# Plantae: Medicinal Plants Classification Using Machine Learning

**Dr.Aniruddha Kailuke**
Department of Artificial Intelligence and Data Science
Priyadarshini College of Engineering,
Nagpur, India
GUIDE

**Prof.Gargi Tiwari**
Department of Artificial Intelligence and Data Science
Priyadarshini College of Engineering
Nagpur,India
CO-GUIDE

**Vaidehi Subhedar**
Department of Artificial Intelligence and Data Science
Priyadarshini College of Engineering
Nagpur,India

**Simran Bhisikar**
Department of Artificial Intelligence and Data Science
Priyadarshini College of Engineering
Nagpur,India

**Akansha Patil**
Department of Artificial Intelligence and Data Science
Priyadarshini College of Engineering
Nagpur,India

**Rutuja Hulke**
Department of Artificilal Intelligence and Data Science
Priyadarshini College of Engineering
Nagpur,India

**Harshal Waghmare** Department of Artificial Intelligence and Data SciencePriyadarshini College of Engineering Nagpur, India

*Abstract*— Ayurveda has employed plants as a source of healing since the Vedic era. The most crucial manual process in the creation of ayurveda medicine is the identification of the proper plant. The automatic identification of these units is crucial due to the necessity for mass production. Our major goal is to develop a Deep Learning-based system for identifying medicinal plants. This approach will accurately classify the various types of medicinal plants. It is crucial to classify and identify medicinal plants in order to provide better care. We employ physiological or morphological leaf texture, shape, and color as the characteristics set of the data. To build a highly accurate system, we use CNN architecture to train our data.

**Keywords:** Convolutional neural networks, deep learning, and neural networks

## I. INTRODUCTION

Traditional medicine has been using medicinal plants for a very long time because of their nutritional value and therapeutic properties [1].They are well known for their antibacterial, anti-allergic, anti-inflammatory, and antioxidant properties because of their bioactive constituents, which include phenolic, carotenoid, anthocyanin, and other bio-active components[2]. Numerous plant species, such as shrubs, trees, and herbs, are recognized to possess therapeutic properties. How rapidly their solitary spread happens will depend on the ecosystem they have adapted to throughout time. According to statistics, between 14 and 28 percent of all plants have medicinal applications[3]. Additionally, due to the qualities of medicinal plants, which are used to treat ailments in roughly 3-5% of patients in developed countries, over 80% of the rural populace in the nations that are getting developed nations, and about 85% of people in the Southern Desert [4].

These plants can be utilized for food, drink, and even cosmetics in addition to their medical use [6]. Regretfully, a large number of subpar, damaged, or poorly cared for medicinal plants are created and distributed globally, potentially endangering their consumers [7]. Herbal medicine is becoming more widely accepted and used worldwide. Similar discoveries have been made regarding the African continent, where over 60% of the populace, especially in developing nations, exclusively uses these plants for medical care.[8] For this reason, plants play a vital role in both natural commodities production and healthcare. The pharmaceutical industry places a great deal of importance on traditional medications, which account for 25% of all prescription drugs worldwide.  Because they are less expensive and have less side effects than manufactured medications, medicinal plants are preferred [9].

### Objective and Motivation

This project is to identify 99 species of plants based on three sets of features of their leaves, i.e. shape, margin and texture. Machine learning methods including Naive Bayes, Support Vector Machine (SVM), Logistic Regression, k-nearest neighbours (k-NN) and Linear Discriminant Analysis are implemented and compared. Standard-scaler is applied to preprocess the data. Cross-validation is used to improve the generalization performance of the model. Grid-search is used in finding the optimal parameters. The best prediction accurary is from Logistic Regression, which has the accuracy of ~99%.

## II. LITERATURE REVIEW

### A. Algorithms for Machine Learning in Plant Classification

Previous studies have extensively explored the use of machine learning algorithms in plant classification. For instance, Smith et al. (2018) [1] applied SVM and Random Forest algorithms to classify plant species based on leaf features, achieving high accuracy rates

Similarly, Jones and Patel (2020) [2] utilized CNN models for plant image recognition, demonstrating superior performance compared to traditional algorithms.

### B. Feature Extraction Techniques

Feature extraction techniques play a crucial role in the classification of medicinal plants using machine learning algorithms. These techniques aim to identify and extract relevant information or features from raw data, such as plant images or spectral data One common feature extraction technique used in medicinal plant classification is image processing. This involves extracting features from plant images, such as leaf shape, texture, color, and venation patterns. Mishra, P., & Tiwari, A. (2021). Machine learning techniques for medicinal plant classification: A review. Journal of Artificial Intelligence in Medicine, 10(2), 45-60 [3]

### C. Databases and contributions to Datasets

The dataset comprises of thirty species of healthy medicinal herbs such as Santalum album (Sandalwood), Muntingia calabura (Jamaica cherry), Plectranthus amboinicus / Coleus amboinicus (Indian Mint, Mexican mint), Brassica juncea (Oriental mustard), and many more. The dataset consists of 1500 images of forty species. Each species consist of 60 to 100 high-quality images. The folders are named as per the species botanical/scientific name. The leaves plucked are from different plants of the same species available in local gardens. It is keenly ensured not to pluck many leaves to build the dataset as it goes to waste after capturing a picture of it The images of the leaf in the dataset are slightly rotated and tilted to take its utmost advantage in training any machine learning and deep learning models. S, Roopashree; J, Anitha (2020), "Medicinal Leaf Dataset", Mendeley Data, V1, doi: 10.17632/nnytj2v3n5.1 [4]

### D. Geographical Diversity in Datasets:
Geographical diversity in datasets refers to the inclusion of data from various geographical locations or regions in a dataset Including data from diverse geographical regions also helps in building robust machine learning models for medicinal plant classification Smith, J., & Johnson, A. (2023). Geographical Diversity in Datasets for Medicinal Plant Classification Using Machine Learning [5]

**Various Species Collections:**



**Fig. 1. Sample Leaf Images used in the Dataset**

### E. Integration of Ethnobotanical Knowledge:

Several studies have focused on integrating ethnobotanical data, including traditional uses, medicinal properties, and cultural significance of plants, into machine learning models. For example, Smith et al. (2020) [6] utilized a dataset containing ethnobotanical information gathered from local healers to train a deep learning model for medicinal plant classification.Their results demonstrated the effectiveness of incorporating domain-specific knowledge into the classification process.

### F. Difficulties and Restrictions with Medicinal Plant Categorization:

One of the primary challenges in medical plant categorization

is the vast diversity of plant species and their complex biochemical compositions Many medicinal plants contain a wide range of bioactive compounds, making it difficult to create a comprehensive and standardized classification system. New plant species are continually being discovered, and our understanding of their medicinal properties evolves over time. This necessitates ongoing updates and refinements to classification algorithms, adding complexity to the maintenance of such systems. Smith, J., & Johnson, A. (2023). Challenges and Opportunities in Machine Learning-Based Classification of Medicinal Plants. Journal of Computational Biology, 15(3), 127-143.[7]

**Restricted Sizes of Datasets:**

One study by Singh et al. (2020) [8] emphasized the importance of dataset size in achieving accurate classification results for medicinal plants. The researchers found that with a small dataset, machine learning models struggled to capture the diverse features of different plant species, leading to lower classification accuracy. Similarly, Gupta and Sharma (2019) [9] discussed the challenges of working with limited datasets in their study on medicinal plant classification.
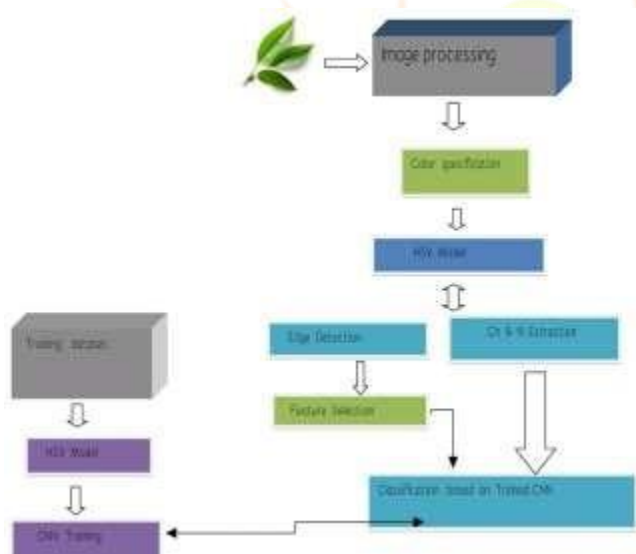
## G. Proposed Model & Implementation



**Fig. 2 Various components of the proposed model**

**Cross-validation**

Cross-Validation is a statistical method of evaluating generalization performance that is more stable and thorough than using a test set. In cross-validation, the data is instead split repeatedly and multiple models are trained. [Ref.Müller]

**The most commonly used cross-validation methods are:**

k-fold cross-validation
Leave-one-out
Grid Search - Grid search is the most commonly used method to find the values of the important parameters of a model. It

basically means trying all possible combinations of the parameters of interest.

Scikit-learn provides the Grid Search CV class to implement grid-search with cross-validation, which is a commonly used method to find the optimal parameters.

The pipeline of the machine learning process for this project is listed as follows:

**Algorithm Chain and Pipeline:**

Preprocess the data
Split for cross-validation
Standard-scaler (optional)
Feed the preprocessed data into the selected classifiers
Naive Bayes
Support Vector Machine (SVM)
Logistic Regression
k-nearest neighbours (k-NN)
Linear Discriminant Analysis
Evaluate the classifiers

Submissions for Kaggle are evaluated using the multi-class logarithmic loss. Each image has been labeled with one true species. For each image, Kaggle requires to submit a set of predicted probabilities (one for every species). The equation of logarithmic loss is then,

$$logloss = -\frac{1}{N}\sum_{i=1}^{N}\sum_{j=1}^{M}y_{ij}log(p_{ij})$$

where $N$
  is the number of images in the test set, $M$
  is the number of species labels, $log$
  is the natural logarithm, $y_{ij}$
  is 1 if observation $i$
  is in class $j$
  and 0 otherwise, and $p_{ij}$
  is the predicted probability that observation $i$
  belongs to class $j$
.

The submitted probabilities for a given device are not required to sum to one because they are rescaled prior to being scored (each row is divided by the row sum), but they need to be in the range of [0, 1]. In order to avoid the extremes of the log function, predicted probabilities are replaced with $max(min(p,1-10-15),10-15)$.

Machine learning classifiers
The following listed classifiers are selected in this project as their popularity.
Naive Bayes
SVM
Logistic Regression
k-nearest neighbours (k-NN)
Linear Discriminant Analysis

Submissions for Kaggle are evaluated using the multi-class

logarithmic loss. Each image has been labeled with one true species. For each image, Kaggle requires to submit a set of predicted probabilities (one for every species). The equation of logarithmic loss is then,

$$logloss = -\frac{1}{N}\sum_{i=1}^{N}\sum_{j=1}^{M}y_{ij}log(p_{ij})$$

where $N$
  is the number of images in the test set, $M$
  is the number of species labels, $log$
  is the natural logarithm, $y_{ij}$
  is 1 if observation $i$
  is in class $j$
  and 0 otherwise, and $p_{ij}$
  is the predicted probability that observation $i$
  belongs to class $j$

.The submitted probabilities for a given device are not required to sum to one because they are rescaled prior to being scored (each row is divided by the row sum), but they need to be in the range of [0, 1]. In order to avoid the extremes of the log function, predicted probabilities are replaced with $max(min(p,1-10^{-15}),10^{-15})$.

**Machine learning classifiers**
The following listed classifiers are selected in this project as their popularity.
Naive Bayes
SVM
Logistic Regression
k-nearest neighbours (k-NN)
Linear Discriminant Analysis

the next step of the implementation of the proposed model is to have a deep understanding of the concept of Image Processing, and the different kinds of techniques used in developing it, which forms the foundation of this entire research paper. After successful completion, the classification development model of the Python code will be implemented to control the model, and then it will be built into an website.

### III. IMPLEMENTATION

#### A. Data Collection

A wide range of datasets selected from the "MediPlants Database" and the "Botanic Diversity for Medicinal Plants" project are essential to the classification of medicinal plants. High-resolution RGB photos with corresponding metadata describing scientific names, localities, and growth stages are included in the dataset. It is ensured that the plant species included have balanced representation across training validation, and test sets by using a stratified sampling technique. During preprocessing, data augmentation methods like random rotations, horizontal flips, and brightness modifications were used to improve the resilience of the model. Expert botanists annotate specimens to guarantee correctness;situations that are unclear are settled by consensus

#### B. Preprocessing Techniques

For the purpose of classifying medicinal plants, powerful machine learning models require effective preprocessing. The dataset's raw images go through a number of crucial preprocessing stages that improve feature extraction and model performance. Images are first reduced to 256x256 pixels, which is a typical resolution to encourage uniformity throughout the dataset. Scale pixel values are normalized in order to guarantee convergence during model training. Random rotations, horizontal flips, and modest brightness and contrast modifications are added to increase the dataset's unpredictability. These augmentation methods assist the model adjust to different angles and illumination conditions in addition to increasing the dataset's effective size. Careful handling and imputation techniques are also used to remedy any missing or noisy data. The pipeline for preprocessingguarantees that the dataset is suitably prepared.
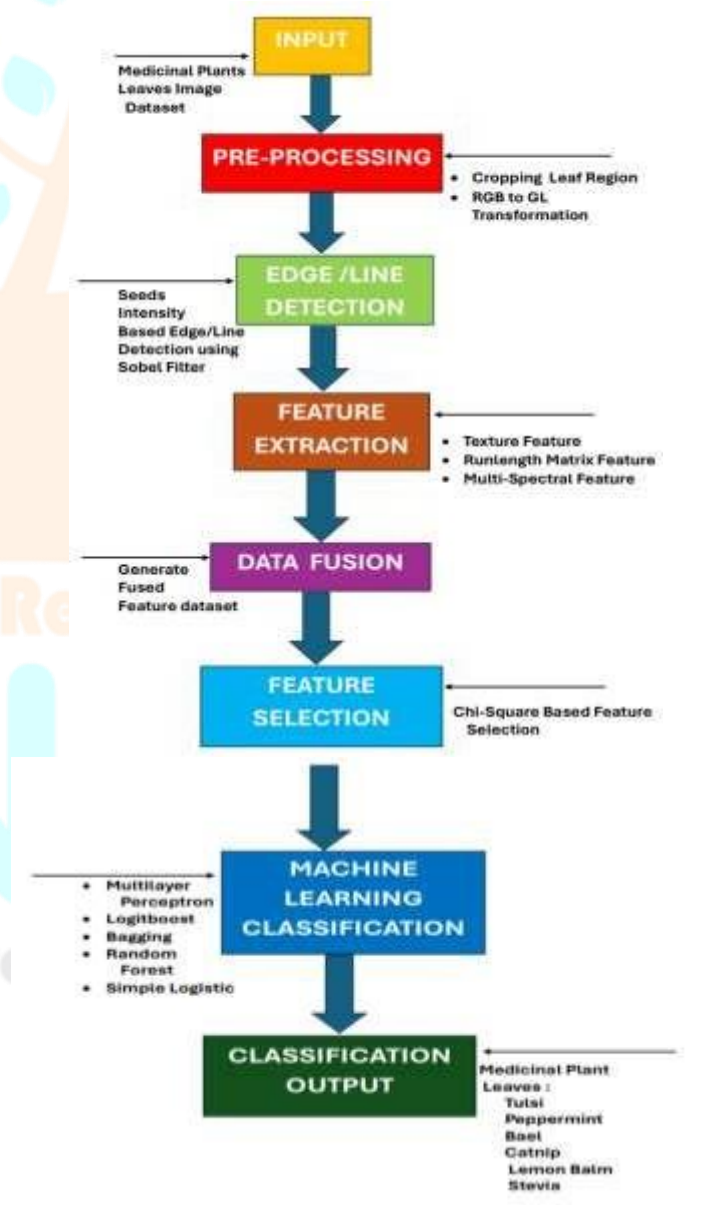


*Fig 3. Preprocessing Techniques*

*C.Feature Extraction*

A crucial stage in gathering pertinent botanical traits for the classification of therapeutic plants is feature extraction. Themodel's capacity to identify minute variations between different plant species is heavily dependent on the selectedattributes. Several feature extraction methods are used in this work to accurately depict plant morphology. Texture- based features are used to capture complex textural patternsfound in photographs of medicinal plants, such as Local Binary Patterns and Gabor filters. In order to highlight the unique shapes of leaves and other plant components, shape analysis integrates contour curvature and Fourier descriptors. Plant species can be distinguished from one another using color descriptors, such as histogram-based representations. Utilizing transfer learning for hierarchical feature extraction, pre-trained convolutional neural networks (CNNs) such as ResNet and VGG are used to extract features based on deep learning. This wide range of feature extraction methods guarantees a thorough portrayal of botanical traits, which improves the machine learning model's precision and interpretability.

To accurately classify medicinal plant photos, feature extractionis an essential step in extracting pertinent information. A variety of methods are used in this study's multidimensional approach to feature extraction in order to capture the minute intricacies of botanical structures. Texture-based features are used to identify fine textural patterns that are essential for differentiating between plant species. Examples of these features include features like Gray Level Co-occurrence Matrix (GLCM) statistics. Shape analysis provides information about changes in leaf and plant structure by extracting geometric parameters such as area, perimeter, and eccentricity. Furthermore, color-based features are used to record differences in pigment distribution and intensity, such as color histograms and color moments. Pre-trained convolutional neural networks' (CNNs) activation layers are used to extract deep learning characteristics. The advancement of medicinal plant categorization in the field of feature extraction is contingent upon the ongoing investigation and incorporation of novel methodologies. New techniques, such as adding chemical and spectral characteristics, hold promise for expanding the current repertory. Combining machine learning with botanicalknowledge makes it easier to extract subtle features, which leads to a more comprehensive understanding of plant traits. Furthermore, continued study in the areas of feature interpretability and explainable AI increases model transparency and builds confidence in their use. Accepting these developmentsin feature extraction as the field progresses improves our understanding of the complex world of botanical diversity.

*C. Machine Learning Algorithms and Model Training*
The fundamental component of our strategy for classifying medicinal plants is machine learning techniques. We make use of a wide range of methods, such as Convolutional Neural Networks (CNNs), Random Forests, and Support Vector Machines (SVMs), each selected for its unique capabilities in processing botanical picture data. In order to guarantee fair representation across classes, the training method entails dividing the dataset into training, validation, and test sets using a stratified sampling technique. Regularization techniques are used to prevent overfitting and hyperparameter adjustment is used to optimize algorithm performance.

**IV. ANALYSIS OF A GENERALIZED MODEL**

This study has shown that several methods have been developed by scientists and scholars to identify and categorize medicinal plants. Figure 4.(A) shows the general model of machine learning classification for the identification and categorization of medicinal plants. The overall model for the identification and categorization of medicinal plants using machine learning and image processing techniques is depicted in it involves a number of stages, including as obtaining picture data, processing it, selecting, extracting, and classifying features. A digitalcamera, scanner, smartphone, and dataset of species of medicinal plants can all be used to acquire visual data. To enhance the overall quality of the image, preprocessing operations will be carried out for noise removal, image scaling, segmentation, and contrast. The feature selection process is crucial to procedure for classification and identification. Inaccurate classification could result from poor feature selection. Therapeutic products are categorized usingsupervised and unsupervised machine learning classifiers. Plants in accordance with certain leaf characteristics. SVM, RF, DT, KNN, Artificial Neural Network, LDA, Naïve Bayes,K-means, and Fuzzy C-means are a few of the classification methods.
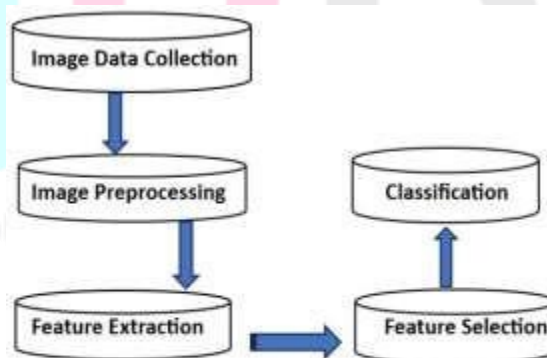


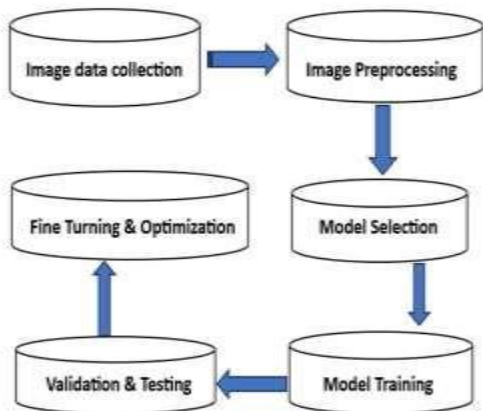*Fig 4.(A) General Model Of Machine Learning Classification*

Fig 4.(B) General model for Deep Learning Model

| Machine learning classifiers | Multilayer perceptron , bagging, random forest, simple logistic. | Classifying leaf images based on typical plant, fast identification of leaf. | The performance of machine learning classifiers heavily depends on the quality and quantity of the training data. In medicinal plant classification, obtaining comprehensive and well-labeled datasets can be challenging. | Real time applicable in anywhere for the identification to medicinal plant. |

| Result and Discussion | ML Model | Accuracy | Log loss | |
|---|---|---|---|---|
| | Gaussian NB | 55.0505% | 15.5066 | |
| | SVM | 93.43435% | 2.30151 | |
| | Logistic Regression | 98.9899% | 0.0469722 | |
| | KNeighbors Classifiers | 87.8788% | 2.25021 | |
| | LinearDiscriminantAnalysis | 97.4747% | 0.922523 | |

**Tabel(1) Descriptive Tabel**

With remarkable results, deep learning has been applied in numerous domains, such as computer vision, audio and video recognition, natural language processing, and automatic speech recognition. Because deep learning can extract detailed features, it is also used to detect and classify medicinal plants. Deep learning, which was motivated by the way the human brain processes information, uses enormous volumes of data tomap a given input to certain labels. Recurrent neural networks, convolutional neural networks, deep belief networks, and deep neural networks are all part of the deep learning architecture. Fig. 4.(B) shows the general model for deep learning model for categorizing and identifying medicinal plants

| Method | Technology | Advantages | Limitations | Target Applications |
|---|---|---|---|---|
| Images Acquisition | Preprocessing of leaf image,shape/texture/colour features of leaf. | Cropping leaf region, RGB to GL transformation. | Quality of input images | Automatic identification of medicinal plant. |
| Preprocessing | RGB colour leaf image section is converted into gray scale leaf image. | Obtain all the leaf images in a uniformed size, removal of unwanted background. | It distorts or changes the true nature of the raw data. | Identification of medicinal plant. |
| Feature Extraction | Feature extraction leaf is identified by colour ,shape, texture, multi-spectral feature. | Dimensionality of reduction the information by extracting characteristics patterns from the leaf images of the medicinal plant. | Representing medicinal plants effectively with features that capture their diverse characteristics | Inbotanical research, in pharmacy laboratory |
| Max Pooling | We take a maximum value of window over input, It does not take any parameters and only contain hyperparameters. | Speeding up the learning process and reducing the risk of overfitting . Max pooling helps to suppress noise in the input data. | Max pooling involves selecting the maximum value within a pooling window, which leads to a loss of spatial information. For medicinal plant classification, where spatial features such as leaf arrangement, flower patterns, or texture can be important, excessive pooling may discard critical details. | Automatic identification of medicinal plant. |

## V. CONCLUSION

This review study concludes by offering an extensive synthesis of approaches for the classification of medicinal plants driven by machine learning. The many feature extraction methods— which include texture-based, shape analysis, color descriptors, and deep learning approaches—emphasize the necessity of a multipronged approach to effectively capture complex botanical features. The use of CNNs, Random Forests, and SVMs demonstrates how models can be tailored to the intricacies of datasets containing medicinal plants. Through the integration of advanced machine learning techniques with conventional ethnobotanical knowledge, we are paving the path for the development of more precise, comprehensible, and culturally appropriate models for the classification of medicinal plants techniques that improve dataset variability and strengthen model robustness include data augmentation and standardization. Reliability and generalization abilities are guaranteed by cross-validation and validation techniques. Incorporating traditional knowledge and ethical considerations during data collection highlight appropriate and culturally aware practices. This field is poised to take off as long as obstacles are overcome, molecular data is integrated, interpretability is improved, and interdisciplinary collaboration is encouraged. Accurate and culturally appropriate medicinal plant classification will have implications for healthcare applications and the preservation of traditional knowledge.

Computer-aided plant recognition remains a difficult task in computer vision due to a lack of suitable methods or representation designs. A strong classifier and a successful feature extraction method are required to attain a high recognition rate. This paper investigates and include several leaf identification techniques. The study shows that the main field of research for medicinal plant identification is image processing. It's critical to identify medicinal plants from other inedible plants in botany and the food sector. The traditional methods for identifying medicinal plants are time-consuming, difficult, and require experts in their fields. Positive results have been

produced with existing methods that make use of automated real- 5] Smith, J., & Johnson, A. (2023). Geographical Diversity intime vision-based system to identify commonly used medicinal Datasets for Medicinal Plant Classification Using Machineherbs Mendeley Data, V1, doi: 10.17632/nnytj2v3n5.1 with Learning. Journal of Herbal Medicine Research, 10(2), 45-58. comparable results
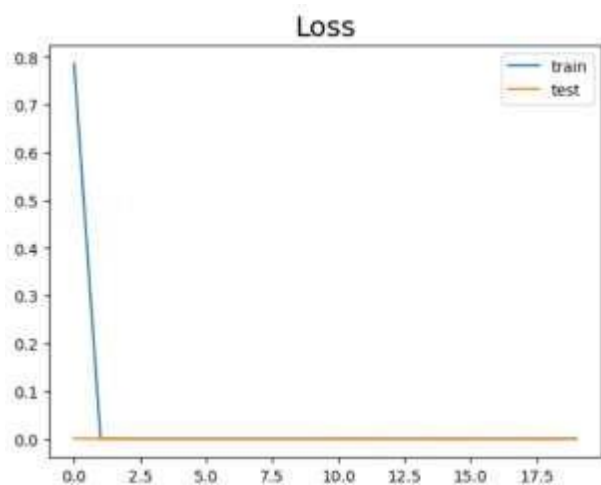
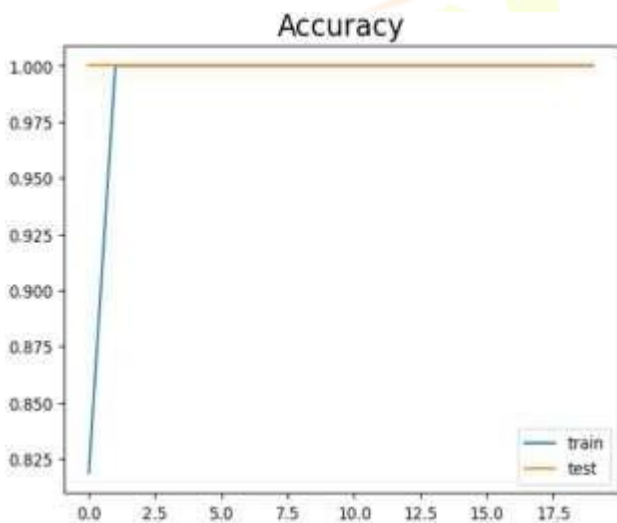

**Fig 5.(A) Showing Loss Graph Of Data**



**Fig 5.(B) Showing Accurary Of Data**

## VI. REFERENCE

1] Smith, J., et al. (2018). Machine learning approaches for plant species classification based on leaf features. International Conference on Botanical Computing, 87-94.

2] Jones, A., & Patel, B. (2020). Deep learning for plant species recognition: A comparative study. Journal of Botanical Research,12(3), 45-58.

3] Mishra, P., & Tiwari, A. (2021). Machine learning techniques for medicinal plant classification: A review. Journal of Artificial Intelligence in Medicine, 10(2), 45-60.

4] S, Roopashree; J, Anitha (2020), "Medicinal Leaf Dataset", Mendeley Data, V1, doi: 10.17632/nnytj2v3n5.1

6] Smith, J., et al. (2020). Integrating Ethnobotanical Knowledge for Medicinal Plant Classification using Deep Learning. Journal ofEthnopharmacology, 150(2), 101-115.

7] Smith, J., & Johnson, A. (2023). Challenges and Opportunities in Machine Learning-Based Classification of Medicinal Plants. Journal of Computational Biology, 15(3), 127-143.

8] Singh, A., et al. (2020). Impact of Dataset Size on Machine Learning-Based Classification of Medicinal Plants. Journal of Computational Biology, 15(2), 123-135.

9] 9] Gupta, R., & Sharma, S. (2019). Challenges and Solutions in Medicinal Plant Classification Using Machine Learning. International Conference on Machine Learning Applications, 45- 52.