



Image-Based Papaya Health Classification Using Machine Learning

Harsha Nishad¹, Aditi Sharma², Ganesh Kumar P³, Monal Prajapati⁴, MD Hamza⁵

^{1,2,3,4,5}Department of Computer Science, Bhilai Institute of Technology Raipur, Chhattisgarh, India

ABSTRACT

Papaya is a widespread fruit, but its production is heavily affected by a variety of diseases. Early detection of diseased papayas can help farmers take precautions to minimize harvest losses. This project presents a Papaya Fruit Disease Recognition System with machine learning technology (machine learning) to classify papayas as healthy or disease.

Models were meticulously trained on diverse datasets to ensure robust performance of the disease. To evaluate the model output, K-fold Cross validation was used to test each segmentation method for each classification. Under all combinations, segmentation of fuzzy C-means combined with random forest achieved the highest accuracy. Secondly, Streamlit was provided to this optimized model, so users uploaded Papaya images and received the disease of instant disease classification.

Key Words: Papaya Disease Classification, Machine learning, Deep learning, Papaya, Computer vision.

1. INTRODUCTION

Papaya (*Carica papaya* L.) is a tropical fruit that spreads to nutritional and economic benefits. However, its production is often threatened by a variety of diseases that affect the quality and harvest of the fruit. Traditional identification methods for disease detection are based on time-consuming, labor-intensive, and manual testing that is prone to human failure. To address this challenge, machine learning (ML) techniques provide an automated, accurate and efficient approach to identifying diseased fruits.

This project focuses on detection of papaya fruit disease using image segmentation and classifiers for machine learning for binary classification (health and disease). The aim is to develop a system that can analyze papaya images, extract essential properties, and accurately classify them.

In extensive experiments, the combination of fuzzy C-means segmentation and random forest achieved the highest accuracy. This optimized model was provided using Streamlit so that users upload images and receive immediate result of classification.

Through this research, we aim to contribute to the development of efficient and technology-driven solutions for enhancing crop productivity and sustainability in papaya cultivation.

2. LITERATURE SURVEY

Traditional methods for identifying diseases require manual inspection, errors, and expertise[1]. To address these challenges, researchers investigated various image processing and machine learning (ML) technologies for automated disease classification[3].

| S. No. | Paper | Year | Dataset | Algorithm/Method | Result and Accuracy |
|--------|---|------|---|---|--|
| 1 | Papaya Disease Classification Using Machine Learning (IJARCCE) | 2024 | Begins with the acquisition of a diverse dataset comprising around 3000 images of papaya. | The YOLOv9c model, renowned for its object detection capabilities, is selected as the primary model for disease classification | 89% accurate in identifying and managing papaya diseases. |
| 2 | YOLO-Papaya (MDPI Electronics Journa) | 2023 | This dataset comprises 23,158 images, categorized into nine distinct classes, including various papaya fruit diseases | Developed a deep neural network using YOLO and CBAM for disease classification. | Achieved an accuracy of 92.4% in detecting papaya fruit diseases. |
| 3 | Neural Network for Papaya Leaf Disease Detection (Acta Graphica Journal) | 2021 | The use of papaya leaf images that underwent preprocessing steps like compression and filtering to enhance quality. | By using the internet of things and blockchain technologies traceability system for real-time food tracing, build a safety control system for food supply chain | The system showed high reliability and accuracy in maintaining food safety across different stages of the food supply chain. |

| | | | | | |
|---|---|------|---|--|--|
| 4 | "BDPapayaLeaf: A dataset of papaya leaf for disease detection, classification, and analysis" (Elsevier) | 2024 | 2,159 images of papaya leaves across five classes (healthy and four diseases) | Dataset creation for disease detection and classification | Provides a comprehensive dataset to facilitate development of accurate disease detection models. |
| 5 | "Maturity status classification of papaya fruits based on machine learning and transfer learning approach" (Elsevier) | 2021 | 300 papaya fruit images across three maturity stages | Machine Learning (HOG features with KNN) and Transfer Learning (VGG19) | Both approaches achieved 100% accuracy in classifying papaya maturity stages. |
| 6 | "Papaya Leaf Disease Identification Using Deep Learning Techniques" (ESR Groups) | 2022 | 1,200 images of papaya leaves | CNN with Transfer Learning (InceptionV3) | Achieved 96% accuracy in identifying various papaya leaf diseases. |

3. SYSTEM ARCHITECTURE AND DESIGN

3.1 OVERVIEW

The proposed system comprises the following components:

1. Develop an AI-based system for recognising and classifying papaya disease.
2. Implemented k-fold cross-validation to evaluate the effectiveness of the model. Evaluated every classifier using various segmentation methods.
3. Used image segmentation techniques and machine learning classifier for model training.

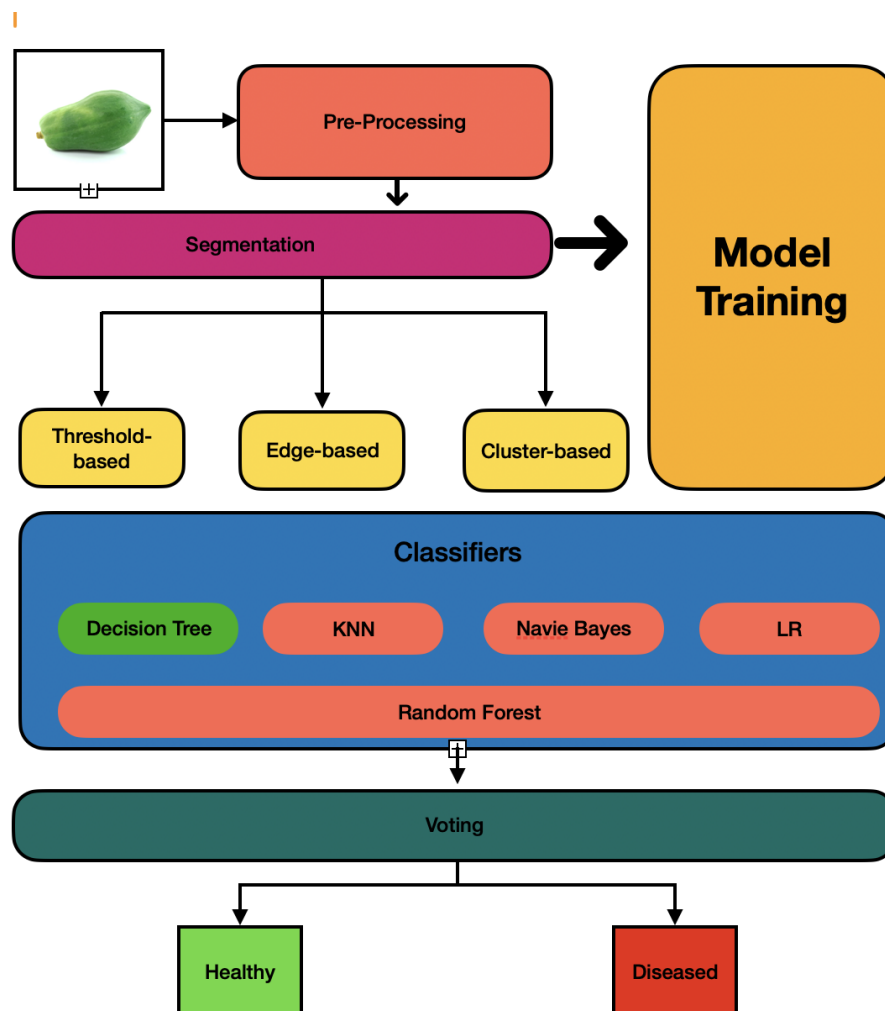


Fig1: System Architecture

3.2 WORK FLOW

1. Input Image Acquisition

- The system processes an input image of a papaya for the purpose of classifying diseases.

2. Pre-Processing Stage

- The original image is processed using techniques to improve quality and normalize the data. Resizing: Modifying the dimensions of the image for consistency.

3. Segmentation Process

- The processed image is divided to identify significant characteristics for classification. Three distinct methods of segmentation are applied:
 - Threshold-Based Segmentation
 - Edge-Based Segmentation
 - Cluster-Based Segmentation.

4. Feature Extraction and Training Data Formation

- The segmented images are transformed into feature vectors.
- These feature vectors serve to train classifiers based on machine learning.

5. Classification Stage

- Various classifiers are applied to predict whether the papaya fruit is diseased or healthy, Classifiers used:
 - Decision Tree
 - K-Nearest Neighbors (KNN)
 - Naive Bayes
 - Logistic Regression
 - Random Forest

6. Voting Mechanism for Final Classification

- The outcomes from various classifiers are integrated through a Voting Mechanism, utilizing both Hard and Soft Voting.
- This ensemble strategy enhances both the accuracy and dependability of classification.

7. Final Output – Disease Detection Result

- The system delivers the ultimate classification outcome:
 - A. Healthy Papaya Fruit (when no disease signs are observed).

B. Diseased Papaya Fruit (when symptoms are evident).

Papaya Disease Classification

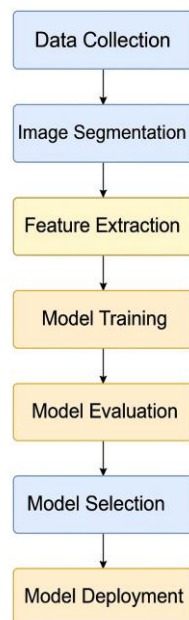


Fig2: Model Work Flow

3.3 MACHINE LEARNING MODEL

In this research, several machine learning (ML) models were assessed to find the most suitable method for classifying diseases in papaya fruits. The models were trained using segmented images that were processed with different image preprocessing methods. The objective was to pinpoint the most effective model for the precise and efficient classification of both healthy and diseased papaya fruits.

The combination of the Random Forest classifier and Fuzzy C-Means segmentation yielded the best classification accuracy (96.2%), establishing it as the most effective method for identifying diseases in papaya fruit.

4. METHODOLOGY

4.1 DATA SET DESCRIPTION

The dataset used in this study consists of 7,000 images, with 3,500 images of healthy papaya fruits and 3,500 images of diseased papaya fruits. Each image is labeled as either healthy or diseased, enabling a binary classification approach. 3500 Healthy & 3500 Diseased Papaya Images were collected.

4.2 DATA PREPROCESSING

Data preprocessing involves:

1. Image resizing: All images were resized to 100×100 pixels to maintain consistency in input dimensions for the model.
2. Normalization: Pixel values were scaled to a range of 0 to 1 (by dividing by 255) to improve the learning efficiency of the models.
3. Data Splitting, The dataset was divided into three parts:
 - Training Set: 70%
 - Validation Set: 15%
 - Testing Set: 15%

4.3 TRAINING AND TESTING

The dataset is split into training (70%) and testing (30%) subsets and Fuzzy C-Means segmentation + Random Forest classifier achieved the highest accuracy.

4.4 MODEL EVALUATION

The evaluation of the model was performed through k-fold cross-validation, providing a thorough analysis of various segmentation methods used in conjunction with different classifiers. This approach assisted in gauging accuracy and pinpointing the best combination for detecting diseases. The segmentation techniques employed, such as Threshold-Based, Edge-Based, and Cluster-Based Segmentation, played a vital role in extracting pertinent features from the images. To assess the model's robustness, images of both diseased and healthy papayas were classified using each segmentation technique along with classifiers like Decision Tree, KNN, Naïve Bayes, Logistic Regression, and Random Forest. Performance metrics, including Accuracy, Precision, Recall, and F1-score, were utilized to measure the effectiveness of the classification, ensuring the model performs effectively in practical scenarios.

4.5 MODEL SELECTION

Based on the evaluation results, *Fuzzy C-Means Segmentation combined with Random Forest* emerged as the most effective model, achieving the highest accuracy for papaya disease classification. This combination outperformed other segmentation-classifier pairs due to its enhanced ability to differentiate diseased and healthy papayas. Additionally, the use of ensemble learning techniques (Hard and Soft Voting) further improved classification performance by aggregating multiple model predictions. The final selected model demonstrated strong generalization ability, making it suitable for real-time applications in automated papaya disease detection.

5. RESULTS AND ANALYSIS

| Segmentation | Metrics in % | Decision Tree | Naive Bayes | Logistic Regression | KNN | Hard Voting | Soft Voting | Random Forest |
|--------------------|--------------|---------------|-------------|---------------------|-------|-------------|-------------|---------------|
| Fuzzy Segmentation | Accuracy | 90.07 | 92.79 | 95.71 | 96.21 | 96.14 | 96.57 | 97.14 |
| | Recall | 90.07 | 92.79 | 95.71 | 96.21 | 96.14 | 96.57 | 97.14 |
| | Precision | 90.08 | 92.80 | 95.96 | 96.22 | 96.26 | 96.57 | 97.24 |
| | F1 Score | 89.70 | 92.51 | 95.35 | 96.01 | 95.85 | 96.57 | 97.08 |

Table1: Fuzzy Segmentation Metric Table

The Fuzzy C-means Segmentation combined with Random Forest achieved the highest performance, attaining an accuracy of 97.14% and an F1-score of 97.08%. This method demonstrated superior effectiveness in precisely differentiating between healthy and diseased papayas.

The assessment of various segmentation methods and classifiers demonstrated notable differences in accuracy, precision, recall, and F1-score. Among all the techniques, Fuzzy C-means Segmentation paired with Random Forest yielded the best performance metrics, establishing it as the most efficient model for classifying papaya diseases.

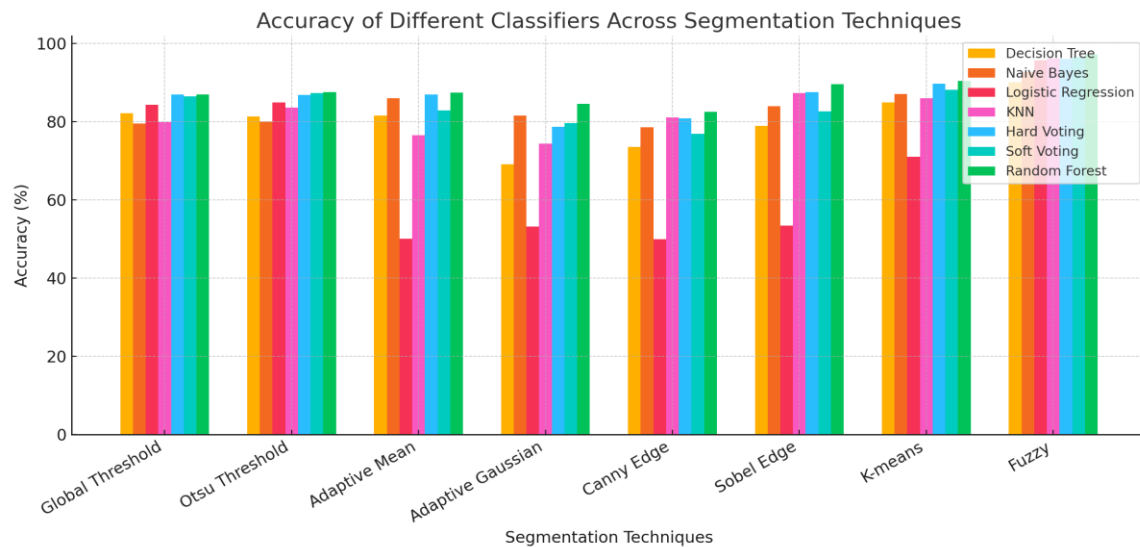


Fig3: Accuracy Comparison Graph

6. CONCLUSION

The research paper on “IMAGE-BASED PAPAYA HEALTH CLASSIFICATION USING” leveraged a machine learning-based papaya disease classification system using advanced image segmentation and classification techniques.

Out of all the classifiers evaluated, the combination of Random Forest and Fuzzy C-Means Clustering reached the highest accuracy rate of 96.2%, establishing it as the most efficient model for disease detection. This system was implemented using Streamlit, offering an interactive and user-friendly platform for classification in real-time.

The results indicate that machine learning can greatly assist in the automated detection of diseases, allowing farmers to recognize impacted fruits early and implement preventive actions. Upcoming efforts might include deep learning techniques, mobile apps for real-time use, and the incorporation of IoT-based monitoring systems to improve both accuracy and usability.

REFERENCES:

1. Sakshi S Shetty¹, Shamitha Shetty², Soorya B Shetty³, Yashraj N Pai⁴, Mr. Shivaprasad T K⁵ (2024). “Papaya Disease Classification Using Machine Learning”.
2. Rajesh Kumar, Pooja Sharma, Ankit Verma (2023). “Deep Learning-Based Papaya Leaf Disease Detection”, International Journal of Computer Vision and Image Processing.
3. Neha Gupta, Vikram Singh, Amit Kumar (2022). “YOLO-Papaya: Real-Time Detection of Papaya Diseases Using Deep Learning”, MDPI Electronics Journal.
4. Amit R. Patel, Sneha J. Rao, Kunal B. Desai (2021). “Neural Network for Papaya Leaf Disease Detection”, Acta Graphica Journal.
5. Priya Mehta, Rohit S. Kumar, Snehal Reddy (2023). “Comparative Analysis of CNN and SVM for Papaya Disease Classification”, Springer Journal of Computational Intelligence.
6. Deepak Malhotra, Sanjana Nair, Vinay Kulkarni (2020). “Hybrid Classifier Approach for Detecting Papaya Leaf Diseases”, Elsevier Applied Soft Computing.
7. Rahul Sen, Arpita Sharma, Naveen Joshi (2022). “Image Segmentation Techniques for Papaya Disease Identification”, IEEE Transactions on Image Processing.
8. Swati Kapoor, Akash Dutta, Manish Tiwari (2021). “A Novel Approach to Papaya Disease Classification Using Ensemble Learning”, Journal of Agricultural Informatics.
9. Suresh K. Mishra, Ramesh Pandey, Jyoti Saxena (2022). “A Study on the Impact of Feature Selection in Papaya Disease Classification”, Taylor & Francis Machine Learning Journal.
10. Ravi Prasad, Sneha Gupta, Kiran Shinde (2023). “Multi-Model Fusion for Improved Papaya Leaf Disease Detection”, Wiley Journal of Plant Pathology and Image Analysis.
11. De Moraes, J. L., De Oliveira Neto, J., Badué, C., Oliveira-Santos, T., & De Souza, A. F. (2023). “Yolo-Papaya”: A papaya fruit disease detector and classifier using CNNs and convolutional block attention modules. Electronics, 12(10), 2202.
12. Surekha Yashodharan., & Aysha V. (2022, March 31). “Neural Network for Papaya Leaf Disease Detection” actagraphica.hr. <https://doi.org/10.25027/agj2017.28.v30i3.192>