm in length. Approximately 1,200 species of nematodes cause plant diseases. Almost every type of plant is susceptible to at least one species of nematode. While most nematodes reside in the soil and feed on tiny roots, some species also live in and feed on bulbs, buds, stems, leaves, and flowers. Nematodes feed on plant juices using a hollow, needle-like mouthpart known as a stylet. When the stylet is inserted into plant cells, the nematode injects a liquid containing enzymes that break down the cell contents, which are then drawn back into the nematode's digestive system. Nematode feeding weakens the plant's resistance, reduces plant vigor and yield, and can provide a point of entry for pathogens causing root rot or wilt.

Plants affected by nematodes often show symptoms similar to those caused by excessive soil moisture, sunburn, frost, mineral imbalances, or insect damage. Typical signs include stunted growth and yellowing of leaves. Nematode infestations can cause plant cells to grow or degenerate. In some cases, both occur. Native nematodes can target cultivated plants when their natural hosts are removed, and others are spread through seedlings, bulbs, tubers, and soil surrounding infected nursery stock. Nematodes may remain active in the soil near roots or in fallow fields. They can enter plants through wounds, natural openings, or by penetrating roots. Some nematodes are endoparasites (feeding within plant tissues), while others are ectoparasites (feeding on the plant's surface). Nematodes need living plant tissues to reproduce, and they are attracted to host roots by heat or chemicals released by the roots. Nematodes typically go through four developmental stages, from egg to adult, with a complete generation cycle taking 20 to 60 days. Some species have only one generation per year, but still produce hundreds of offspring.

The development of nematodes is influenced by factors such as the growing season, temperature, water and nutrient availability, and soil properties. Other pests, such as viruses, protozoans, mites, flatworms, and other nematodes, as well as nematode-parasitic bacteria, also play a role in nematode populations. Crop rotation, previous cropping practices, and the application of toxic substances to the soil also affect nematode populations. Additionally, plant species, variety, age, and nutrition all influence nematode development.

1.3 Importance of Plant Disease Detection:

It is Important for Correct Plant Disease Identification? Disease control initiatives may result in a waste of time and resources without accurate identification. Additional plant losses could result from the application of disease control strategies that are inadequate to handle the disease-causing agent. Infectious parasites including nematodes, fungi, oomycetes, viruses, and bacteria are the root cause of plant illnesses. Because a large range of organisms can cause a variety of symptoms (Figure 1), accurate pathogen identification is essential to creating a management plan. Injury vs. Illness It's critical to comprehend

the distinctions between a plant injury and a disease. A sudden injury results from an outside force over a brief period of time.

1.4 Techniques For Disease Detection:

1.4.1 Machine Learning Methods

K-Nearest Neighbour (KNN) Algorithm for Machine Learning

- One of the simplest Machine Learning algorithms, K-Nearest Neighbor is based on the Supervised Learning approach.
- The K-NN algorithm makes the assumption that the new case and the data are comparable to the cases that already exist, and it places the new instance in the category that is most similar to those cases.
- The K-NN algorithm saves all the information that is accessible and categorises additional data points based on similarity. This means that utilising the K-NN method, fresh data can be quickly and accurately sorted into a suitable category.
- Although the K-NN approach can be used for both classification and regression

problems, classification challenges are where it is most frequently applied.

- Since K-NN is a non-parametric technique, it makes no assumptions about the underlying data.
- It is also known as a lazy learner algorithm since it saves the training dataset rather than learning from it immediately. Instead, it uses the dataset to perform an action when classifying data.
- KNN algorithm simply stores the dataset during the training phase and subsequently classifies new data into a category that is quite close to the new data.
- **Example:** When training, the KNN algorithm simply stores the dataset; when it receives new data, it then classifies that data into a category that is quite similar to the new data.

Support Vector Machine Algorithm

One of the most popular supervised learning algorithms is called the Support Vector Machine (SVM), and it is employed to solve Classification and Regression problems. However, it is mostly used in Machine Learning Classification issues. The objective of the SVM method is to find the best decision boundary or line that can classify the n-dimensional space, allowing us to classify additional data points with ease in the future. A hyperplane is the name for this optimal boundary. To assist in creating the hyperplane, SVM selects the extreme vectors and points. Support vectors, which are used to represent these extreme scenarios, are the basis of the SVM methodology. View the graphic below to see how a choice classifies two separate groups.

Types of SVM

There are two types of svm:

- **Linear SVM:** The phrase "linearly separable data" describes information that can be split into two categories using just one straight line. This type of data is classified using Linear SVM, and the classifier that is utilised is called the Linear SVM classifier.
- **Non-linear SVM:** A dataset is considered to be non-linear if it cannot be classified using a straight line, in which case the classification technique used is called a non-linear SVM classifier.

- **Hyperplane and Support Vectors in the SVM algorithm:**

Hyperplane: In n-dimensional space, the classes can be divided into a variety of lines or decision borders; nevertheless, it is necessary to choose the best decision boundary for categorising the data points. This ideal boundary is known as the SVM hyperplane.

Given that the dataset's features define the hyperplane's dimensions, a straight line will be the hyperplane if there are just two features (as in the example image). In addition, if there are three features, hyperplane will be a two-dimensional plane.

Support Vectors: The closest data points or vectors near the hyperplane and those that have the most bearing on the hyperplane's position are called support vectors. Because they support the hyperplane, these vectors are referred to as support vectors.

- **Random Forest Algorithm**

Preferred machine learning algorithm Random Forest is a part of the supervised learning strategy. It might be applied to ML issues that call for both regression and classification. It is built on the idea of ensemble learning, which is a method for integrating many classifiers to solve complex issues and enhance model performance.

Random Forest, as the name indicates, is a classifier that increases the projected accuracy of the dataset by averaging numerous decision trees applied to different subsets of the provided data. Instead of depending just on one decision tree, the random forest gathers forecasts from each decision tree and predicts the result based on the votes of the majority of projections.

- **Logistic Regression in Machine Learning**

- Logistic regression is one of the most well-known Machine Learning algorithms that falls under the umbrella of Supervised Learning. It is used to forecast the categorical dependent variable using a specified set of independent variables.
- Logistic regression may be used to forecast the outcomes of a categorical dependent variable. The outcome must thus be a discrete or categorical value. It offers the probabilistic values that lie between 0 and 1 rather than the precise values between 0 and 1. It can be either True or False, 0 or 1, or Yes or No.

- The use of logistic regression and linear regression differs significantly. In contrast to linear regression, which is used to address classification issues, logistic regression addresses regression issues.
- Instead of fitting a regression line, we use a logistic function with a "S" shape that predicts two maximum values (0 or 1) in logistic regression.
- The logistic function's curve displays the probability of a number of events, including whether or not the cells are cancerous, a mouse.

- **Logistic Function (Sigmoid Function):**

- A mathematical function called the sigmoid function is employed to convert anticipated values into probabilities.
- It transforms any real value between 0 and 1 into another value.
- Since the logistic regression's value must lie within the range of 0 and 1, it can never go above or below this limit, resulting in a "S"-shaped curve. The sigmoid function or logistic function is another name for the S-form curve.
- We apply the threshold value idea in logistic regression, which establishes the likelihood of either 0 or 1. Examples include values that incline to 1 over the threshold value and to 0 below it.

- **Assumptions for Logistic Regression:**

- A categorical variable must be the dependent one.
- Multicollinearity should not exist in the independent variable.

- **Type of Logistic Regression:**

Logistic regression may be divided into three types according to the categories:

- Binomial: A dependent variable in a binomial logistic regression may only be one of two potential kinds, such as 0 or 1, Pass or Fail, etc.
- Multinomial: The dependent variable in multinomial logistic regression may be one of three or more potential unordered kinds, such as "cat," "dogs," or "sheep."
- Ordinal: In an ordinal logistic regression, the dependent variables may be categorised as "low," "Medium," or "High."

CHAPTER 2

LITERATURE SURVEY

1. Unsupervised image translation using adversarial networks for improved plant disease recognition.

Reference: Nazki et al.

Dataset: 2789 tomato plant disease images

Technique used: Generative Adversarial Network And Deep CNN

Output: Accuracy= 86.1%

Advantages:

Better demonstrating of information appropriation (pictures more honed and more cleared). GANs can prepare any sort of generator organization.

Disadvantages:

Difficult to prepare, unstable training process. Require many guidelines to obtain satisfying results. Mode Collapse issue.

2. Cucumber leaf disease identification with global pooling dilated convolutional neural network.

Reference: Zhang et al.

Dataset: Acquisition of 600 cucumber sick leaves of 6 regular cucumber leaf infected

Technique used: GPDCNN

Output: Accuracy = 94.65%

Advantages:

GPDCNN is more robust than different strategies.

Disadvantages:

Completely associated layer has such a large number of parameters which decreases the speed of preparing (training) and effectively bring about over-fitting.

3. Multilayered Convolution neural network for the Classification of mangoleaves infected by Anthracnose Disease.

Reference: SINGH CHOUHA N et al.

Dataset: Captured images at SMVDU, Katra

Technique used: Multilayer convolutional neural network (MCNN)

Output: Accuracy = 97.13%

Advantages:

The essential advantage of MCNN diverged from its paradigms is that it therefore perceivesthe critical features with no human administration.

Disadvantages:

MCNN has a few layers then the training process takes a ton of time if the PC doesn't compriseof a good CPU.

4. Sunflower leaf diseases detection using Image Segmentation based onParticle swarm optimization.

Reference: Vijai Singh

Dataset: Capture Sunflowers leaves.

Technique used: Particle Swarm Optimization Algorithm.

Output: Accuracy = 98%

Advantages:

The upsides of PSO are that PSO is easy to implement and there are scarcely any boundariesto change.

PSO perform in a way that is better than the GA as for computational efficiency.

Disadvantages:

PSO is one of the well known techniques, however its application for the issue isn't confoundedbecause of the simple characteristics.

5. Deep Convolutional neural network based detection system for real time corn plant disease recognition.

Reference: Mishra et al.

Dataset: Plant Village dataset.

Technique used: Deep Convolution Neural Network

Output: Accuracy = 88.46%

Advantages:

With little dependence on pre-processing, this algorithm requires less human effort. It is actually a self-learner, which makes the preprocessing phase, easier.

Disadvantages:

It requires an enormous dataset to process and train the neural organization.

6. Performance analysis of deep learning CNN models for disease detection in plants using image segmentation

Reference: Sharma et al.

Dataset: Tomato healthy and infected leaves images

Technique used: Convolution Neural Network

Output: Accuracy = 98.6%

Advantages:

Perhaps the greatest favorable position of CNN is the programmed extraction of highlights by handling straightforwardly the crude pictures.

Disadvantages:

CNNs don't have arranged outlines which are a fundamental component of human vision.

7. Tomato Leaf Disease Detection using Convolution Neural Network.

Reference: Agarwal et al.

Dataset: Images taken from Plant Village dataset.

Technique used: Images taken from Plant Village dataset.

Output: Classification Accuracy= 76% to 100%, Average Accuracy for disease=91.2 %

Advantages:

The Storage space needed by proposed model was of order of 1.5MB where as pre prepared models had additional room need of around 100MB appropriately demonstrated the upside of the proposed model over pre-trained models.

Disadvantages:

A CNN is essentially more slow because of an activity, for example, pooling.

8. Seasonal Crops Disease Prediction and classification Using Deep Convolutional Encoder Network

Reference: Khamparia et al.

Dataset: Plant Village Dataset

Technique used: Deep Convolution Encoder Network

Output: Accuracy = 97.50%

Advantages:

Softmax classifier is used at output layer. It returns the probabilities of each class if there should arise an occurrence of a multiorder model, and the target class should have high probability.

Disadvantages:

This method lack a mechanism to map deep layer feature maps to input dimensions.

9. Deep Neural Networks Based Recognition of Plant Diseases by Leaf Image Classification

Reference: Sladojevic et al.

Dataset: Capture images by agricultural experts. **Technique used:** Deep Convolution Neural Network **Output:** Accuracy= 96.3%

Advantages:

DCNNs included image and object classification, face detection, and image segmentation. DCNN have more hidden layers especially more than 5, which increases the accuracy.

Disadvantages:

CNN don't encode the position and direction of an object.

Absence of ability to be spatially invariant to the input information.

10. A Review of Machine Learning Approaches in Plant Leaf Disease Detection and Classification

Reference: MAJJI V APPLALANAIDU, G. KUMARAVELAN.

Dataset: plant village dataset

Technique used: Color Co-occurrence Matrix(CCM), Gray Level Co-occurrence Matrix(GLCM), Minimum Enclosing Rectangle(MER), Color Co-occurrence Matrix(CCM), CCM, GLCM, Discrete Wavelet Transform(DWT) Scale Invariant Feature Transform(SIFT)

Objective: This review provides a comparative analysis of various state-of-the-art ML and DL algorithms to identify and categorize plant leaf diseases.

11. Research on machine learning framework based on random forest algorithm

Reference: Qiong Ren , Hui Cheng and Hai Han.

Technique used: Random forest algorithm

Objective:

This article examines and analyses the machine learning framework based on the random forest algorithm with the goal of enhancing the random forest algorithm's current restrictions. It also creates and implements a number of machine learning frameworks.

12. Random Forest with Adaptive Local Template for Pedestrian Detection

Reference: Tao Xiang, Tao Li, Mao Ye, and Zijian Liu. **Dataset:** TUD Pedestrians, INRIA pedestrians **Technique used:** Random forest

Output: Accuracy= 90.8%

Objective: Detection of pedestrians in cluttered environments. The main concept of our approach is to combine several weak classifiers that are specified by adaptive local templates and to do it using Random Forest. Iteratively and layer-by-layer, the forest is constructed. The splitting functions in the forest are learned using the adaptive local templates, and when they are all of the same depth, they produce a weak classifier. By minimising a global loss and adding each new weak classifier, sample weights are updated. The suggested technique achieves the state-of-the-art or competitive performance, according on the final experimental findings on two difficult pedestrian datasets.

13. Improving the Random Forest Algorithm by Randomly Varying the Size of the Bootstrap Samples

Reference: Md Nasim Adnan.

Technique used: Random Forest

Objective: By using the Random Subspace technique on bootstrap samples for large dimensional datasets, the Random Forest algorithm provides quite a variety of decision trees as the basic classifiers. By introducing more variation among the decision trees, we may improve the ensemble accuracy. Every time a decision tree is generated as the basis classifier in Random Forest, the size of the bootstrap files stays the same. To improve the accuracy of the forest, we suggest changing the size of the bootstrap samples at random within a predetermined range. We do a thorough experiment using a variety of datasets from the UCI Machine Learning Repository. The outcomes given in this research demonstrate the enormous potential of our method.

14. An Ensemble Random Forest Algorithm for Insurance Big Data Analysis

Reference: Ziming Wu, Weiwei Lin, Zilong Zhang and Angzhan Wen.

Technique used: Random Forest

Objective: This paper analyses the imbalance distribution of insurance business data, concludes the preprocessing algorithms of the imbalance dataset, and proposes an ensemble random forest algorithm based on Apache Spark which can be used in the large scaled imbalanced classification of insurance business data. The experiment results showed that the ensemble random forest algorithm is more suitable in the insurance product recommendation or potential customer analysis than the traditional tra Preprocessing for unbalanced classification algorithms might make use of the suggested bootstrap undersampling approach in conjunction with KNN. Together with bootstrap sample preprocessing, ensemble learning techniques may be able to speed up learning even further. They also provide a useful comparison to other unbalanced data mining strategies.

15. leaf and skin disease detection using image processing

Reference: Manjunath Badiger ,Varuna kumara ,Sachin CN shetty,Sudhir poojary

Technique used: K-means algorithm and SVM classifier

Output: Accuracy = 96.3%

Advantages:

Easy to understand and implement. Can handle large datasets well. Disadvantages of K-Means Sensitive to number of clusters/centroids chosen.

Disadvantages:

It requires to specify the number of clusters (k) in advance. It can not handle noisy data and outliers. It is not suitable to identify clusters with non-convex shapes.

16. plant disease detection using machine learning

Reference: Niveditha M, pooja R, prasad Bhat N, shashank N

Technique used: HOG, Random forest

Output: Accuracy = 70%

Advantages:

Work well for small resolutions. Typically does detection via classification, i.e. uses a binary classifier.

Disadvantages:

More time consuming to construct than a frequency polygon.

17. plant disease detection using CNN

Reference: Nishant Shelar, Suraj shinde, Shubham sawant, Shreyas dhmal

Technique used: CNN **Output:** Accuracy = 96% **Advantages:**

local spatial coherence in the input (often images), which allow them to have fewer weights as some parameters are shared. This process, taking the form of convolution makes them especially well suited to extract relevant information at a low computational cost.

Disadvantages:

Classification of Images with different Positions, Adversarial examples, CoordinateFrame, Other minor disadvantages like performance.

18. Pest detection in crop using video and Image processing

Reference: Madhuri Devi Chodey, Dr. Noorilla Shariff c, Gauravi Shetty

Technique used: K-means algorithm and SVM classifier

Output: Accuracy = 96.3%

Advantages:

SVM works relatively well when there is a clear margin of separation between classes. SVM is more effective in high dimensional spaces.

Disadvantages:

SVM algorithm is not suitable for large data sets. SVM does not perform very well when the data set has more noise i.e. target classes are overlapping.

19. Image Classification Using Resnet-50 Deep Learning Model

Reference: Aryan Garg **Dataset:** STL-10 **Technique used:** Resnet-50

Output: Accuracy= 76.229%

Advantages:

Networks with large number (even thousands) of layers can be trained easily without increasing the training error percentage.

Disadvantages:

High computational complexity - Residual neural networks can often require significant processing power and may not be suitable for certain tasks.

20. Deep Learning in Image Classification using Residual Network

(ResNet) Variants for Detection of Colorectal Cancer

Reference: Devvi Sarwinda , Radifa Hilya Paradisa , Alhadi Bustamam ,Pinkie Anggia

Dataset: Warwick-QU

Technique used: Deep Residual Network (ResNet)

Output: Accuracy= 73% -88%

Advantages:

ResNets help in tackling the vanishing gradient problem using identity mapping.

Disadvantages:

High memory requirements - Residual networks require large amounts of memory in order to store the necessary parameters and weights

CHAPTER 3

IMAGE PROCESSING

3.1 Introduction to Image Processing:

Image Processing: A Digital Approach

Image processing is the technique of converting a physical image into a digital format, allowing for manipulation, enhancement, or information extraction. In image science, it refers to any method where the input is an image—like a photograph or video frame—and the output may be a modified image or a set of features derived from it.

While analog and optical methods exist, digital image processing is the most commonly used. The process starts with image acquisition, using devices like digital cameras or optical scanners to capture the initial input. The image is then processed to improve its quality or extract meaningful data.

This step is crucial in many deep learning and computer vision applications, where preprocessing can significantly enhance model performance. In the entertainment industry, image processing is often used for tasks like adding or removing objects from pictures.

Typically, images are treated as two-dimensional signals and processed using conventional signal processing techniques. Subsections of images, known as regions of interest (ROIs), may be analyzed individually, as they often contain meaningful clusters of visual elements.

One practical application is defect detection, where image processing helps identify surface flaws by focusing on the damaged areas. As a result, it has become a rapidly growing technology across diverse industries, especially in manufacturing.

Image processing generally involves three main steps:

1. Image acquisition using a digital device.
2. Analysis and manipulation, including data compression, enhancement, and detecting patterns invisible to the human eye.
3. Output generation, where the processed image or a report is produced based on the analysis.

3.1.1 Purpose of Image Processing:

Image processing is classified into five categories. They are as follows:

1. Visualization - Pay consideration to the articles that exist not apparent.
2. Image polishing and re-establishment - Towards improve the quality of an image.
3. Measurement of pattern - Regulates the size of discrete things in an image.
4. Image acknowledgement - Categorize things in an image.

3.1.2 Types of Images:

An image can be represented as a two-dimensional function, where the spatial coordinates (x, y) define the position, and the intensity (or gray level) at each point represents the amplitude of the image at that location. This intensity indicates the brightness or darkness of the pixel.

When the values of this function are made discrete—both in spatial coordinates and amplitude—the result is called a digital image. The process of using a digital computer to manipulate or analyze such images is known as digital image processing.

In a digital image, each pixel corresponds to a specific coordinate in the 2D plane and holds one or more numerical values that describe its intensity or color. These values are referred to as samples.

In programming languages like Python, image processing typically deals with two main types of images:

- Binary images, where each pixel is either black or white.
- Grayscale images, where pixel values represent varying shades of gray.

3.1.2.1 Binary Image:

This type of images are stored in a logical array. It are also referred to as bi-level or two-level images. (This idea is known as black and white, or B&W). Certain i/o equipment, for instance laser printers and computer displays, can only be utilised with bi-level images.



Fig.3.1 Binary Image

A binary image will be shown in the figure 3. In this, each pixel undertakes one of only two isolated principles: 1 or 0.

3.1.2.2 Gray Scale Images:

This type of image is frequently collected of ranging from dark at the lowest concentration to white at the highest. This is because the value of each pixel represents a particular trial. Grayscale images are frequently produced when measuring the intensity of light at each pixel in a single band of the electromagnetic spectrum. (For example, infrared, visible light, ultraviolet, and so on). Grey scale images are frequently created by gauging the strength of images at towards each pixel. It is for visual displays are commonly saved with 8 bit per sample pixel in order to record 256 intensities, or shades of grey. A grayscale copy is displayed in the figure 3.2.



Fig.3.2 Gray scale Image

3.1.3 Kinds of Image Processing:

The two types of image processing used are:

- Analog
- Digital

3.1.3.1 Analog Image Processing:

Analog Image Processing refers to the use of analog methods to work with two-dimensional analog signals, particularly in the fields of image science and computer science. Unlike digital image processing, these techniques are applied to physical formats such as printed images or photographs.

In analog or visual image processing, analysts rely on interpretive techniques to examine and evaluate images. The effectiveness of this process depends heavily on the analyst's expertise and the specific area being studied.

Visual interpretation plays a key role in analog image processing, where human observation and judgment are essential. Analysts often combine their subject knowledge with limited data to extract meaningful insights from images.

3.1.3.2 Digital Image Processing:

Computer-based numerical image modification is made possible through digital processing techniques. Raw data collected from satellite imaging sensors often contains defects, which must be corrected through multiple processing steps to achieve accurate and usable results.

When using digital methods, all data typically undergoes three main stages:

1. Pre-processing
2. Enhancement and display
3. Information extraction

In this process, digital computers are employed to handle and process the image. A scanner-digitizer is used to convert the physical image into a digital format, which can then be manipulated.

Digital image processing is defined as the execution of a series of operations on the numerical representation of an object to produce the desired outcome. Starting from an original image, it generates a modified version, effectively transforming the input into a new and enhanced image.

3.2 METHODOLOGY FOR IMAGE PROCESSING:

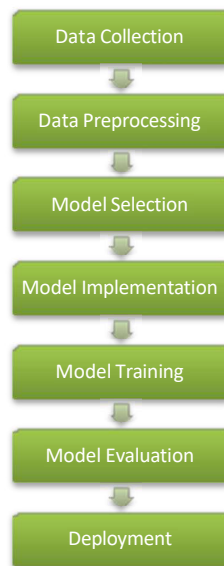


Fig.3.3 Flow Chart

Steps involved:

1. Data Pre-processing:

The initial stage is data preprocessing.

Data preprocessing involves several steps, including: 800 images of leaves from the classes Diseased and Healthy are loaded into the system during training. A Python image processing library called OpenCV is used to convert images from RGB to BGR format. Since OpenCV only accepts images in BGR color format, RGB images must be converted accordingly. When converting from BGR to HSV, luma (image intensity) is separated from chroma (color information).

1. Image Segmentation:

Image segmentation is the process of dividing a digital image into multiple segments, also referred to as image regions or objects (sets of pixels). More precisely, image segmentation assigns a label to each pixel in the image so that pixels with the same label share specific characteristics. It is essential in this context to separate the leaf image from the background and to perform color extraction.

2. Feature Mining:

Feature mining is the process of transforming raw data into mathematical features that can be used while retaining the important information from the original dataset. Compared to applying machine learning directly on raw data, this approach produces better results. To extract comprehensive features from the image, three types of feature descriptors are used:

- a) Color: Color Histogram
- b) Shape: Hu Moments
- c) Texture: Haralick Texture

Once the features are extracted, they are combined into a single feature set.

3. Model Training:

For a better understanding of the device, the labels are numerically encoded depending on the photographs in the folder. Two sections of the dataset have been separated. They are divided 80/20 as the training set and the testing set. Data pre-processing should include feature scaling so that it can manage extremely variable magnitudes. Extraction of Features: The features are

taken from the photos and saved in an HDF5 file. Modeling: The following five machine learning models are used to train the model:

- a) Random Forest
- b) Logistic Regression
- c) KNN
- d) Naive Bayes
- e) SVM

Once model is trained , The 10 k fold cross validation technique is now being recycled to authorize the model.

2. Prediction and Testing:

The prediction is to be done whether the leaf is Diseased or Healthy. Prediction is done using confusion matrix which determines accuracy,precision,f1 score and recall for the applied algorithms.

3.2.1 Confusion Matrix

A matrix is active to evaluate how well a given set of test data performs the categorization models. Only after the true values of the test data are known can it be resolute. It is also referred to as an error matrix since it displays the errors in the model performance as a medium. The subsequent list of Confusion matrix features includes:

Real values are the real values for the in case data, whereas projected values are the values that the model expects.

- True Negative: The typical projected no, and the actual or else real value also indicated no.
- True Positive: The model correctly predicted yes, and the outcome matched that prediction.
- False Negative: This error is now and then referred to as a Type-II error and occurs when the model forecast no but the actual value was yes.
- False Positive: Even though the model predictable Yes, the actual result was No. Another name for it is a Type-I error.

Table 3.1 Confusion Matrix

| | | Actual Values | |
|------------------|--------------|---------------|--------------|
| | | Positive (1) | Negative (0) |
| Predicted Values | Positive (1) | TP | FP |
| | Negative (0) | FN | TN |

3.2.2 Evaluation Parameters

Accuracy:

The proportion of correctly classified data instances over all data instances is known as accuracy.

$$\text{Accuracy} = \frac{\text{TN} + \text{TP}}{\text{TN} + \text{FP} + \text{TP} + \text{FN}}$$

Precision:

A good precision should rather be 1. Only when the numerator and denominator are equal, or when $\text{TP} = \text{TP} + \text{FP}$, does precision become 1, which also implies that FP is zero.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

Recall:

Recall should ideally be 1 for a good classifier. Recall becomes 1 only when the numerator and denominator are equal i.e $\text{TP} = \text{TP} + \text{FN}$, this also means FN is zero. As FN increases the value of denominator becomes superior than the numerator and recall value decreases.

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

F1-Score:

$$\text{F1-Score} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

F1 Score becomes 1 only when precision and recall are both 1. F1 score converts high only when both precision and recall are high. F1-score is the harmonic mean

of precision and recall and is a better measure than accuracy.

3.2.3 Introduction to Machine Learning:

Machine learning is a branch of artificial intelligence (AI) that focuses on identifying patterns within data and building models that are intuitive and practical for users. While it is a subfield of computer science, machine learning differs significantly from traditional computational methods.

In conventional computing, systems operate based on explicitly programmed instructions to perform calculations or solve problems. In contrast, machine learning allows computers to learn from data inputs, using statistical analysis to produce outputs that fall within a defined range.

Machine learning tasks are generally divided into broad categories, based on how the system learns from data or the type of feedback it receives during the learning process. The following are the three main groups into which machine learning implementations fall: -

- a) Supervised Learning
- b) Unsupervised Learning
- c) Reinforcement Learning

a) Supervised Learning:

This technique involves building a mathematical model based on a dataset that includes both input and corresponding output values. The dataset, known as training data, is composed of multiple examples. Each example is represented as a feature vector (an array of values), and the full set of training data is organized into a matrix. Supervised learning algorithms use this data to develop a function that can predict outputs for new, unseen inputs. This is done through repeated optimization of an objective function. A well-trained model can accurately estimate outputs for inputs that were not part of the training data. When an algorithm improves its prediction accuracy over time, it is considered to have learned how to perform the task.

b) Unsupervised Learning:

These algorithms work with datasets that contain only input values, identifying patterns or groupings—such as clusters—within the data. Since the data is not labeled, categorized, or classified, the algorithms learn without direct supervision. Rather than learning from feedback, they detect patterns and respond to the presence or absence of these structures in new data.

Unsupervised learning is commonly used for tasks like density estimation, such as determining a

probability mass function. It also plays a significant role in data summarization and interpretation. Clustering, a form of classification, involves dividing a set of observations into subgroups (called clusters) so that items within the same group are similar according to certain criteria, while items in different groups are dissimilar. Different clustering algorithms make varying assumptions about data structure, often using similarity measures. These measures might include **internal cohesion** (similarity among members of a cluster) and **separation** (differences between clusters). Some methods also rely on estimated data density or connectivity to form clusters.

c) Reinforcement Learning:

Reinforcement learning is a type of machine learning that focuses on training software agents to make decisions in an environment in order to maximize a total reward over time. The environment is often modeled as a Markov Decision Process (MDP). While many reinforcement learning methods are based on dynamic programming, they are especially useful when the exact mathematical model of the MDP is unknown or too complex to define. These algorithms learn optimal behavior through interaction with the environment rather than relying on predefined models. Reinforcement learning is widely used in applications such as autonomous driving and training agents to compete in games.

3.2.4 Algorithms used:

3.2.4.1 Random forest classifier:

RF is a eminent machine learning process from the directed learning technique. It can be applied to both sorting and deterioration problems in machine learning. It indeed is created on the perception of cooperative learning, which is a process that involves combining multiple classifiers to solve a compound problematic and expand the typical act.

"RF is a classifier that comprises a total of decision trees on amenities for low of the given and takes the mean to improve the analytical skill of that dataset," as the name indicates. Instead of relying on a particular decision tree, the rf takes the predictions from each tree and envisages the ultimate output constructed on the mainstream division of forecasts.

Random Forest is effective. RF mechanism in two-phase major is to create by merging N

decisiontree, and second is to make calculations for each tree shaped in the first segment.

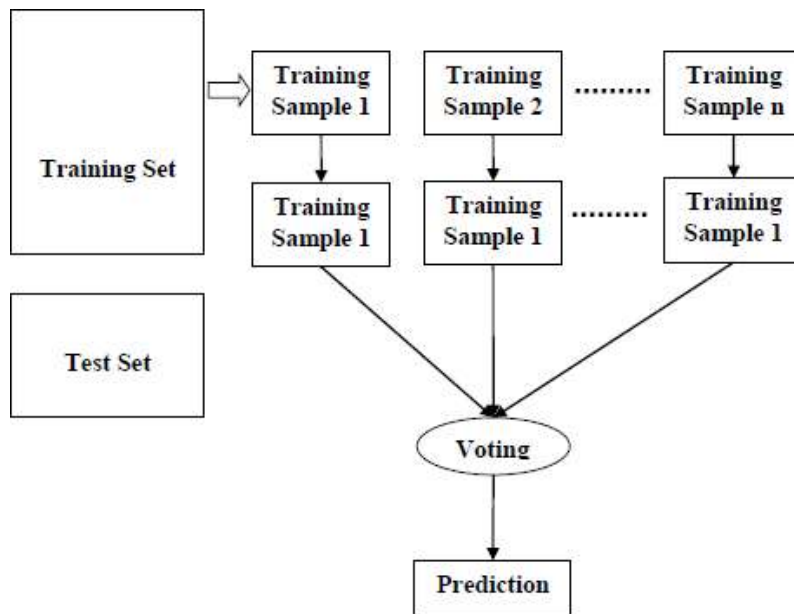


Fig.3.4 Working of Random Forest Algorithm

The method can be illuminated in the below stages and figure:

Step 1: Pick K data points at arbitrary from the preparation set.

Step 2: Build decision trees for the chosen data points .

Step 3: Determine the size N for the numeral of decision trees you want to construct.

Step 4: Reverse steps 1 and 2.

Step 5: Find the conventions of every decision tree for new records arguments and allot the new data points to the segment that obtains the peak divisions.

3.2.4.2 Logistic Regression:

- It is a prevalent Machine Learning system that falls under the Machine learning Learning method. It is used to estimate the definite reliant on variable from a set of self-determining variables.
- The output of a category independent variable is prophesied by logistic regression. As a result, the outcome must be categorical or discrete. It can be 0 or 1, true or False, and so on, but instead of giving the exact values , it gives the probabilistic values that drop between 0 and 1.

- In the field of logistics, In Logistic regression, as a substitute of appropriate a regression line, we fit an "S" shaped logistic function, which foresees two extreme values (0 or 1).

3.2.4.3 Support vector machine(SVM):

- Support Vector Machine (SVM) is a prominent Supervised Learning algorithm that is used for Classification and Regression problems. However, it is chiefly used in Machine Learning for Classification problems.
- The SVM algorithm's goal is to find the top line or decision boundary for categorising n-dimensional space so that we can easily place new data facts in the accurate class in the future. A hyperplane is the best decision frontier.

3.2.4.4 K-Nearest Neighbour:

- It is a modest Machine Learning procedure that uses the Supervised Learning system.
- It adopts correlation between the fresh case/data and existing cases and places the new case in the classification that is most similar to the existing categories.
- It supplies all accessible data and uses similarity to catalog new data points. This means when new data is generated, it can be rapidly secret into a well-suited category using this algorithm.
- This technique can be located used for both deterioration and classification, but it is most frequently recycled for organization difficulties.

3.2.4.5 Naive Bayes:

- This approach is a managed knowledge system that practices the Bayes theorem to solve arrangement problems.
- It is mainly used in text labeling with a large training dataset.
- This is a simple and useful Classification algorithm that aids in the expansion of firm machine learning prototypes adept of making quick predictions.
- It is a classification algorithm, which funds it expects created on an object's probability.

3.2.5 TECHNOLOGIES USED

3.2.5.1 Python:

Python is an advanced, universal-persistence software design language that is interpreted. Python's core idea emphasises encryption readability, as evidenced by its extensive use of indentation. Its language structures and object-oriented approach are planned to support systems analyst in writing clear, rational

It is dynamically typed. It is compatible with a selection of computing paradigms, counting structured (primarily procedural), object-oriented, and functional programming. Python is frequently mentioned to as a "strategies included" language due to its general typical library. Guido van Rossum began developing Python as a beneficiary to the ABC programming language and it was first released in 1991 as Python 0.9.0. Python 2.0 was released in 2000, and it included innovative topographies like list comprehensions and a trash collection system that used orientation totaling. Python 3.0, released in 2008, was a significant revision of the language that is not completely completely well-suited, and much Version 2 code does not run basic on Python 3. Python 2 was phased out with version 2.7.18 in 2020.

Python is a language for programming that supports multiple paradigms. Various of its features sustenance functional programming and aspect-oriented programming (including metaprogramming and metaobjects (magic methods)). Object-oriented programming and structured programming are fully braced. Many other frameworks, such as design by contract and logic programming, are maintained by delays.

3.2.5.2 Libraries

Python's large ordinary collection, which is widely viewed as one of its greatest assets, provides tools suitable for a inclusive series of responsibilities. Many regular formats and protocols, such as HTTP, are stayed for Internet-facing applications. It has modules for developing user interfaces with graphics, connecting to interactive databases, generating pseudorandom numbers, arithmetic with arbitrarily chosen decimals, attempting to manipulate regular expressions, and unit analysis.

Some parts of the source file are covered by specifications (for example, the Web Application Gateway Interface (WSGI) integration follows PEP), but most modules are not. Their code, interior documentation, and test suites define them. However, because the majority of the codebase is merge Python code, only very little segments require modification or revising for alternative.

3.2.5.3 Pandas

It is a computer software library printed for the Python programming language for data handling and examination. In individual, it offers data gatherings and acts for handling geometric tables and time series. It is permitted software unrestricted under the three clause BSD license. The title is consequent from the term "board files", an econometrics tenure for

data sets that include remarks over numerous epochs for the same individuals. Its name is a play on the phrase "Python data analysis" itself. Wes McKinney started construction what would turn out to be pandas at AQR Capital while he was an investigator from 2007 to 2010.

Features:

- Data Frame object for data handling with unified indexing.
- Apparatuses for understanding and characters data between in-memory data edifices and different fileformats.
- Data alignment and integrated handling of missing data.
- Reforming and twisting of data sets.
- Statistics organization column insertion and deletion.
- Assemblage by appliance allowing split-apply-association operations on data sets.

3.2.5.4 Scikit-learn

Scikit-learn is a free software machine learning reference library for the Python programming language.

It features various classification, regression and gathering algorithms, including support vector machines, random forests, k-means and DBSCAN, and is intended to interoperate with the Python geometric and scientific libraries NumPy and SciPy.

Scikit-learn is written mainly in Python and uses NumPy comprehensively for high-performance linear algebra and arrangement operations. Furthermore, some essential algorithms are written in python to advance performance.

3.2.5.5 NumPy

NumPy is a Python library that augments backing for large, multidimensional arrays and matrices, as well as a huge gathering of high-level mathematical operations to operate on these arrays. Jim Hugunin created NumPy, Numeric, with help from various other developers. Scientist by providing features of the able to compete Num array into Numeric

with major improvements. NumPy is offers a bunch with numerous contributors.

Features:

- NumPy is aimed at the Python C Python mention implementation, which is a non-enhancing bytecode interpreter.
- Algorithms written for this description of Python are frequently much slower than compiled equivalents.
- NumPy addresses the slowness issue in part by so long as multidimensional arrays as well as functions and hands that operate proficiently on arrays without the need for rewriting.

3.2.5.6 Seaborn

It is a data visualisation library. It offers a high-level interface for forming visually pleasing and informative quantitative graphics.

Seaborn aids in data investigation and command. Its plotting functions operate on data structures and arrays containing entire datasets, performing the necessary concept mapping and statistical aggregation internally to generate informative plots.

Its dataset-oriented, declarative API allows you to focus on what the individual aspects of the plots mean rather than the details of how to draw them. There is no single best method for visualising data. Different plots are better suited to answering different questions. Seaborn's consistent. When statistical values are estimated, seaborn uses bootstrapping to compute confidence intervals and draw error bars representing the estimate's uncertainty.

Statistical analyses require knowledge about the distribution of variables in your dataset. The seaborn function `displot()` supports several approaches to visualizing distributions.

3.3 RESULTS

3.3.1 Confusion matrix Result

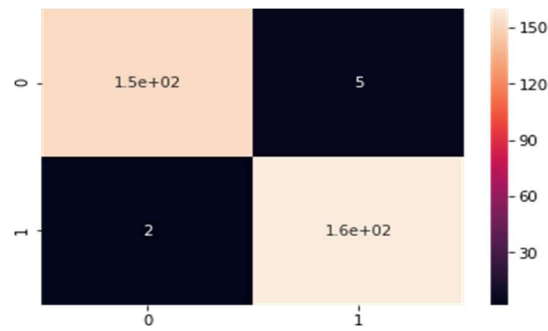


Fig.3.5 Plot of Confusion Matrix

In the above confusion matrix, 5 images are predicted as false negative and 2 images are predicted as false positive.

3.3.2 Comparison Table for Different Machine Learning Algorithms

| Parameters | Random forest | Logistic Regression | KNN | Naive Bayes | SVM |
|------------|---------------|---------------------|--------|-------------|------|
| Precision | 0.98 | 0.88 | 0.96 | 0.88 | 0.94 |
| Recall | 0.98 | 0.86 | 0.96 | 0.86 | 0.94 |
| f1-score | 0.98 | 0.86 | 0.96 | 0.86 | 0.94 |
| Accuracy | 0.9812 | 0.9265 | 0.9562 | 0.8578 | 0.94 |

Table 3.2 Comparison Table For Different Machine Learning Algorithms
comparison table for different algorithms such as RF, LR, SVM, Naive Bayes, KNN
with parameters like Precision, Recall, F1-score, Accuracy.

Random Forest Algorithm results better accuracy of all algorithms with 98%.

3.3.3 Boxplot Comparison for Different Algorithms

Models are trained on different plant leaves. The machine is given 1600 images of leaves for each class Diseased and Healthy in order to train different models.

A Boxplot comparison is plotted for different machine learning techniques. The figure

shownbelow depicts the accuracy vs different machine learning algorithms.

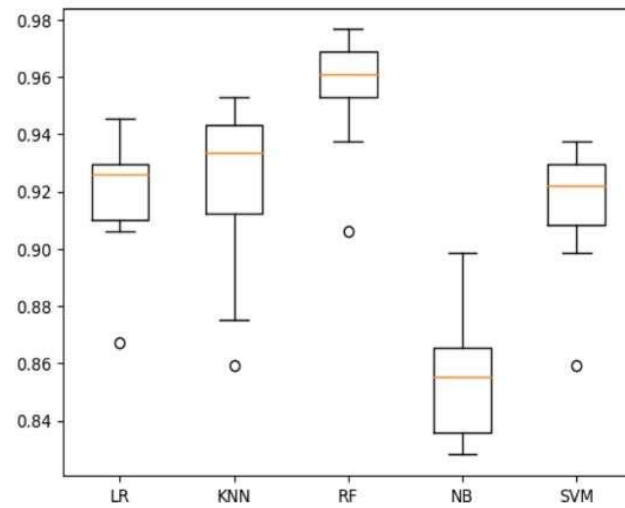


Fig.3.6 Boxplot Comparison

3.3.4 Results of Plant Leaf Images:

Loading Original Image- For training, the machine is given a total of 1600 images of leaves, This 1600 images are divided into two classes namely Healthy and Diseased. The image is changed from RGB format to BGR.

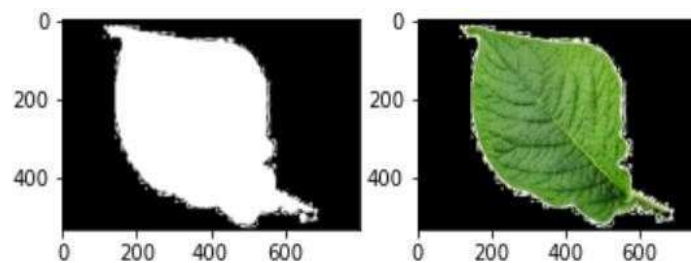


Fig.3.7 Healthy Leaf

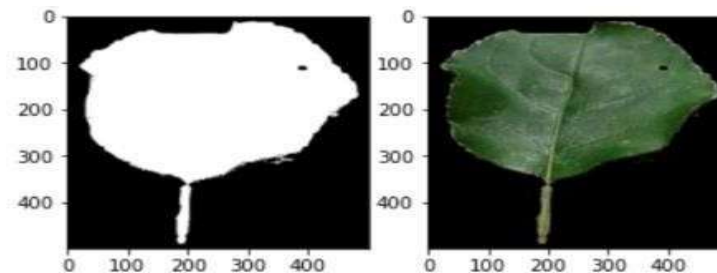


Fig.3.8 Diseased(Low)

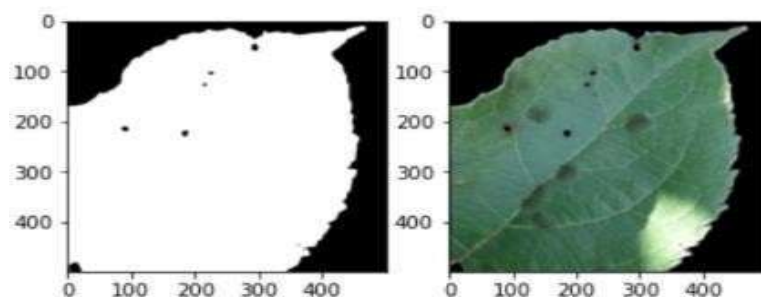


Fig.3.9 Diseased(Medium)

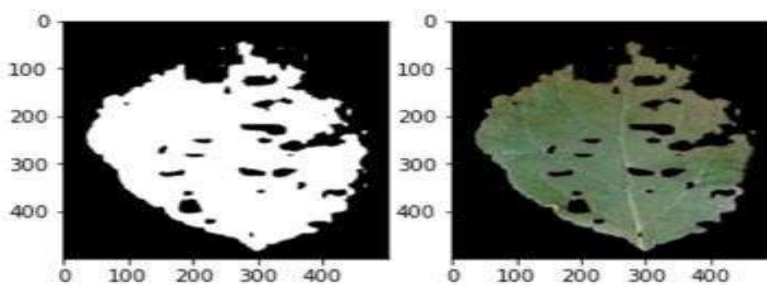


Fig.3.10 Diseased(High)

CHAPTER 4

METHODOLOGY

4.1 Introduction:

There are approximately 570 million farms around the world, and over 90% of the global population has some connection to agriculture. Farmers are responsible for producing a significant portion of the world's food—about 80%. However, crop diseases remain a major cause of reduced productivity and quality in broadacre farming. These diseases may result from bacteria, viruses, nematodes, and other pathogens, affecting both the above-ground and underground parts of plants.

Various factors, such as climate change and declining pollinator populations, pose additional threats to food quality and quantity, directly impacting farmers' livelihoods. In their efforts to combat plant diseases, farmers often incur losses in money, time, and other resources. With the ongoing environmental changes, early recognition and diagnosis of these plant diseases is more critical than ever.

Identifying the names of diseases and understanding how to control them remain persistent challenges in agriculture. Some plant diseases are difficult to diagnose based solely on visible symptoms. Therefore, accurate detection is vital not only for disease prevention but also for ensuring the long-term sustainability of the agricultural sector.

To address these challenges, technological innovation is needed—especially through the integration of automation, information, and communication technologies in both farming and commercialization. Without these advancements, reliance on imported food may increase, driving up costs and potentially impacting public health.

Interventions such as the use of fungicides, disease-specific chemicals, and pesticides can help manage crop diseases—especially if implemented early. Having an early warning system for plant illnesses would significantly enhance the effectiveness of such applications and boost agricultural productivity.

In recent years, smartphone-assisted disease detection has become a reality, thanks to the global spread of smartphones and advancements in machine learning. Machine learning and deep learning algorithms are increasingly being applied to improve the accuracy and efficiency of plant disease detection. Numerous studies have explored various methods and technologies in this domain, including traditional machine learning techniques widely used in image processing and plant disease identification.

4.1.1 Convolutional Neural Networks:

Image classification classifies the photos as a dataset of unprocessed pixels and defines the images as objects. The CNN model of neural networks makes us to derive image with more accurate representations. Contrary to this, image recognition techniques, we specify the picture characteristics ourselves, CNN begins with the image's pixel data which is raw data, trains the model is trained, and then extracts many number of features automatically for improved categorization.

4.1.2 Transfer learning:

It uses a model which was trained for single job as the foundation for a model on another. When the other task is identical to the previous task or when there was a dearth of data for the second task, this can be helpful. The model can learn more quickly and efficiently on the second challenge by starting with the learnt features from the first task. Because the model will have already picked up broad features that are likely to be helpful in the second task, this can also aid in preventing overfitting. It can sometimes be unavailable or impractical to train a network from starting because it needs a lot of data information, power, and time. It is a method for applying to the neural networks which are trained previously (such Alexnet, Inception net, and VGG16) to new tasks by altering the final classification layer. As we move deeper in the network, layers learn towards learn different

patterns more specialised to the work which they were they were trained. The earliest levels learn relatively general features. The models which were trained before are on a vast Compared to a neural network created from scratch, images can learn these abstract properties more effectively.

Block Diagram of Transfer Learning:

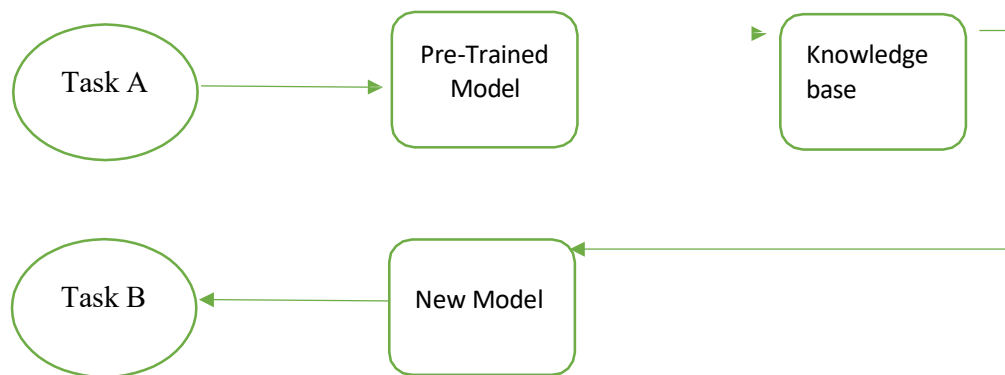


Fig.4.1 Transfer Learning

4.1.2.1 Need of Transfer learning:

Transfer learning is required in machine learning when a pre-trained model is used to solve a new problem that is different but related to the original problem for which the model was trained. This technique is particularly useful when dealing with complex, hard-to-solve problems that require a significant amount of training data and computational resources. By leveraging the knowledge and expertise gained from pre-training on a large dataset, transfer learning can significantly reduce the amount of data and time required to train a new model for the target.

4.1.2.2 Advantages of transfer learning:

- **Quicken the learning process:** By employing a model that has already been trained, the model may learn the another task more accurately, fastly and efficiently since it is already familiar with the characteristics and patterns in the data.
- **Improved performance:** Because the model can use the first task's information to its advantage, transfer learning can result in improved performance on the second task.

- Handling a tiny dataset: Since the model has already learnt general features that are likely to be helpful in the second task, transfer learning can assist prevent overfitting when there is a dearth of data for the second task.

4.1.3 Resnet50:

Residual Network, often known as ResNet, is a network that supports ongoing learning. Resnet50 stands for Residual Network with 50 layers. This network's pretrained version was trained on greater than 10 lakh photos from the database and is used as the foundation for a variety of tasks. It was the victor of the 2015 challenge on imagenet. In comparison to models like CNN, InceptionV3, MobileNetV2, GoogLeNet, and others, ResNet50 has a higher accuracy. The best model for identifying plant diseases is Resnet50, which has been demonstrated in numerous journals and studies. This Resnet50 model is also one of the Top 4 Pre-trained Models for picture classification. The ResNet50 model is a transfer learning model, which implies that it doesn't require intensive computing because it takes into account these three factors.

The main modernization with ResNet was that it empowered us to train very deep neural systems with more than one hundred fifty films. A main downside of convolutional networks is the problem of "Vanishing Gradient ". This pointedly diminishes throughout the technique called backpropagation, consequently masses will alteration slightly. Therefore this working to grow round that. It contains SKIP CONNECTION method.

Connection(skip): Connecting the convolutional block's output to the initial input. A direct link known as a skip connection neglects the some layers of the prototypical. The outcome is not the equivalent because of this connection. Without this connection, input X is multiplied with the weights of the layer trailed by adding a term called bias. This technique made the network to evade certain layers and reuse features from former stages of the network, which helps avoid the issue of vanishing gradients and allows for deeper and more accurate training of the model. Essentially, skip connections help confirm that important features are conserved and passed through the network, even as the data travels through multiple layers of the model.

4.1.3.1 Residual block:

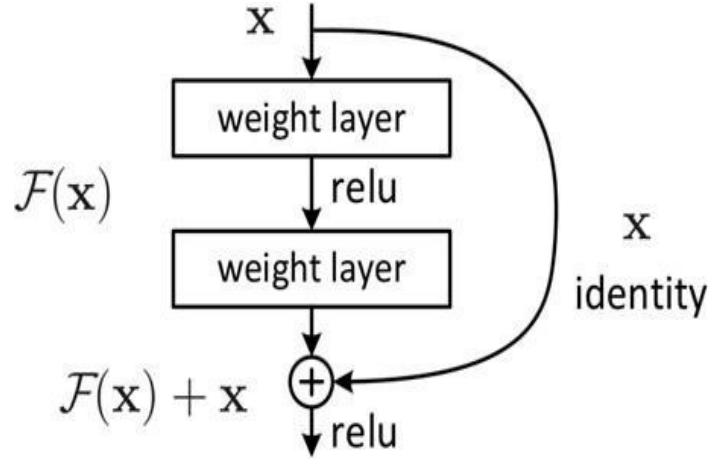


Fig.4.2 Residual block

The skip connection is the most vital idea in play here, as can be seen in Fig 4. 1. In essence, a skip connection is an uniqueness mapping in which the outcome of one film is added to the previous layer input. A straightforward residual function can be summarised like this using the preceding figure:

When x - is applied input and Here $f(x)$ -the layer's outcome, the block's output can be articulated as like this:

$$y = f(x) + x$$

This is the furthestmost straightforward characterization of a block. Now, there can be approximately circumstances where the outcome from the film and the input which is identical have diverse proportions. For kindly, if we took a CNN where we distinguish that subsequently intricacy operation, the size of the input is condensed (proportionally), then adding effort to it is a problematic. So, what here can be finished is that in the connection (skip), we improve some procedure or meaning (in this case intricacy operation) such that the input is reformed or constituted to the obligatory magnitudes.

So, the description can be modernized as follows:

$$y = f(x, \{w_i\}) + w_s * x$$

At this point, Wi is the constraints specified to the layers,ws tenure can be done with convinced intricacy conformation to make input and output magnitudes identical.

4.1.3.2 Block diagram of Resnet50:

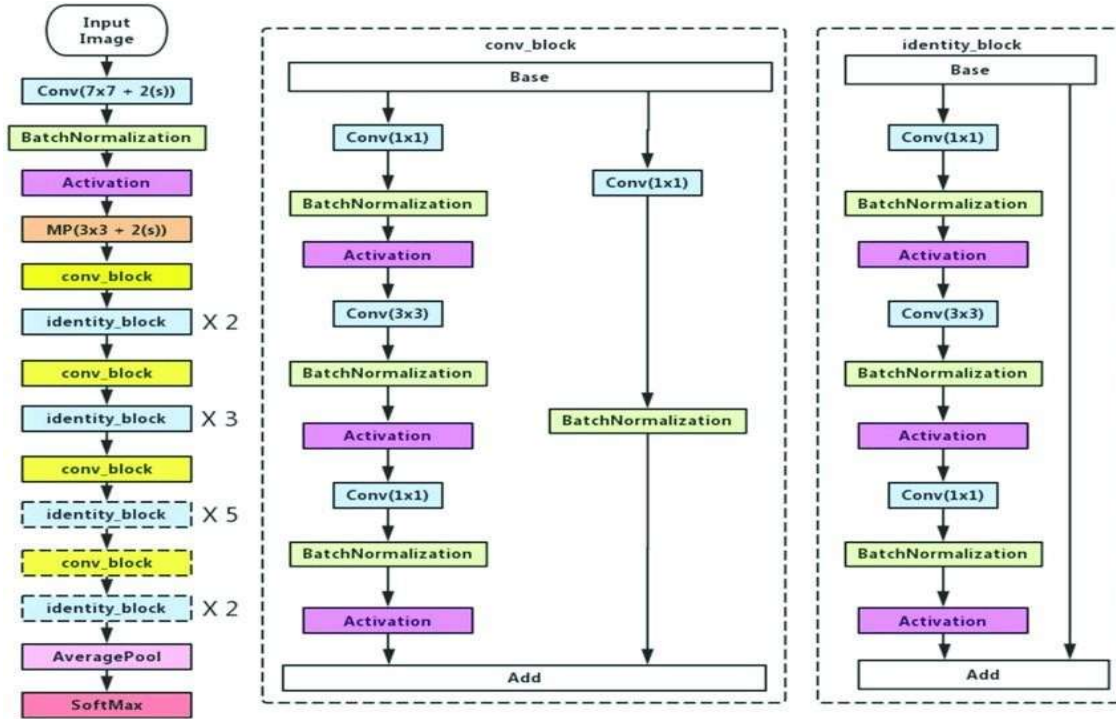


Fig.4.3 Resnet50 Architecture

The 50-layer ResNet construction embraces the ensuing elements, as shown in the below:

- A seven*seven kernel intricacy together with sixty four additional kernels with a two-stride.
- A maximum pool layer with a two stride.
- Nine additional layers-three*three, sixty four kernel intricacy, additional with one*one, sixty four kernels, and a third with one*one, two hundred fifty six kernels. These three layers are recurrent three epochs.
- Twelve additional layers with one*one, one hundred kernels, three*three , one hundred twenty eight kernels, and one*one, five hundred twelve kernels, recapitulated four many times.
- Eighteen extra films with one*one, two hundred fifty six cores, and two cores three*three, two hundred fifty six and one*one, one thousand twenty four, recapitulated six times.

- Nine additional films with one*one, five hundred twelve cores, three*three , five hundred twelve cores, and one*one, two thousand forty eighty cores recapitulated three times.

There are 2 core categories of lumps are castoff in a Res Net, reliant chiefly on if the in/outcome magnitudes are the equivalent and diverse.

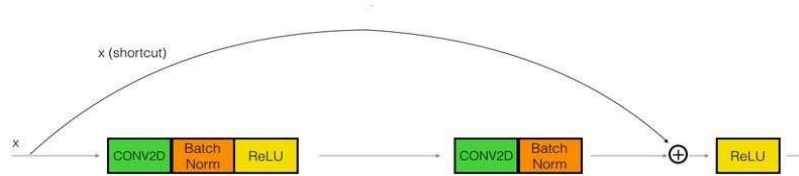


Fig.4.4 Identity Block

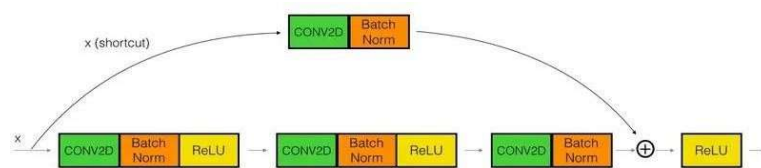


Fig.4.5 Convolution Block

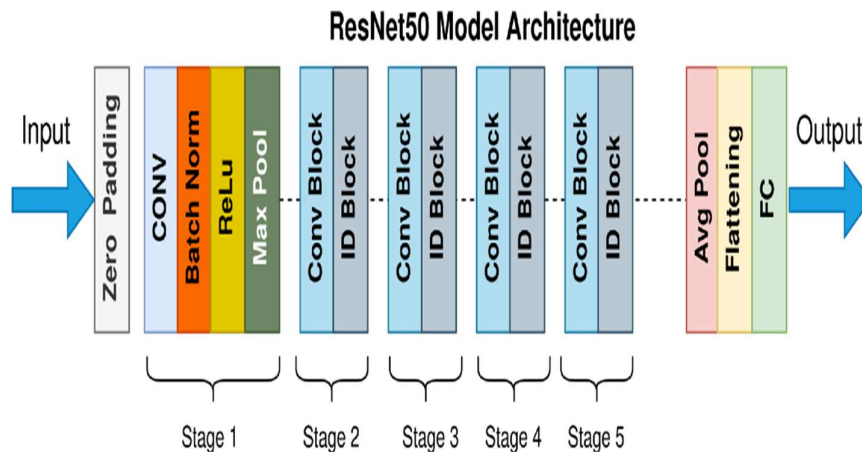


Fig.4.6 Resnet50 Architecture

1. Identity Block: This lump is typical lump used in Residual Nets and resembles to the case wherever the input stimulation has the identical measurement as the productivity stimulation.

2. Convolutional Block: We can use this kind of lump when the input,output magnitudes do not be similar. The transformation with the similar block is that nearby a CONVOLUTION layer which is two dimensional in the shortest and easiest pathway.

The ResNet-50 prototypical comprises of five steps individually with a intricacy and Identity block. Each intricacy lump has three intricacy layers and individually identical lump likewise has three intricacy layers. It contains 25 million trainable criterion.

4.1.3.3 Special characteristics of ResNet-50:

The construction of ResNet-50 is based on the perception shown, with one significant exemption. The holdup element is castoff in the fifty layer Res Net. A holdup enduring lump, likewise referred to as a "holdup", uses eleven intricacies to expurgated on the count of strictures and matrix developments. It styles individual layers training suggestively quicker. As a substitute by means of a load of 2 steps, it services 3 number of layers. Kkip connection that is added to the ResNet design greatly improved the presentation of the network with various layers. ResNets are fundamentally just other networks with a few minor amendments. The construction follows the same purposeful steps as CNN or other systems, but an spare step is encompassed to address apprehensions like the problem like vanishing gradient, among others.

4.1.3.4 Compensations of ResNet:

- Systems with bulky quantity (even 1000s) of layers can be trained effortlessly deprived of cumulative the training error proportion.
- It helps in embark upon the vanishing gradient delinquent using identity mapping.

4.2 Methodology:

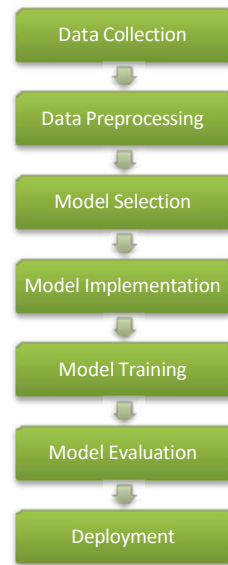


Fig.4.7 Methodology

4.2.1 Dataset:

4.2.1.1 Training Dataset:

It is castoff to train and suitable the ML prototypical, is the principal subcategory of the unique dataset. The ML procedures learns how to do prophecies for the specified action, data is first abounding into them for training purpose. Whether we are consuming supervised learning or unsupervised, the training data varies.

In unsupervised learning, inputs do not tagged with the apposite outputs, henceforth the training data contains not labeled data points. In order to yield prophecies, models must citation patterns from the afforded datasets which are training. In otherside, labels are encompassed in the training data for administered learning in command to help the model be trained and prophecies made.

The type of training data that we offer to the model is exceedingly accountable for the model's accuracy and prophecy ability. It means that the recovering the superiority of the training data, the improved will be the presentation of the model. Training data is

roughly in excess of or equal to 50% of the whole data for an ML development or deep learning project.

4.2.1.2 Testing Dataset:

As soon as the prototypical has been accomplished using the dataset(training), it is stint to assessment it consuming the test dataset. This dataset evaluates the prototypical's recital and assurances that it can oversimplify efficaciously to newfangled or unmapped dataset. The assessment dataset is a unlike subcategory of the original data from the training set of data. When the model working out is finished, it exploits it as a standard because it has some comparable geographies and a alike class distribution of probability. A glowing-systematized dataset called test data affords evidence for each type of setting the prototypical strength happenstance in the authentic creation. Characteristically, the set of data which is test varieties up 25–30% of the inclusive innovative dataset.

At this time, we may also scrutinize and disparity the challenging accuracy with the training accuracy, or, more unambiguously, the exactness of our model when pragmatic to the test dataset in assessment to the learning set of data. The model is well-thought-out to have overfitted if its accurateness on training data is sophisticated its accurateness on testing data. In a point to make exact assumptions, the testing data must moreover abundantly or moderately resemble to the innovative dataset.

4.2.1.3 Requirement of Splitting Dataset:

One of kind decisive phases in data pre-processing is distributing the set of data into learn and examine sets. By making use of this, we may augment the recital of our prototypical and hence deliver improved obviousness. We can contemplate of it as if we accomplished our model with a learning set of data and then established it with a test dataset that is utterly dissimilar from the training dataset, at which point our model would be impotent to diagnose the associations amongst the topographies. As a consequence, the model's concert will agonize if it is trained and tested on two dissimilar datasets. Thus, it is decisive to rift a dataset into a train and a test set.

In this way, we can straightforwardly appraise the presentation of our prototypical. Such as, if it accomplishes glowing with the learning data, but does not accomplish glowing with the test set of data, then it is appraised that the prototypical may be overfitted. For excruciating the dataset, we can use the traintestsplint purpose.

4.2.1.4 Difference between Training Data and Testing Data:

- The main discrepancy between working out data and challenging data is that working out data is the subdivision of innovative data that is used to train the prototypical, whereas testing data is used to check the accurateness of the prototypical.
- The training dataset is normally grander in proportions associated to the testing dataset. The universal fractions of piercing train and test are 80 and 20, 70 and 30, or 90 and 10.
- Training data is glowing acknowledged to the prototypical as it is castoff to train the prototypical, however testing data is like nor seen data /new-fangled data to the model.

4.3 Image Pre-Processing:

Pictures need to be treated before they can be used for prototypical learning and disturbance. This include, but is not inadequate to, variations in shade, magnitude, and direction. Pre-processing is ended to progress the picture's quality so we can analyze it more efficaciously. Through pre-processing, we are able to get rid of undesired falsifications and increase certain attributes that are crucial for the request we are evolving. Those potentials could amend based on the application. For software to effort correctly and deliver the required results, an image might be preprocessed.

4.3.1 Importance of Image Pre-Processing:

To formulate portrait statistics for prototypical input, preprocessing is obligatory. For illustration, convolutional neural networks' copiously associated layers necessitated that altogether the metaphors be in arrangements of the equivalent size.

If the input descriptions are very large, dwindling the proportions of these metaphors will greatly diminution the expanse of time looked-for to train the model without suggestively distressing model performance.

Even nevertheless geometric renovations of images (like alternation, clambering, and transformation) are regarded as as pre-processing procedures, the goalmouth of pre-processing is an upgrading of the twin data that overwhelms inadvertent falsifications or augments some image geographies decisive for succeeding processing.

4.3.2 Criterion of Image PreProcessing:

- For accomplishing recovering consequences from the practical prototypical in machine learning developments the presentation of the statistics has to be in an appropriate style. Some quantified model needs statistics in a quantified presentation, for specimen, Random Forest procedure does not provision insignificant morals, consequently to accomplish random forest procedure valueless morals have to be accomplished from the inventive underdone set of data.
- Additional characteristic is that the set of data might be configured in this type of way that more than singular ML and Deep Learning procedure are accomplished in single set of data, and best output of among is selected.
- Since metaphors occur in diverse presentations, i.e., ordinary, duplicate etc., we required to take to contemplation and homogenize before alimentating them into a system.

4.3.3 Image pre-processing techniques:

There are different image preprocessing techniques present. They are:

1. Grayscale alteration
2. Standardization
3. Data Expansion
4. Image tuning

1. Grayscale alteration: It is purely adapting metaphors from tinted to B&W. It was customarily cast off to diminish totaling complication in ML procedures. Meanwhile furthest depictions no requirement of color to be familiar, it is prudent to use scale i.e gray, which condenses the numeral of pixels in a twin, consequently, tumbling the calculations obligatory.

2. Standardization: Also mentioned to as statistics rescaling, it was the development of jutting picture data pixels to a defined previous range. It is frequently used on diverse set of data, and you hunger to standardize all to put on the equivalent procedures over them. Standardization is usually pragmatic to alter an pictures values ideals to a distinctive or high accustomed intellect.

It's assistances embrace:

- Impartiality transversely on all images – so mounting all pictures to an identical assortment of [zero, one]consents all descriptions to subsidize correspondingly to the full damage somewhat than when otherimages have more and less pixels assortments stretch durable and scrawny waste, individually.
- Delivers a typical lr- Subsequently more pixel metaphors necessitate a low lr and less pixel pictures high lr, scaling again benefits deliver a ordinary erudition rate.

3. Data expansion:It is the progression of making negligible adjustments to prevailing information to surge the situation assortment without assembling newfangled data.It is a procedure used for enlarging a dataset. Ordinary data expansion procedures includes parallel & erect spinning,harvesting, shearing, etc. Performing data augmentation helps in avoiding a neural network from learning extraneous features. This outcomes in better model presentation. Standard data expansion techniques include parallel & erect flipping, rotation, harvesting etc. There are 2 forms of expansion:

1. Off_line expansion - Castoff for insignificant datas. It is functional in the statistics pre-processing stage.
2. On-line expansion- Used for bulky datasets. It is customarily functional in instantaneous.

4. Standardizing images: Standardization is a technique used to ensure that images have similar sizes by rescaling them. It involves adjusting the data so that it has a average of zero. This process enhances the quality and consistency of the data..

4.3.4 Steps for loading custom dataset for Deep Learning Models:

1. Sweeping the image file:The arrangement of the dossier can be JPEG, BMP, etc.
2. Resizing the twin to contest the input size.
3. Renovate the metaphor pixels to datatype which is float.
4. Standardize the metaphor to have picture element values scrambled down between zero and one from zero to two hundred fifty five.
5. Metaphor data should be either a numpy array or a tensor object.

4.4 Data augmentation:

It is a procedure cast-off in ML and deep learning to surge the amount and variety of information available for training models. This technique involves creating new data from existing data by applying transformations such as flipping, rotating, scaling, or cropping. The goal of data augmentation is to advance the performance and oversimplification of representations by divulging them to a broader choice of training information.

By generating new data with different variations, data augmentation can help models learn to recognize objects from different perspectives and in different lighting conditions. This technique can also reduce overfitting by providing more diverse examples for the model to learn from.

Some common types of data expansion practices include random cropping, parallel flipping, upright flipping, zooming, and color shifting. These techniques can be applied towards several categories of information including pictures, text, and audio. In summary, data expansion is a powerful technique for improving the quality and robustness of machine learning models by creating additional training data with different variations and reducing overfitting.

4.4.1 Data(Image) Augmentation steps:

There is no ordinary set of expansion steps that are instantaneously going to progress the recital of the model on which you are employed. In fact, approximately expansions may have a negative impact on your model presentation.

1) Grayscale: Color fluctuations are an specimen of duplicate alterations that may be functional to all pictures (learn and examine) or arbitrarily rehabilitated in learning solitary as expansions. Normally, gray scaling is a color variation added to all metaphors. Although we may contemplate “high indication is constantly good; we should demonstrate the prototypical color,”we might perceive more appropriate prototypical recital when pictures are scaled. In addition, color is sometimes not as pertinent to a model. If you use greyscale, you don't need to apprehension about gathering images for every color of an object; your model will learn more wide-ranging features about an object that do not depend on color. Color pictures are deposited as different color values, whereas grayscale images are only deposited as a range. This means for CNNs, our model only requirements to labor with single matrix per single picture, not 3.

- 2) Random Flips: Flipping an image about its x- or y-axis pushes our model to recognise that an object does not necessarily have to be read from left to right or up to down; flipping may be irrational in order-dependent circumstances, such as deciphering text.
- 3) Random Rotations: Rotating a picture is especially critical when a prototypical will be utilised in a position that is not fixed, such as in a phone applications. It can be difficult because it, causes pixels to be damage on the boundaries of our pictures and needs mathematics to apprise any jumping cases.
- 4) Random brightness and exposure: Image brightness is adjusted to be arbitrarily sunnier and shadier is most useful if a prototypical must operate in a diversity conditions. It is critical to evaluate the highest and tiniest light levels in the apartment.

4.4.2 Image augmentation in Keras:

The Keras ImageDataGenerator class is a speedy and relaxed way to improve your images. It provides numerous expansion methods, such as standardization, rotation, brightness alterations, and many more. The fundamental advantage of the Keras ImageDataGenerator class is that it is designed to provide real-time data expansion. That is, while your model is still being trained, it is providing augmented images in real time. The ImageDataGenerator class provides the prototypical with new copies of the images at each epoch. It does not, however, include it in the original amount of photographs; instead, it merely distributes the updated images. If that were the case, the model would have seen the original images more than three times, which would clearly be excessive for our prototypical. Another advantage of ImageDataGenerator is that it devours fewer space. This is due to the reason, we would load entirely of the photos at once if we may not utilise this class. Nevertheless, when we practice it, we freight the pictures in consignments, which hoards a lot of retention.

4.4.2.1 Techniques present in Image_Data_Generator class:

1. Random alternation: Image alternation is a popular approach that allows the prototypical to develop insensitive to entity orientation. By fleeting a number in the range argument to the Image_Data_Generator class, you can randomly spin images by any degree between zero and 360degrees. When the image is spin, some values travel separate the picture, leaving an empty space that must be filled. You can fill this in a variety of ways, such as a continuous value or nearest values. The fill_mode option specifies this,

and the defaulting rate is nearest which merely swaps the unfilled region with the closest standards.

2. Random Shiftings: It is possible that the thing is not continually in the centre of the picture. To resolve this difficulty, we can move the pixel values of the picture parallelly or steeply by adding a constant number to all pixels. The `Image_Data_Generator` class accepts the parameter `shift height range` for a vertical picture shift and `breadth_shifting_range` for a horizontal image shift. If the value is a floating quantity, it indicates the proportion of the image's width or height to shift. Otherwise, if it is an integer number, the breadth or elevation are simply shifted by that more pixels.

3. Random Flips: It is also an excellent improvement technique that can be given to a wide variety of things. For flipping along the vertical or horizontal axis, the `ImageDataGenerator` class has options `horizontal_flip` and `vertical_flip`. This method moreover, should be tailored to the object in the photograph. Vertically overturning an vehicles ,it would be absurd associated to responsibility it with a object like a football.It quantified that, I'm working to jesting my picture in both ways only to show how the expansion is working.

4. Random Brightening: It aimlessly fluctuates the intensity of the picture. It is an excellent expansion strategy so that our object will not always be in optimal lighting circumstances. As a result, we must train our model on images obtained in a variety of lighting conditions.The `ImageDataGenerator` class's `brightness_range` option can be used to control the brightness. It selects a brightness shift value from a lean of 2 floating points. Standards fewer than one make the picture darker, where it is greater than one make it brighter.

5. Random Zooming: It randomly zooming in or out of the picture. For zooming, the `ImageDataGenerator` class's `zoom range` parameter allows a float value.We should deliver a lean of 2 quantities that define the lesser and higher bounds. If we provide a floating parameter, zooming shall take place in the range $[\text{one minus zoomrange}, \text{one plus zoomrange}]$.If the quantity is less than one, the image will zoom in. Any value larger than one causes the image to zoom out.

4.4.2.2 ImageDataGenerator methods:

There are a few methods in the ImageDataGenerator class which implement augmentation.

1. Flow_from_directory: While the neural network model is learning on the training data, the flow_from_directory() method allows you to read photos directly from the directory and enhance them. This technique wants photos from various classes to be present in different files but within the same parent folder. The following are some critical parameters of this method:

- directory: This is the pathway to the parent binder which contains the subfile for the dissimilar class images.
- Size of the target: Given the input image size.
- Mode of the color: rgb is set for coloring pictures else scaling gray.
- Scope of batch: Size of the groups of information.
- Mode of class: The setting is for 1 dimensional labels(binary), however the unconditional setting is for two-Done parameters.
- Seeding :Result is set to reproduction

2. Flowing_from_frames of data:

The flow-from-frame data() is additional prodigious technique in the Image_Data_Generator session that consents you to unswervingly supplement images by understanding its term and board value from a data frame. This originates very conveniently when we had many pictures deposited inside the equivalent file .This method also has a few parameters.

- Frames of data: It comprehends the picture labels and board values.
- manual: The pathway to the file that encompasses all the pictures.
- xcol: Support name in the Data Frame that had twin variable terms.
- ycol: Support name in the Data Frame that has the goal values.
- classmode: Fix to binary values is for one dimensional binary parameters whereas unconditional is for two dimensional 1-hot values.
- Size of target: Input pictures size.

- Size of batch: Batches size
- seed: Set to reproduce the result.

4.5 Features extraction using Resnet50:

Feature extraction is part of the dimensionality reduction process, which splits and condenses a starting collection of raw data into smaller, easier-to-manage groupings. As a result, processing will be more straightforward. The fact that these massive data sets contain a large number of diverse variables is their most important aspect. Processing these variables necessitates a significant amount of computational resources. Feature extraction helps to extract the best feature from enormous data sets in order to efficiently reduce the amount of data by selecting and combining variables into features. These features are straightforward to utilise while correctly and uniquely describing the actual data set.

4.5.1 Need of Feature extraction:

When you have a large data set and need to reduce the number of resources without sacrificing any critical or relevant information, the feature extraction technique comes in handy. Feature extraction assists in reducing the amount of redundant data in a data source. Finally, data reduction allows the prototypical to be built with low mechanism struggle while simultaneously growing learning space and generalisation processes in the ML method.

4.5.2 Building the Resnet50 model:

We can build the Resnet50 model by using the following default code by passing the necessary parameters based on the requirement.

```
tf.keras.applications.resnet50.ResNet50(
    include_top=False,
    weights='imagenet',
    input_tensor=None,
    input_shape=None,
    pooling=None,
    classes=14,
    **kwargs
)
```


Here `weights='imagenet'`, we can load the pretrained ImageNet weights. Otherwise `weights=None`. So initializing the model with random weights. `include_top=False` to not include the final pooling and fully connected layer in the original model. Different Pooling layers can be added based on the requirement and a dense output layer to the ResNet-50 model. Number of Classes parameter depends on the dataset.

4.6 Training the model:

The Python deep learning modules `keras.fit()` and `keras.fit_generator` can be used to train our machine learning and deep learning models. These two functions can achieve the same result. The primary difference is that `fit` is employed when the entire training dataset fits into memory without the need for data augmentation. The `fit_generator` is used when fitting a huge dataset into memory or when data augmentation is required. Before training, the model must be built. After training the model, we can use the `keras.fit` method to determine its validity.

Keras.Fit method: We developed iterators for picture enhancement. We input it to the neural network, which augments it. All we need to do is pass the iterator, together with epochs, size of the batch, and remaining essential parameters, into the `fit` method used to the network prototypical. A Convolutional Neural Network (CNN) prototypical will be used. The `fit` method applies the prototypical to figures generated in batches by a Python generator. The following are `fit` method arguments.

- The first argument is used for iterating for the training of different images that we get from the `flow()`.
- Epoch makes the count of forward and backward authorizations of the working out information.
- Number of epochs is an imperative parameter postulates the quantity of batches of images that are in a single epoch. It is habitually engaged as the measurement of the original set of data alienated by the length of batch.
- Validation data took the corroboration set of data or the corroboration generator outcome from the generator technique.
- Steps of Validation is comparable to number of steps in epoch, but for corroboration data. This can be cast off when we are expanding the justification set of pictures as well.
- Batchsize : it can gross any integer value or NULL and by default, it will be set to 32 generally. It postulates no. of models per gradient.

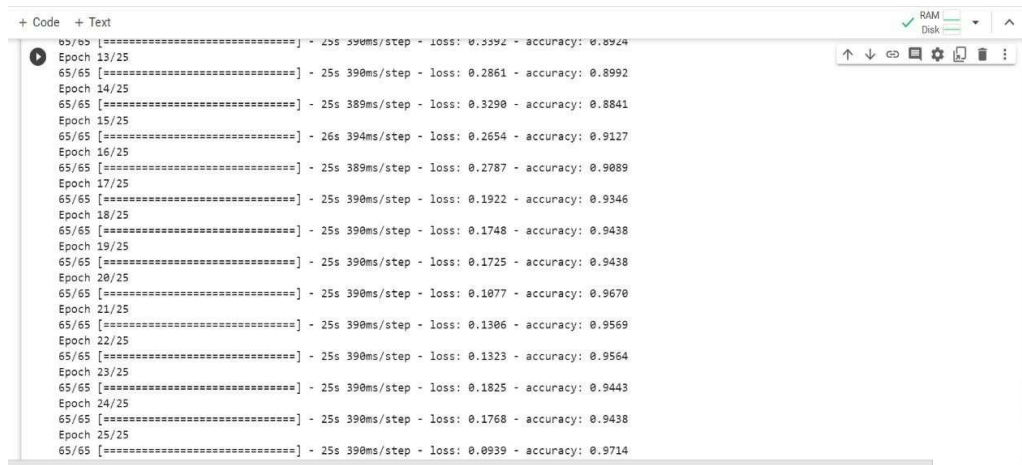


Fig.4.8 Training accuracy

```

[ ] print("Evaluate on test data")
    results = model.evaluate(x_test, y_test, batch_size=32)
    print("test loss, test acc:", results)

Evaluate on test data
22/22 [=====] - 5s 148ms/step - loss: 0.9971 - accuracy: 0.7602
test loss, test acc: [0.9971484541893005, 0.7601743936538696]

```

Fig.4.9 Validation accuracy

4.7 Results:

The dataset is given as an input to perform Pre-processing and Data Augmentation. After that features are extracted using Resnet50 by pooling and model is trained for the given dataset(70%trained) and for the model, video is given as input and gives output as healthy or diseased leafas shown in Fig.4.10, Fig.4.11, Fig.4.12, Fig.4.13, Fig.4.14 and Fig.4.15.

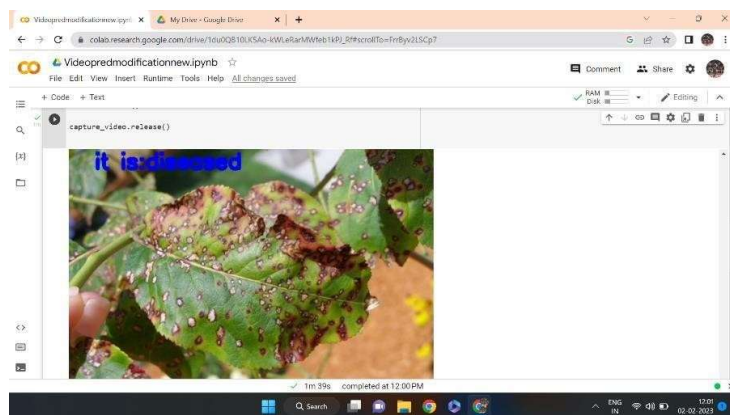


Fig.4.10 Diseased leaf

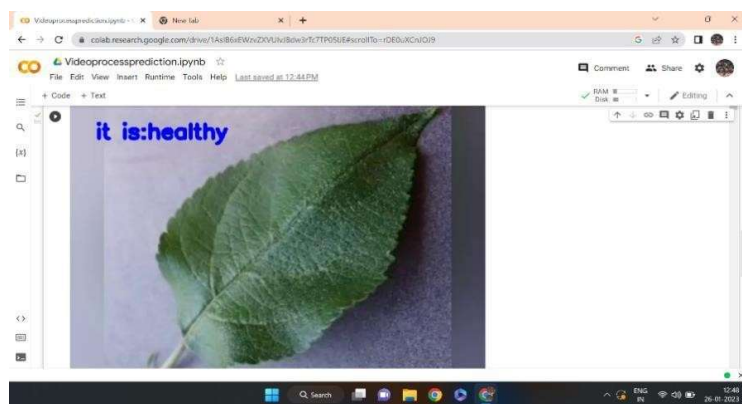


Fig.4.11 Healthy leaf

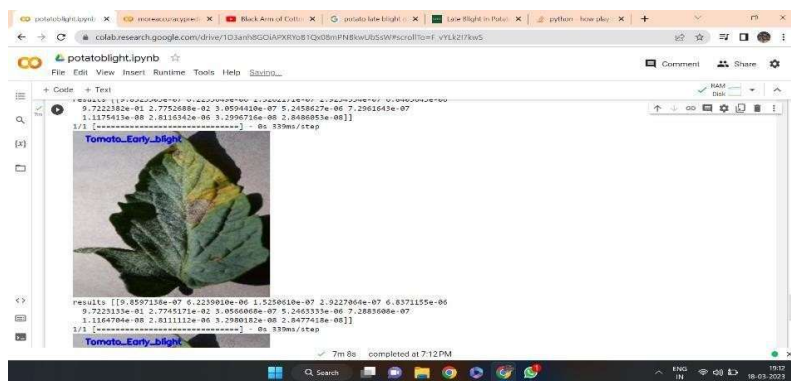


Fig.4.12 Diseased with tomato early blight

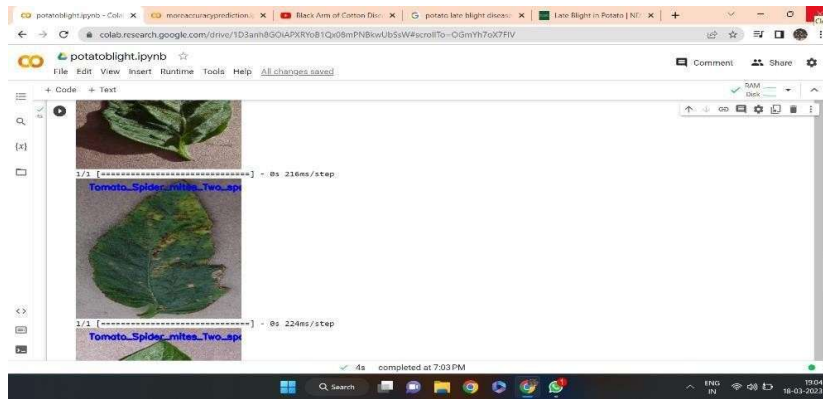


Fig.4.13 Diseased with Tomato spider mites

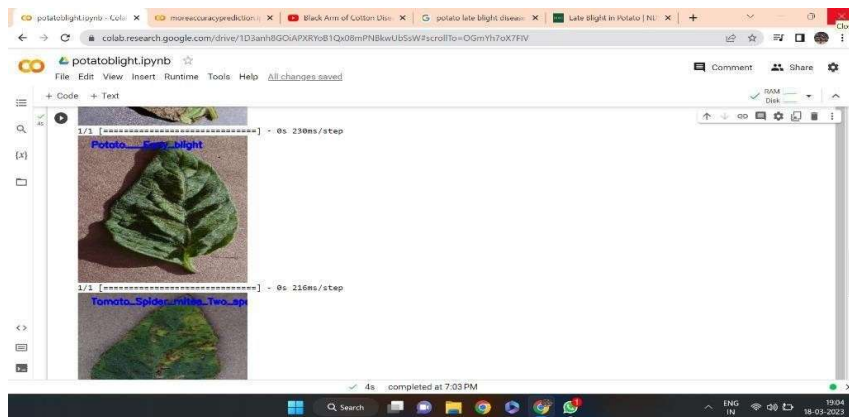


Fig.4.14 Diseased with Potato early blight

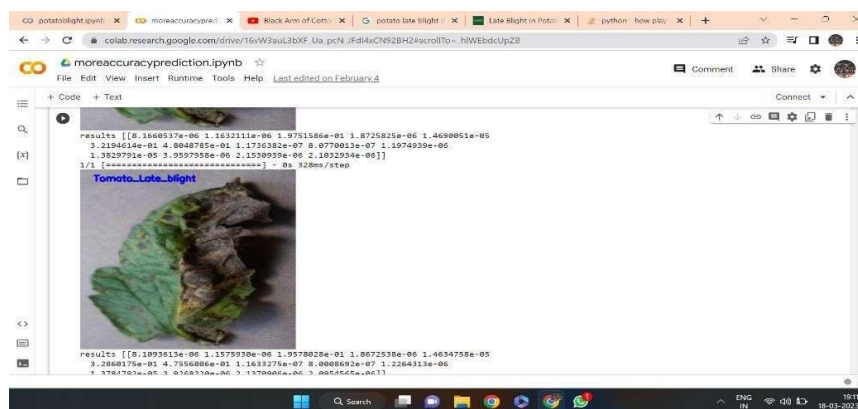


Fig.4.15 Diseased with Tomato late blight

CHAPTER 5

CONCLUSION AND FUTURE WORK

Plant leaf disease detection using machine learning algorithms has shown promising results in recent years. Machine learning algorithms can accurately classify plant leaves into healthy and diseased categories, thus helping farmers to identify and treat plant diseases early on, and preventing crop damage and yield loss.

In this project, we have explored various machine learning algorithms such as Logistic Regression, KNN, Naive Bayes, Support Vector Machines (SVMs), and Random Forests for plant leaf disease detection. We have also examined the performance of these algorithms on different datasets and evaluated their accuracy, precision, recall, and F1 score.

We also explored the use of ResNet-50, a deep convolutional neural network, for plant leaf disease detection using machine learning algorithms for video datasets. ResNet-50 is a powerful algorithm that can extract high-level features from images, making it well-suited for plant leaf disease detection. We evaluated the performance of the algorithm on various datasets and achieved high accuracy in detecting plant diseases.

Overall, the results of our experiments show that Random Forests Classifier outperform other algorithms in terms of accuracy and generalization on image datasets and ResNet-50 for video datasets. Therefore, these algorithms can be used as a reliable and efficient algorithm for plant leaf disease detection.

The use of automated monitoring and management systems are gaining increasing demand with technological advancement. In the agricultural field, loss of yield mainly occurs due to widespread disease. Mostly the detection and identification of the disease is noticed when the disease advances to a severe stage, therefore, causing the loss in terms of yield, time and money. The proposed system is capable of detecting the disease at an earlier stage as soon as it occurs on the leaf, hence saving the loss and reducing the dependency on the expert to a certain extent is possible. It can provide the help for a person having less knowledge about the disease. Depending on these goals, we have to extract the features corresponding to the diseases.

REFERENCES

1. Haseeb Nazki, Sook Yoon, Alvaro Fuentes, Dong Sun Park “Unsupervised image translation using adversarial networks for improved plant disease recognition” Published by Elsevier B.V,(2020).
2. Shanwen Zhang, Subing Zhang, Chuanlei Zhang, Xianfeng Wang, Yun Shi “Cucumber leaf disease identification with global pooling dilated convolutional neural network” Published by Elsevier B.V, (2019).
3. Uday Pratap Singh, Siddharth Singh Chouhan, Sukirty Jain, And Sanjeev Jain “Multilayer Convolution Neural Network for the Classification of Mango Leaves Infected by Anthracnose Disease” (2019).
4. Vijai Singh “Sunflower leaf diseases detection using image segmentation based on particle swarm optimization” 2019 Published by Elsevier,(2019).
5. Sumita Mishra, Rishabh Sachan, Diksha Rajpal “Deep Convolutional Neural Network based Detection System for Real-time Corn Plant Disease Recognition” 2020 Published by Elsevier B.V, (2019).
6. Parul Sharma, Yash Paul Singh Berwal , Wiqas Ghai “Performance analysis of deep learning CNN models for disease detection in plants using image segmentation” open access 2019 Published by Elsevier B.V, (2019).
7. Mohit Agarwal, Abhishek Singh, Siddhartha Arjaria, Amit Sinha, Suneet Gupta “Tamato Leaf Disease Detection using Convolution Neural Network” 2019-2020 Published by Elsevier,(2019).
8. Aditya Khamparia, Gurinder Saini, Deepak Gupta, Ashish Khanna, Shrasti Tiwari, Victor Hugo C. de Albuquerque “Seasonal Crops Disease Prediction and Classification Using Deep Convolutional Encoder Network” ,(2019).
9. Srdjan Sladojevic, Marko Arsenovic, Andras Anderla, Dubravko Culibrk and Darko Stefanovic “Deep Neural Networks Based Recognition of Plant Diseases by Leaf Image Classification”, Volume 2016 Hindawi Publishing Corporation,(2016).
10. Majji V Applalanaidu, G. Kumaravelan “A Review of Machine Learning Approaches in Plant Leaf Disease Detection and Classification” IEEE,(2021).
11. Qiong Ren , Hui Cheng and Hai Han “Research on machine learning framework based on random forest algorithm” : AIP Conference Proceedings,(2017).
12. Tao Xiang, Tao Li, Mao Ye, and Zijian Liu “Random Forest with Adaptive Local Template for Pedestrian Detection” Hindawi Publishing Corporation,(2015).

13. Md Nasim Adnan “Improving the Random Forest Algorithm by Randomly Varying the Size of the Bootstrap Samples”Adnan,(2014).
14. Ziming Wu, Weiwei Lin, Zilong Zhang and Angzhan Wen “An Ensemble Random Forest Algorithm for Insurance Big Data Analysis”IEEE,(2017).
15. Manjunath Badiger, Varuna kumara,Sachin CN shetty,Sudhir poojary “Leaf and skin disease detection using image processing” Global Transactions Proceedins,(2022).
16. Niveditha M, Pooja R, Prasad Bhat N, shashank N, “Plant disease detection using machine learning” IEEE (2021).
17. Nishant Shelar ,Suraj shinde ,Shubham sawant ,Shreyas dhumal “Plant disease detection using CNN ” Turkish Journal of Computer and Mathematics Education, (2021).
18. Madhuri Devi Chodey, Dr.Noorilla Shariff C, Gauravi Shetty “Pest detection in crop using video and Image processing” IJRASET (2020).
19. Aryan Garg“Image Classification Using Resnet-50 Deep Learning Model”Analytics vidya,(2022).
20. Devvi Sarwinda , Radifa Hilya Paradisa , Alhadi Bustamam ,Pinkie Anggia “Deep Learning in Image Classification using Residual Network (ResNet) Variants for Detection of Colorectal Cancer” International Conference on Computer Science and Computational Intelligence,(2020).