

Advancements in Machine Learning and Deep Learning for Plant Disease Detection: A Comprehensive Survey of Techniques, Datasets and Future Directions

Suma G.R

Assitant Professor

Dept. of Information Science & Engg
SSIT,SSAHE university,Tumkur

Dr. Ghouse Ahamed Z

Assitant Professor

Dept. of Electronics and Communication Engineering
SSIT,SSAHE university,Tumkur

Abstract— Global food security is increasingly threatened by plant diseases, which significantly reduce crop yields and negatively impact agricultural economies. Traditional manual detection methods are not scalable due to their reliance on expert knowledge and labor-intensive processes. This survey examines recent advancements in machine learning (ML) and deep learning (DL) techniques for automated plant disease detection [1], providing an in-depth analysis of the potential these technologies hold for transforming agricultural practices- The survey explores a variety of ML and DL approaches [2] that have been developed to detect plant diseases with high accuracy. It delves into the types of datasets used for training these models, highlighting the importance of large, diverse, and well-labeled datasets for improving model performance. Additionally, the survey reviews various performance metrics employed to evaluate the effectiveness of these techniques, such as accuracy, precision, recall, and F1 score. Looking towards the future, the survey suggests several research directions to enhance the efficacy of ML and DL in precision agriculture

Keywords— Plant Disease Detection, Plant Disease Classification, Machine Learning, Deep Learning, Image Processing.

INTRODUCTION

Agriculture is a foundational sector globally, critical to the sustenance and economic stability of societies. However, plant diseases remain a key factor threatening crop health and productivity, leading to significant losses in yield and financial detriment. Traditional manual disease identification methods, while effective to some extent, have substantial drawbacks in terms of scalability, cost, and time. These methods are labor-intensive and require a high level of expertise, which is not always available, particularly in large-scale farming operations. Consequently, their utility is limited, and they often fail to provide timely identification and intervention.

The advent of automated detection systems using machine learning (ML) and deep learning (DL) techniques [1] offers a promising solution to these challenges. By leveraging vast amounts of data and sophisticated algorithms, these technologies can transform plant disease detection. Automated systems can detect diseases early with high accuracy, enabling timely and effective interventions that are crucial for maintaining crop health and maximizing yield. This survey examines the advancements in ML

and DL for automated plant disease detection, providing a comprehensive overview of the current state of the field. It explores various approaches, including traditional ML techniques such as support vector machines (SVM), decision trees, and random forests, as well as more advanced DL methods [4] like convolutional neural networks (CNNs) and recurrent neural networks (RNNs). Each approach's strengths and weaknesses are analyzed, along with the types of plant diseases they are most effective in identifying. The survey also delves into the datasets used for training these models. High-quality datasets are essential for developing robust and accurate models. The survey highlights the importance of having large, diverse, and well-annotated datasets which can capture the variability in plant diseases across different regions and conditions. Additionally, the survey reviews various performance metrics used to evaluate these models, such as accuracy, precision, recall, F1 score, and the area under the receiver operating characteristic curve (AUC-ROC). Despite the significant progress made, several challenges remain. One major challenge is the generalizability of these models. Often, models trained on specific datasets do not perform well when applied to new data from different regions or under different environmental conditions.

Objective: This paper surveys recent advancements in machine learning (ML) and deep learning (DL) for plant disease detection, providing a comprehensive examination of the various models, datasets, performance evaluations, challenges, and opportunities within this rapidly evolving field.

METHODOLOGY

Search Strategy: To ensure a comprehensive review of the literature, we employed a meticulous search strategy [3] across several academic databases. The primary databases utilized were Science Direct and PubMed, known for their extensive repositories of scientific and medical research.

Our search terms were carefully selected to encompass a broad range of relevant topics within the intersection of agriculture and advanced computational techniques. The terms used included "plant disease detection," "machine

learning agriculture,"[10] and "deep learning for plant health." These keywords were chosen to capture studies focusing on the use of machine learning and deep learning technologies [5] in identifying and diagnosing plant diseases. By utilizing this rigorous search strategy, we aimed to compile a comprehensive and representative body of literature, providing a solid foundation for our review of the current state-of-the-art in plant disease detection using machine learning[10] and deep learning techniques.

Inclusion Criteria: This review specifically focused on research published [15] between 2018 and 2023, ensuring that the studies selected were relatively recent and reflective of the latest trends and advancements in the field. Priority was given to studies with direct practical applications in the agricultural sector, particularly those involving the implementation of machine learning (ML) and deep learning (DL) techniques [15] for disease detection in real-world settings. This selection criterion aimed to identify research that not only explored theoretical models but also showcased real-world case studies, field trials, and practical implementations that demonstrate the effectiveness of ML/DL-based solutions for early disease diagnosis, pest control, crop management, and overall agricultural productivity. Studies that involved large-scale datasets, sensor networks, and integration with farming technologies were highly valued. The goal was to capture evidence of how these technologies are being utilized to address current challenges in agriculture, such as increasing crop yields, minimizing losses, and enhancing sustainability through smarter disease monitoring and management

Data Analysis: Studies were classified based on ML/DL methods, dataset used, and reported accuracy, precision, and recall scores(survey).

LITERATURE REVIEW

4.1 Machine Learning Approaches

Feature-Based Classification: traditional machine learning (ML) models [4] M. Shrivastava and Pradhan proposed a "Machine Learning Paradigm for Rice Plant Disease Classification" such as Support Vector Machines (SVM) [14], k-Nearest Neighbors (KNN), and Random Forest, commonly employ a feature-based approach for disease classification in agriculture. In this framework, features like color, shape, texture, and other visual properties are manually extracted from leaf images or other plant parts to differentiate between healthy and diseased plants. These features are typically chosen based on domain knowledge, such as understanding the specific symptoms of a particular disease. For example, leaf spots, discoloration, or changes in texture may be indicative of certain fungal or bacterial infections.

Limitations: Despite their computational efficiency and simplicity, traditional machine learning (ML) models, such as Support Vector Machines (SVM), k-Nearest Neighbors (KNN), and Random Forest, face significant limitations when applied to more complex and real-world datasets. One of the primary challenges is their relatively lower accuracy when working with datasets that contain a high degree of variability or complexity. This is especially true in agricultural applications, where environmental conditions, lighting, plant varieties, and disease stages can vary considerably. These models, relying on manually

extracted features like color, shape, and texture, may struggle to capture the full complexity of disease symptoms across different conditions, leading to reduced classification performance. A critical limitation is their sensitivity to noise and image variations. In the field, images of plants [7] can be affected by factors such as varying lighting conditions, occlusions (e.g., overlapping leaves or plant structures), background clutter, and noise in the data due to sensor limitations. As a result, traditional ML models may misclassify images or fail to generalize well to unseen data, especially if the training data was collected in controlled environments that do not replicate real-world conditions. This sensitivity to noise makes these models less robust and less reliable in practical, dynamic settings. For example, Singh et al. (2020) highlighted this issue in their research, reporting a high classification accuracy of up to 94.6% when using traditional feature-based models on controlled datasets, where the environmental conditions were consistent and the images were relatively clean and well-annotated. However, they also observed a significant drop in performance when these models were applied to field conditions. In their field surveys, the accuracy decreased, indicating that the models struggled to adapt to the more variable and unpredictable nature of real-world agricultural environments. This drop in efficacy was attributed to factors such as inconsistent image quality, the presence of different types of diseases with overlapping symptoms, and the influence of changing environmental conditions.

This demonstrates a key limitation of traditional ML models: their reliance on handcrafted features and controlled environments often fails to account for the complexities of field conditions, leading to lower performance in practical applications. In addition, the models require constant tuning and adjustment of features, which can be labor-intensive and require ongoing domain expertise to maintain accuracy over time, espec proposed axially as the conditions change. This challenge has led to a growing interest in more advanced techniques, such as deep learning, which can automatically learn more robust and flexible features, making them better suited for real-world agricultural applications.

CNN Architectures: Convolutional Neural Networks (CNNs) [8] B. K. Mohanty, D. Salathé proposed a "Using Deep Learning for Plant Disease Detection," have revolutionized the field of image classification and are particularly powerful in disease detection tasks within agriculture. Unlike traditional machine learning models, CNNs, such as ResNet, VGG, and EfficientNet, do not rely on manually crafted features but instead learn spatial hierarchies directly from raw image data. This ability to automatically extract meaningful features from images allows CNNs to eliminate the need for domain-specific feature engineering, making them highly flexible and adaptable to a wide range of disease detection tasks.

ResNet (Residual Networks) is one of the most well-known CNN architectures[13], primarily designed to address the vanishing gradient problem in deep neural networks. ResNet uses residual connections, which allow information to bypass layers, facilitating the training of much deeper networks. This architecture has been particularly effective in image classification tasks, including plant disease detection, where it can learn to identify complex patterns

across multiple layers of abstraction. By capturing fine-grained details in early layers and high-level semantic features in deeper layers, ResNet is highly efficient at distinguishing between subtle differences in disease symptoms, even when the images are noisy or affected by environmental factors.

VGG (Visual Geometry Group) networks, known for their simple and uniform architecture, use a series of convolutional layers followed by fully connected layers. Although VGG networks tend to be computationally expensive, they have been widely used in many applications due to their straightforward design and effective performance. For agricultural disease detection, VGG's ability to learn hierarchies of spatial features from images enables it to accurately classify plant diseases by recognizing patterns in textures, colors, and shapes that are often present in diseased plants. EfficientNet is another CNN architecture [6] A. Chen et al. proposed a "DL for Tea Disease Recognition Using LeafNet," that has gained attention for its efficiency in terms of both computation and accuracy. EfficientNet introduces a compound scaling method, which uniformly scales the depth, width, and resolution of the network, allowing for better performance with network, allowing for better performance with fewer parameters. This makes EfficientNet particularly advantageous in resource-constrained environments, where computational power may be limited. Its high accuracy and efficiency make it a popular choice for real-time disease detection tasks, where fast and precise results are critical for timely intervention in agricultural practices.

One of the key advantages of CNN architectures like ResNet, VGG, and EfficientNet is their ability to automatically learn relevant features from image data, which reduces or completely eliminates the need for domain expertise in feature extraction. In traditional machine learning approaches, the selection and extraction of features often require significant knowledge of the specific disease symptoms and plant characteristics. However, CNNs can learn complex patterns and subtle details in the data that may not be easily identifiable by human experts. This self-learning capability is particularly beneficial in the context of plant disease detection, where visual symptoms can vary greatly across different plant species, growth stages, and environmental conditions. Furthermore, CNNs have been shown to perform exceptionally well in large-scale image datasets, which are common in agricultural applications. By leveraging large volumes of annotated images, CNNs can train on diverse datasets and improve their generalization capabilities, leading to better performance in real-world applications. For instance, in field settings, CNNs can recognize a broad spectrum of diseases, even when images are taken under varying light conditions or with different camera setups, making them far more robust and reliable compared to traditional feature-based models.

The ability of CNNs to handle complex, high-dimensional image data with minimal pre-processing makes them an attractive choice for modern disease detection systems in agriculture. Their flexibility, ability to reduce human intervention, and high accuracy in diverse conditions are key factors driving their widespread adoption in the industry.

Transfer Learning: Transfer learning [3] S.A. Ramcharan, B. Baranowski, et al., proposed a "Cassava Disease Classification using MobileNet," has become a critical technique in overcoming the challenge of limited labeled data in specific domains, such as agricultural disease detection. In transfer learning, a model that has been pre-trained on a large, diverse dataset (such as ImageNet, which contains millions of images from a wide range of categories) is fine-tuned or adapted for a specific task with a smaller, domain-specific dataset. This approach leverages the generalizable features learned by the model during its pre-training on a large dataset and applies them to new tasks, making it particularly useful when there is a lack of sufficient labeled data in specialized fields like agriculture. In plant disease detection, transfer learning allows for the use of pre-trained models, which have already learned to identify fundamental patterns in images—such as edges, textures, and shapes—which are common across many types of images, including plant leaf images. By fine-tuning these models on a smaller dataset specific to plant diseases, transfer learning enables effective disease detection without the need to train a model from scratch, which would require extensive computational resources and large amounts of annotated data.

For instance, Ramcharan et al. (2021) demonstrated the effectiveness of transfer learning in plant disease detection by applying MobileNet, a lightweight convolutional neural network (CNN), pre-trained on the ImageNet dataset, to cassava disease classification. They achieved an accuracy of 80.6%, which is a promising result considering that cassava disease datasets are relatively limited in comparison to broader image datasets like ImageNet. This shows that transfer learning can adapt well to agricultural applications, where gathering large, labeled datasets can be a major challenge. By transferring knowledge from a broad, generalized dataset to a more specific, niche dataset (cassava diseases, in this case), the model can efficiently classify plant diseases even with limited available data. The success of transfer learning in this context highlights several advantages. First, it reduces the need for large amounts of annotated data specific to agricultural diseases, a common obstacle in the field. Second, it accelerates model training and improves generalization, as the pre-trained model has already learned useful low-level features that are transferable across tasks. This enables faster development of disease detection systems without the time and resource-intensive process of training deep learning models from scratch. Additionally, pre-trained models like MobileNet are computationally efficient, which is essential when deploying these systems in resource-constrained environments, such as rural farms with limited access to high-performance computing resources.

I. DATASETS AND EVALUATION METRICS

Popular Datasets:

- PlantVillage: The Plant Village dataset [5] is a well-known open-source resource, particularly valuable for the study of plant disease detection using machine learning and computer vision. It consists of over 50,000 images across 38 plant species and various plant diseases, making it one of the largest publicly available datasets [18] for this field. These images cover a wide range of symptoms and disease stages, providing a comprehensive training resource for models focused on plant health monitoring.

CASCADA: This dataset provides a collection of real-world images depicting a variety of crop diseases across different environmental conditions, growth stages, and regions. This diversity makes it highly relevant for practical, real-world applications, such as the development of machine learning models for crop disease detection in agriculture. The images are captured under natural settings, ensuring that the dataset reflects the complexities and challenges faced by farmers, such as variations in lighting, weather, and disease progression.

However, despite its usefulness, the dataset is relatively small in size, which may limit the performance of machine learning algorithms trained on it. A larger dataset would allow for more robust model training and generalization. Furthermore, to enhance the dataset's coverage, additional augmentation techniques (such as rotating, cropping, and adjusting lighting) are necessary. These augmentations would help to artificially expand the dataset and improve its applicability for tasks such as disease identification, classification [11], and prediction under varying field conditions. Additional data collection or synthetic data generation could further bolster the dataset's scope and usability in real-world agricultural settings (survey).

V. EVALUATION METRICS

Accuracy, precision, and recall are foundational metrics [19] used to evaluate the performance of classification models, whether in binary or multi-class settings. Here's a deeper look at each:

1. Accuracy

• **Definition:** Accuracy is the ratio of correctly predicted instances (both true positives and true negatives) to the total number of instances in the dataset.

• **Formula:**
$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{Total Instances}}$$

$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{True Positives} + \text{True Negatives} + \text{False Positives} + \text{False Negatives}}$$

2. Precision

• **Definition:** Precision measures the proportion of positive predictions that are actually correct. In other words, it indicates how many of the instances the model predicted as positive truly belong to the positive class.

• **Formula:**

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

• **Use Case:** Precision is critical in situations where false positives are costly or undesirable. For example, in medical testing, predicting a disease when the person is healthy (false positive) might lead to unnecessary treatments.

• **Example:** In a spam email detection model, if the model identifies 80 emails as spam but only 60 of them are actually spam, the precision would be 75% (60/80).

3. Recall (Sensitivity or True Positive Rate)

• **Definition:** Recall measures the proportion of actual positives that are correctly identified by the model. It tells us how well the model detects the positive class.

• **Formula:**

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

• **Use Case:** Recall is important in situations where missing a positive instance is critical. For instance, in disease detection, it's better to have a higher recall (even at the cost of lower precision) to ensure that as many positive cases as possible are identified, even if it means some false positives.

4. Precision-Recall Trade-Off

• Often, precision and recall are inversely related: as you increase precision, recall tends to decrease, and vice versa. A model with high precision might miss some positives (low recall), while a model with high recall might generate many false positives (low precision).

• This trade-off is usually managed using the F1 Score, which is the harmonic mean of precision and recall:
$$F1 \text{ Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

• The F1 score provides a single metric that balances both precision and recall, useful when you need to optimize for both.

5. Accuracy, Multiclass Extension

• In a multi-class classification problem (with more than two classes), precision, recall, and F1 score can be computed for each class separately (often referred to as "one-vs-rest" approach), and the average across all classes can be taken:

• **Macro average:** The average of precision, recall, or F1 score across all classes without considering the class imbalance.

• **Weighted average:** The average where each class's score is weighted by the number of true instances in that class, making it sensitive to class imbalance.

These metrics are integral in evaluating and comparing the effectiveness of different models, especially in imbalanced datasets, where certain classes might be underrepresented. By using these metrics thoughtfully, you can make more informed decisions about model performance and application.

6. Intersection over Union (IoU) and Structural Similarity Index (SSIM)

are both important metrics used in image processing and computer vision tasks, but they serve different purposes and focus on different aspects of image quality and accuracy.

Intersection over Union (IoU):

- Purpose: IoU is primarily used in tasks like image segmentation to measure the accuracy of the predicted segmentation against the ground truth. It quantifies the overlap between the predicted segmentation mask and the actual mask.

- Calculation: IoU is calculated as the ratio of the intersection area of the predicted and ground truth masks to the union area of both masks. Mathematically:
$$\text{IoU} = \frac{\text{IntersectionArea}}{\text{UnionArea}}$$

$$U = \text{Union Area}$$

Structural Similarity Index (SSIM):

- Purpose: SSIM is a perceptual metric used to evaluate the structural similarity between two images. It is especially useful for assessing the quality of image preprocessing tasks like denoising, compression, and enhancement.

- a. Purpose: IoU is primarily used in tasks like image segmentation to measure the accuracy of the predicted segmentation against the ground truth. It quantifies the overlap between the predicted segmentation mask and the actual mask.

- b. Calculation: SSIM compares local patterns of pixel intensities, accounting for luminance, contrast, and structure. It is calculated using:
$$\text{SSIM}(x, y) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)}$$
 where:

- μ_x and μ_y are the mean pixel intensities of images x and y ,
- σ_x^2 and σ_y^2 are the variances,
- σ_{xy} is the covariance,
- C_1 and C_2 are constants to stabilize the division.

VI. CHALLENGES AND FUTURE DIRECTIONS

Data Constraints: One of the key challenges [20] in machine learning (ML) and deep learning (DL) [9] is the lack of large, diverse datasets, which directly impacts the performance, generalization, and robustness of models when deployed in real-world scenarios. ML and DL models typically require vast amounts of labeled data for training to capture the complexity and variety of inputs they might encounter in practical applications. Without sufficiently large and varied datasets, models may struggle with overfitting, biased predictions, or poor generalization to unseen data. This issue is particularly pronounced in specialized fields where acquiring real-world data is costly,

time-consuming, or even impractical. For instance, in healthcare, autonomous driving, or rare event prediction, collecting high-quality, diverse datasets is challenging due to privacy concerns, resource limitations, and the inherent rarity of the events themselves.

Potential Solutions:

1. **Synthetic Data Generation:** One promising solution to address data constraints is the use of synthetic data generation. Generative Adversarial Networks (GANs) are one of the most widely used techniques for generating synthetic data. By training a GAN to mimic the distribution of real-world data, it can create realistic, labeled datasets that may be used for training models in scenarios where real data is scarce or sensitive. This method is particularly useful in domains like image recognition, where synthetic images can be generated to supplement existing datasets.
2. **Data Sharing Initiatives:** Another approach to overcoming the lack of large, diverse datasets is through data-sharing initiatives.
3. **Augmentation and Transfer Learning:** In situations where acquiring new data is not feasible, techniques like data augmentation [9] and transfer learning can help to expand the utility of existing datasets.
4. **Environmental Variability:** Environmental variability refers to the numerous external factors that can introduce unpredictability or noise into the data,

VI. CONCLUSION

Machine learning (ML) and deep learning (DL) [2] have the potential to revolutionize plant disease detection by automating the identification and monitoring of plant health, which is critical for sustainable and efficient agricultural practices. These technologies enable farmers to detect diseases early, improve crop yields, and minimize the use of pesticides, leading to more environmentally friendly farming methods. While DL models [16] show high accuracy, they also require significant resources and large datasets. Innovations in data synthesis, model optimization, and real-time analysis are essential for advancing DL in agriculture. Future work in plant disease detection using machine learning (ML) and deep learning (DL) should indeed place a strong emphasis on model interpretability and scalability to drive the widespread adoption of these technologies, especially in farming communities. These two factors are critical in making ML and DL models more accessible, usable, and trusted by farmers, particularly those in resource-limited regions.

REFERENCES

1. Y. M. Alsakar et al., "Plant Disease Detection and Classification Using Machine Learning," IEEE Trans. Comput. Vision Agri., vol. 3, no. 2, pp. 130-139, 2023.
2. N. Jain, A. A. Pujari, and R. Kumar, "Automated Identification of Plant Leaf Diseases Using Transfer Learning," IEEE Access, vol. 9, pp. 225-238, 2022.
3. S. A. Ramcharan, B. Baranowski, et al., "Cassava Disease Classification using MobileNet," Frontiers Plant Sci., vol. 12, pp. 122-132, 2021.
4. M. Shrivastava and P. Pradhan, "Machine Learning Paradigm for Rice Plant Disease Classification," J. Plant Pathol., vol. 11, no. 4, pp. 304-316, 2020.
5. L. Lee et al., "Plant Leaf Disease Detection Using Deep Learning," Comp. Electron. Agric., vol. 177, pp. 45-55, 2019.
6. A. Chen et al., "DL for Tea Disease Recognition Using LeafNet," Symmetry, vol. 11, no. 3, pp. 1-11, 2019.
7. S. Khirade, "Plant Disease Detection Using Image Processing," IEEE Int. Conf. Comput. Vision, pp. 768-771, 2018.
8. B. K. Mohanty, D. Salathé, "Using Deep Learning for Plant Disease Detection," Front. Plant Sci., vol. 7, no. 2, pp. 1-10, 2018.
9. P. Ghosh and K. Singh, "Hyperspectral Imaging for Plant Disease Detection," IEEE Int. Conf. Hyperspectral Imag., pp. 102-110, 2019.
10. L. Wang et al., "Agricultural Crop Disease Identification with Machine Learning," IEEE Access, vol. 6, pp. 132-141, 2018.
11. A. Al Bashish, "Detection and Classification of Plant Diseases," Int. Conf. Sig. Imag. Proc., pp. 1-4, 2018.
12. J. Shrestha and G. Deepak, "Plant Disease Detection Using CNN," IEEE Signal Proc. Conf., pp. 102-107, 2020.
13. S. Singh, M. Bhowmik, "Detection of Potato Diseases with Multiclass SVM," IEEE Can. Conf., pp. 120-126, 2018.
14. P. Das et al., "Machine Learning in Crop Disease Detection," IEEE Technol. Innov. Agri., vol. 10, no. 2, pp. 140-151, 2020.
15. N. A. Sakr, "AI for Plant Disease Classification," Sensors, vol. 21, no. 11, pp. 1-14, 2021.
16. S. Devendran, "Efficient Plant Disease Detection using CNNs," Turk. Comput. Math., pp. 1-13, 2021.
17. K. Wang et al., "Image Quality Analysis with SSIM and IoU," IEEE Comput. Vision Agric., vol. 4, no. 2, pp. 1-12, 2019.
18. J. Pujari et al., "Detection of Fungal Diseases Using Image Processing," IEEE Trans. Agric., pp. 12-25, 2019.
19. M. Barbedo, "Challenges in Plant Disease Detection in Agriculture," Biosyst. Eng., vol. 9, no. 3, pp. 50-60, 2021.