# EARLY DETECTION OF PLANT DISEASES USING A HYBRID ENSEMBLE FEATURE SELECTION WITH DEEP NEURAL NETWORK FOR MODERN AGRICULTURE

# Rajendra Prasad Bellapu<sup>1</sup>, Ramashri Tiruamala<sup>2</sup>, Rama Naidu Kurukundu<sup>3</sup>

<sup>1,3</sup> Department of ECE, Jawaharlal Nehru Technological University Anantapur College of Engineering, Ananthapuramu, Andhra Pradesh, India
<sup>2</sup> Department of ECE, Sri Venkateswara University College of Engineering, Tirupati, Andhra Pradesh, India

Corresponding Author

# **Rajendra Prasad Bellapu**

Department of ECE, Jawaharlal Nehru Technological University Anantapur College of Engineering, Ananthapuramu, Andhra Pradesh, India *Email: rajendra.14ph0416@gmail.com* 

#### Abstract

Plants play a significant role in Indian agriculture as well as the economy of the country. However, the expected growth of plants are affected by diseases, which may cause complete damage to leaf, fruits, flower, and stem, which also leads to economic losses in agriculture. Therefore, plant disease detection is an essential task for improving crop quality and production process. Researchers developed popular techniques, namely Support Vector Machine (SVM) and Convolution Neural Network (CNN), to recognize plant diseases. However, the classification accuracy is diminished due to the high curse of dimensionality with redundant data. Feature selection techniques are developed to address these issues, but single feature selection techniques, namely ReliefF, F-score are unstable in nature, which affects the classification accuracy for various subsets of features. So, as to settle all these issues, a hybrid ensemble feature selection technique is introduced in this research study. The input images are pre-processed using a multi-scale retinex algorithm, where the segmentation of leaf images is carried out by using Kernel Fuzzy C Means (KFCM), and affected area segmentation is carried out by using the multilevel Otsu Thresholding technique. The features are extracted using a hybrid feature extraction technique, and optimal features are selected using the ensemble feature selection technique with Mutual Information (EFS-MI). Finally, Deep Neural Network (DNN) is developed to categorize the healthy and affected leaves of Plant Village Dataset (apple and potato) and collected dataset (rice and groundnut). The experimental results proved that the proposed DNN achieved 98.77% of accuracy while existing multi-class SVM (M-SVM) achieved 97.03% of accuracy on potato data.

**Key words:** Agriculture; Deep Neural Network; Hybrid Ensemble Feature Selection; Mutual Information; Plants; Redundant Data; Segmentation; Support Vector Machine

#### Introduction

Agriculture is seen as an important sector in the Indian economy and an income source for many individuals. Agriculture is the basic need for human existence. In developing countries like India, the production of food products such as wheat, fruits, and vegetables need to be maximized to meet human needs. In addition, the quality of the products must meet specific quality standards to maintain the safety of human health and welfare 1,2. The importance of plants has continued to evolve in modern life, and many researchers from various scientific disciplines connected directly or indirectly with plants. Plants influence the climate and the ecosystem. They have many uses, such as agriculture, environment, energy, health, medicine, etc. 3. However, farmers also face water shortages, natural disasters, plant diseases, and many other challenges. In general, the leaves of the plant are essential components and symbolize the properties of the whole plant 4; it is the first source for identifying most plant diseases. The approach for early diagnosis of plant diseases is a significant

move towards precision agriculture. Recognition of plant diseases at its early stages may increase the probability of recovery and reduce damage to crops 7. Using effective imaging techniques 8, Rice Blast (RB) and Bacterial Leaf Blight (BLB), early and late blight, and several other viruses caused by fungi and bacteria can automatically be detected. Therefore, this study focuses only on the recognition of plant diseases based on the properties of the leaves.

Precision farming uses the latest knowledge to improve decision making 9. Early and accurate identification is an essential aspect of disease monitoring 10. Leaf symptoms are essential information for identifying diseases that occur in different types of plants. Most of these diseases can be judged by the naked eye of a person skilled in the art based on their symptoms. However, plant disease identification is expensive due to the lack of specialists and high cost 11. In this regard, in collaboration with experts in the field of agriculture, computer researchers have proposed several algorithms for the automatic recognition of plant and fruit viruses 12. The significant details of plant leaf diseased images can be attained by using the learning algorithms. Leaves are segmented according to some common spectral characteristics of image objects like color, size, shape, spatial relationship with neighbouring pixels, and texture. In general, the strategy for processing images consists of two steps; Essential attributes are initially extracted from the input leaf images, and an effective classifier is used in a second step to classify the images into healthy and diseased images 13.

In this study, a multi-scale retinex algorithm is used to pre-process the input data, which is given as input to the segmentation process. Two kinds of segmentations are carried out in this study using KFCM and Multilevel Otsu Thresholding for leaf and affected areas of pre-processed data. To choose the best subset of features to enhance classification 7 accuracy, hybrid feature extraction techniques and ensemble feature selection methods are used. Finally, DNN is used as a classifier for plant disease classification. The experiments are carried out in terms of six significant parameters on the Plant Village dataset (apple and potato) as well as on the collected real-time dataset (rice and groundnut).

The organization of the paper includes: Section 2 presents the study of existing techniques from the year 2018 to 2020. Section 3 denotes the problem statement, where the explanation of the proposed methodology is given in Section 4. The validation of the proposed ensemble feature selection, and the proposed classifier on two datasets, are described in Section 5. Finally, the conclusion of the study, with its future work, is represented in Section 6.

# Literature Review

This section defines the study of existing plant leaf disease recognition techniques on a standard dataset. In addition, the advantages of existing techniques, along with their limitations, are also presented.

Khan, et al., 14 focused on the identification and classification of several fetal diseases built on correlation coefficients and depth features (CCDF). In the first step, with the use of a hybrid system, the contrast of the input image was initially increased, followed by the proposed CC-based partitioning approach that segregated infected areas from the backdrop. However, the CCDF method focused only on segmentation and feature selection techniques, which required feature extraction techniques for improving accuracy.

Sharma, et al., 15 developed the CNN model for full images as F-CNN and segmented images as S-CNN model. In contrast to the F-CNN model, the S-CNN model shows 98.6% better accuracy when tested on single data for 10 disease groups. A limitation of the S-CNN model, however, was that it might not work in regions with more than one symptom of the disease. Another drawback was that it was susceptible to the segmentation quality; it requires manual inputs for segmentation.

Khan, et al., 16 applied a new method to identify and detect apple disease, which includes three pipelines: pre-processing, segmentation of lesion, and feature extraction with classification. In the first stage, the apple plant leaf lesions were improved with a hybrid method, which was a combination of 3D filters, de-

correlation, 3D Gaussian filters, and 3D media filters. The lesions were then segmented by a correlation method, and their results were optimized. Finally, the properties of histogram features and color features were combined with a parallel fusion method. The extracted properties have been optimized by GA. Khan, et al., achieved better results in terms of accuracy and execution time than traditional SVM. In order to reduce the computational cost, deep learning techniques are required in this study.

Kamal, et al., 17 detected the plant disease using a model with depth-wise separable convolution architecture, which was based on leaf images. In this study, two versions of depth-wise separable convolution were considered. The model was trained and tested on a subset of 82,161 healthy images from the publicly available Plant Village dataset, which includes 55 different healthy and diseased plants. This model suffers in terms of execution speed due to its depth wise separable nature.

Akram, et al., 18 applied the CNN model to classify diseases of different fruits. First, the spatial capabilities were extracted with the assistance of previously trained spatial models such as AlexNet, and VGG, which were later refined with the help of transfer learning concepts. A multilevel fusing method based on the entropy-controlled threshold was also proposed and used to measure the average of the selected features. However, one of the difficulties with using the zero-fill concept when merging parallel features decreased the classification accuracy and also increased the computation time.

Chouhan, et al., 19 implementation of a method called Bacterial foraging based Radial Basis Function Neural Network (BRBFNN) to automatically identify and classify plant leaf diseases. Optimization of bacterial foraging was used to ensure optimal RBFNN weight, to identify areas infested with various diseases on plant leaves, and to increase the speed and accuracy for their classification. The region growth algorithm improved network efficiency by finding and grouping starting points with common points for the feature removal process. BRBFNN worked on fungal diseases of apple dataset collect from plant village dataset. BRBFNN was superior in computational efficiency for detecting and classifying diseases but only worked with fungal diseases.

Arsenovic, et al., 20 implemented a PlantDiseasesNet with a two-phase network containing PDNet-1 and PDNet-2. PDNet-1 was responsible for finding the leaves of plants according to species, while PDNet-2 was responsible for sorting the leaves of these plants. The trained model achieved 93.67% accuracy in the given dataset and proved effective in challenging environments. Accuracy has been improved through the usage of other sources of information such as weather, field location, and stage of the plant. However, the trained model was worked only on the single leaf features, which was the drawback of the study.

Al-bayati, et al., 21 exposed the illness of plant leaf using early and late fusion of two classifiers: Modified Optimized DNN (MODNN) with evolutionary grasshopper feature optimization (GOA), Speeded Up Robust Features (SURF), and Modified CNN (MCNN). Using the early and late fusion of classifiers, the accuracy of this method for Plant leaves has developed and enhanced. The experiments used the parameters such as Recall, Precision, F-measure, Error, and Accuracy to describe the validation of the model. From the results, it has seen using early fusion was better than late fusion due to the early integration of training instances while using the early fusion.

# **Problem Statement**

In the field of agriculture, global problems such as climate change, the appearance of many features on one single image, viewing images in complex backgrounds, and in multi-featured areas is not simple to recognize the disease from that single image. The appearance of the pathological area on the border of the image, the different properties of the leaf features, the presence of tone that changes in lighting, and the similarity of colors between different features are challenges. Due to the low contrast and noise, it is challenging to identify the micro-calcification and a large amount of valuable information in an image. The shape, colour, and texture of each lesion are different so, it is difficult for any model to detect and classify without proper preprocessing of plant diseased leaf images. Researchers using computer vision and machine learning methods have encountered the following problems in this area:

(a) During Pre-processing: Separating leaf elements from an image which contains different symptoms, and most of the background areas are similar,

(b) During Feature Extraction: Features define the exact characteristic of diseases so reducing the dimensionality of features for classification is a big task since unnecessary features require a lot of computing resources to process them and also affects the accuracy of the model.

(c) During Feature Selection: Selection of irrelevant features reduces the training speed of the model, reduces accuracy, and leads to an overfitting problem. The solutions for these challenges for any machine learning or deep learning model in a computer-based vision system are:

(a) The exact classification of symptoms is only performed by a clear distinction between affected pixels and healthy pixels, which is determined by the correct pre-processing step;

(b) Selecting the subset of features which are relevant using a relevant feature selection technique.

## **Proposed Methodology**

This section explains the solutions for the early diagnosis of plant diseases for increasing crop yields. In this study, two datasets, namely the Plant Village dataset, collected dataset for rice and groundnuts, are used, then KFCM and Multilevel Otsu Thresholding techniques are used for the segmentation process. Then, the features of the segmented images are extracted by using hybrid feature extraction techniques. Then, the relevant information is selected by ensemble feature selection techniques. Finally, the DNN model is used for the classification of plant diseases for selected diseases. **Figure 1** presents the proposed method's overall block diagram.



Figure 1 Overall Block Diagram of the Proposed Methodology

# 4.1. Dataset Collection

The proposed method uses collected datasets of rice and groundnut along with the Plant Village dataset (https://www.kaggle.com/emmarex/plantdisease) for leaf disease prediction. The rice images are collected from the Agriculture Research Station (ARS), Nellore, and groundnut are collected from the Regional Agricultural Research Station (RARS), Tirupati. The collected images are captured by using a Canon 32M pixel DSLR camera. In the standard dataset, many images for various types of classes (i.e., apple, potato, citrus, corn, tomato, etc.) are included. In this research study, two main classes of a standard dataset, namely apple and potato, are considered, where **Table 1** defines the dataset description in terms of the total number of images for each class. **Figure 2** illustrates the sample images of the Plant Village dataset, where **Figure 3** presents the sample images of the collected real-time dataset used in this research work.

-		
Standard Dataset / Collected Dataset	Classes	Total number of image
Apple	Scab	630
	Black Rot	621
	Cedar Rust	275
	Healthy	1645
Potato	Early Blight	1000

Late Blight

Table 1 Data Description of Standard and Collected Real Time Dataset

1000

	Healthy	152
Rice	Powdery Mildew	250
	Septoria	250
	Wheat Rust	250
	Healthy	250
Groundnut	Healthy	250
	Affected	250



Apple Scab



Potato Early Blight



Apple Black Rot





Apple Cedar Rust **Apple Healthy** 



Potato Healthy

Figure 2 Sample Images from Plant Village Dataset

The dimensions of the images used are 512×512, which is given as input for the pre-processing technique. In order to evaluate the proposed model, the model is tested and compared with both the standard dataset and manually collected dataset. Ground truth images for manually collected datasets are built carefully by using photoshop express software.



**Groundnut Rosette** 

Groundnut Alternaria spot



Groundnut Leaf scorch



**Groundnut Healthy** 

## Figure 3 Sample Images from Collected Real-Time Dataset

#### 4.2. Pre-Processing

The contrast of the input images are improved by applying the Multi-scale Retinex (MSR) algorithm. The weighted sum of the output Single Scale Retinex 22 is defined as the output of the MSR, where the formula for MSR is expressed in Eq. (1-3)

$$R_{MSR_i} = \sum_{n=1}^{N} w_n R_{n_i} = \sum_{n=1}^{N} w_n \left[ \log I_i(x, y) - \log(F_n(x, y) * I_i(x, y)) \right]$$
(1)

$$F_n(x, y) = C_n \exp[-(x^2 + y^2)/2\sigma_n^2$$
(2)

$$I_i(x,y) = S_i(x,y)r_i(x,y)$$
(3)

Where the input image on the i - th color channel is illustrated as  $I_i$ , normalized surround function is defined as F, illumination represents as  $S_i$ , scene reflectance is denoted as  $r_i$ , N denotes the number of scales,  $w_n$  is the weight of each scale, and the filter standard deviation is defined as  $\sigma$ .

The general appearance of the color elements of an image is determined by the notion of a gray world, which states that for an image with a color variation, the average value of the red component, green component, and blue component of the image must be measured up to the total value of gray. Images that violate the gray world's assumptions, e.g., images can be affected by a particular color, the retinex process described above results in grayscale images resulting in reduced color saturation. To address this issue, color restoration has been incorporated into MSR. It is proposed to change the multiplication by the color retrieval function in the MSR output. Compute the chromaticity coordinates by using Eq. (4)

$$I'_{i}(x,y) = \frac{I_{i}(x,y)}{\sum_{j=1}^{S} I_{j}(x,y)}$$
(4)

For the  $i^{th}$  color channel, where S is the number of spectral bands. S = 3 for RGB color space. The restored color MSR is given by Eq. (5)

$$R_{MSRCR_i}(x, y) = C_i(x, y)R_{MSR_i}(x, y)$$
(5)

Where,  $C_i(x, y) = f(I'_i(x, y))$  is the *i*<sup>th</sup> channel of the color restoration function (CRF). In this preprocessing technique, the CRF provides the finest overall colour restoration and is defined by Eq. (6)

$$C_i(x, y) = \beta \log[\alpha I'_i(x, y)]$$
(6)

Where  $\beta$  is a gain constant, and  $\alpha$  controls the strength of the nonlinearity. Figure 4 presents the sample images of pre-processed apple and groundnut images.



Figure 4 Sample Pre-processed Images using Multi-scale Retinex Algorithm

International Journal of Modern Agriculture, Volume 9, No.4, 2020 ISSN: 2305 - 7246

#### 4.3. Segmentation

After pre-processing, RGB (Red, Green, Blue) is converted into YCbCr (Luminance, Chrominance), where KFCM is used to process the cb plane, and Multilevel Otsu thresholding is used to process the cr plane. The reason for choosing YCbCr rather than RGB is that in image standards such as JPEG, MPEG1, MPEG2, and MPEG4, YCbCr easily gets rid of redundant colour information 23. Due to its transformation simplicity and clear separation of luminance and chrominance components makes YCbCr color space is used for segmenting the affected leaf region. In this study, segmentation is carried out by two techniques; KFCM is used for leaf segmentation, and multilevel Otsu thresholding is used for affected area segmentation, which is explained as below:

### 4.3.1. Kernelized Fuzzy C Means (KFCM)

For noiseless images, a traditional FCM gives better segmentation results. But, the FCM neglects to classify the noisy information because of the inconsistencies of the element information, which prompts to assigning the membership esteems to become erroneous. This is the fundamental reason behind inappropriate segmentation that happens during the processing of a noisy image by the FCM 24,25.

To overcome the concern of the standard FCM, the KFCM algorithm is designed. By using nonlinear mapping capacity, the KFCM changes over the input information in the plane of images into advanced dimensional element space. The complex and nonlinear distinct issue in the information plane can be changed over with the assistance of the mapping capacity into linearly separable in future space. At this point, the FCM can perform its task with the determined element space. The objective function of the KFCM is presented in the following Eq. (7).

$$J_m = \sum_{i=1}^{c} \sum_{k=1}^{N} u_{ik}^p \|\varphi(x_k) - \varphi(v_i)\|^2 = 2 \sum_{i=1}^{c} \sum_{k=1}^{N} u_{ik}^p \left(1 - k(x_k - v_i)\right)$$
(7)

Where *p* represents a real number, indicates fuzziness,  $u_{ik}^p$  is the data point membership  $\sum_{i=1}^{c} u_{ik} = 1$ and  $v_i$  is the cluster centroid. From the above Eq. (7), *c*, N denotes cluster and data point numbers. Here  $\varphi$  is the mapping function. Here, the linear high dimensional feature space is transformed from the non-linear mapping of the image plane by using the Gaussian Kernel Function (GKF) given by equation (8).

$$K(x, y) = \exp(-d(x, y)^2 / \sigma^2)$$
 (8)

By applying the KFCM, the pre-processed images are segmented from the background (i.e., leaf area segmentation), shown in **Figure 5**.



Figure 5 Segmented Leaf Images Using Kernelized Fuzzy C Means Algorithm

4.3.2. Multilevel Otsu Thresholding

International Journal of Modern Agriculture, Volume 9, No.4, 2020 ISSN: 2305 - 7246

The multilevel Otsu thresholding is a simple and effective algorithm for affected area segmentation of input data because it utilizes only the values of maximum variance of the classes. At first, the intensity level L of the normalized image  $I_{norm}$  is calculated by using Equation (9).

$$PH_i^e = \frac{H_i^e}{M}, \sum_{i=1}^M PH_i^e = 1, e = \begin{cases} 1,2,3 & \text{if } RGB\\ 1 & \text{if } grayscale \end{cases}$$
(9)

Where,  $PH_i^e$  is denoted as distribution probability,  $H_i^e$  is a pixel value between the intensity level from *i* to *e*, *M* is stated as the number of pixels in the normalized image  $I_{norm}$ , *IN* is denoted as intensity level ( $0 \le i \le L-1$ ), and *E* is indicated as image components (RGB or grayscale). In probability distribution, the histogram value is normalized using Equation (10).

$$w_o^E(th) = \sum_{i=1}^{th} PH_i^e, w_1^E(th) = \sum_{i=th+1}^{L} PH_i^e$$
(10)

Where,  $\sigma_1^e = w_0^e (\mu_0^e + \mu_T^e)^2$ ,  $\sigma_2^e = w_1^e (\mu_1^e + \mu_T^e)^2$ ,  $\mu_0^e$  and  $\mu_1^e$  are represented as the average rate for class variants 1 and 2,  $\sigma^{2e}$  is indicated as variants between the classes *C*,  $\sigma_1^e$  and  $\sigma_2^e$  are indicated as class variants 1 and 2 25. Hence, the objective function is estimated by using Equation (11).

$$J(th) = \max(\sigma^{2e}(th)), 0 \le th_i \le L - 1, i = 1, 2, 3, \dots K$$
(11)

Where,  $th = th_1, th_2, \dots, th_{K-1}$  is indicated as a vector that consists of multiple thresholds. By increasing the objective function, the Otsu between class variance function is maximized to obtain an optimal threshold level of depth gesture image for better segmentation. **Figure 6** shows the sample images for affected area segmentation on apple and groundnut images.





These segmented leaves and affected areas are given as input to the feature extraction for extracting the features.

#### 4.4. Feature Extraction

In this research study, hybrid feature extraction techniques include GLCM, color features, and Local Ternary Pattern (LTP), which is explained as follows:

**4.4.1. Gray Level Co-occurrence Matrix (GLCM)**: The spatial dependence of gray levels in segmented images are calculated by using GLCM; the total number of gray levels in the segmented images are equal to the total number of rows and columns. Co-occurrence matrices are constructed in four spatial orientations  $(0^{0},45^{0},90^{0} \text{ and } 135^{0})$ . In this study, various textural features are calculated using GLCM like Entropy 5, Sum entropy 8, Difference entropy, Inverse difference (INV), Inverse difference with normalized and Inverse difference moment normalized 27,29.

**4.4.2. Local Ternary Pattern (LTP)**: LBP is widely used in various computer vision applications due to its simplicity and reliability against the variations of illumination. However, the classification performance is limited because it is highly sensitive to noise 26. Therefore, LTP is designed to solve this problem, which encodes the third position with pixel variance. Small pixel differences are easily overwhelmed by noise. The middle pixel and its neighbouring pixel difference are encoded as trinary 27. To reduce the dimensionality, this trinary code has been divided into two binary codes: positive LBP and negative LBP 28.

**4.4.3.** Color Features: Color is one of the most used features because it provides minimal storage, high-speed retrieval, and simple computational processing. In addition, the color moment method has the smallest feature vector size and the least complexity in computation. Therefore, it can be seen as a suitable parameter for generating feature vectors that can be used for classification purposes 29. Information about the color distribution in the image can be obtained using lower-order moments. In this study, the vectors represent the properties of mean, variance, and skewness, i.e., first, second, and third-order, kurtosis, and standard deviation 30 in RGB, YCbCr, and HSV (Hue, Saturation, Value) planes.

# 4.5. Feature Selection

On a large dataset, the training method takes a long time due to unnecessary dimensional features and curses that lead to the degradation of model performance. The selection-based method is used to select the relevant subset features and rejects the unnecessary features in data to overcome these difficulties. Four types of feature selection are considered, including the filtered, wrapper, embedded, and hybrid approach. The main problem with filter-based feature selection is that they don't consider unnecessary data/redundancy in the selected features. The single filter-based method provides low classification 25 accuracy because it has biasness on the selected features. This study includes an ensemble feature selection technique to select the optimal subset of features that increase the prediction accuracy during classification to overcome these difficulties. Five main features are selected in this study that consist of Relief-F 31, Pearson Correlation coefficient (PCC) 32, F-score 33, Infinite Feature Selection (IFS) 34 and Term Variance (TV) 35. The proposed method is a collective approach that joins subsets gained from different filters using mutual information of feature-classes and features-features, as shown in **Figure 7**.



Figure 7 Proposed Hybrid Feature Selection Technique

In **Figure 7**, the subsets FS1, FS2, FS3, FS4, and FS5 are selected by the five feature selection techniques. According to the MI 36, the selected subset features are combined by the combiner that depends on

the features class and feature. A combiner which uses MI, takes into account all the first ranked features derived from subsets which are selected from feature selection techniques and generally use the same features; if all the ranked-features are the same, then optimal features are considered from the common features. Based on the experimental study, the method uses the value for the user-defined threshold value as  $\alpha = 0.75$ . In this ensemble method, the "combiner" plays a significant role in synchronizing the various feature selection approaches. This study focuses on minimizing the repetition of selected functions by incorporating feature classes and interactions between features using MI for the ensemble feature selection technique. **Table 2** provides the information for reduced feature-length on training and testing images using proposed ensemble feature selection techniques.

Datasets	Without Feature Selection		With Feature Selection		
	Training Images	Testing Images	Training Images	<b>Testing Images</b>	
Apple	2220*597	951*597	2220*179	951*179	
Potato	1506*597	666*597	1506*179	666*179	
Rice	120*597	80*597	120*179	80*179	
Groundnut	70*597	30*597	70*179	30*179	

Table 2 Proposed	<b>Ensemble Feature</b>	Selection Techni	ques for Optimal	Features Length

These optimal feature-length are given as input for final classification (DNN), which is explained in the next section.

# 4.6. Classification using Deep Neural Network

The affected regions are classified by using the autoencoder DNN classifier based on selected features. In this research, the proper selection of the autoencoder DNN will be a suitable solution for the classification process when there is no prior knowledge about the distribution data. An autoencoder DNA usually acts as a feedforward network and is not a pre-learning method with greedy level-by-level learning. Data in DNN flow from input to output without a loop function. The main advantage of the Autoencoder DNN classifier is the ability to reduce losses. The automatic encoder is a coding frame that, as shown in **Figure 8**, consists of a network of neurons with several hidden layers.



Figure 8 Structure of DNN

International Journal of Modern Agriculture, Volume 9, No.4, 2020 ISSN: 2305 - 7246

The SoftMax level uses a different trigger function where the non-linearity applied to the previous level can be different. The SoftMax activation function is expressed in Eq. (17).

$$h_{i}^{l} = \frac{e^{w_{ih}^{l}l-1+b_{i}^{l}}}{\sum_{j} w_{i}^{l}h^{l-1+b_{i}^{l}}}$$
(12)

Where  $w_i^l$  is  $i^{th}$  row of  $W^l$  and  $b_i^l$  is  $i^{th}$  bias term of the final layer. This research can employ  $h_i^l$  as an estimator of P(Y = i|x). Where Y is the connected label of input data vector x. In this case, four output neurons at the SoftMax layer can be interpreted to identify the plant disease on apple, potato, rice, and groundnut.

#### **Results and Discussion**

The proposed system is experimented with using MATLAB (version 2018a) with a 3.0 GHz Intel i3 processor, 1TB hard disc, and 8 GB of RAM. To validate the effectiveness of the proposed feature selection and classifier system, it is compared with the existing systems on the publicly available plant village dataset and collected dataset. The parameters such as Accuracy, sensitivity, F1-measure, specificity, Matthews Correlation Coefficient (MCC), and Threat score (TS) or critical success index (CSI) are used for validation, which are indicated in the Equation (13 - 18).

$$MCC = \frac{TP \times TN - FP \times FN}{\left(\sqrt{(TP + FN)(TP + FP)(TN + FN)(TN + FP)}\right)} \times 100$$
(13)

$$F - score = \frac{2TP}{2TP + FN + FP} \times 100 \tag{14}$$

$$Sensitivity = \frac{TP}{FN+TP} \times 100$$
(15)

$$Accuracy = \frac{TN+TP}{TN+TP+FN+FP} \times 100$$
(16)

$$Specificity = \frac{TN}{FP+TN} \times 100 \tag{17}$$

$$CSI = \frac{TP}{TP + FN + FP} \times 100 \tag{18}$$

Where True positive (TP) represents the affected plants correctly identified as plant disease, falsepositive (FP) illustrates the healthy plants incorrectly identified as plant disease. True negative (TN) describes the healthy plants correctly identified as healthy plants, and False negative (FN) presents the affected plants incorrectly identified as healthy plants.

#### 5.1. Performance Analysis of Proposed Ensemble Feature Selection

In this section, the performance of ReliefF, TV, IFS, PCC, F-score, and proposed ensemble or combiner feature selection along with DNN classifier are considered on Plant Village as well as collected datasets. **Figure 9** provides the comparison of experimental results of proposed EFS-MI and different feature selection techniques with the DNN classifier on the Apple dataset.

From **Figure 9**, it is evident that the performance of the proposed ensemble feature selection technique with mutual information is better than the single feature selection technique alone. For instance, the PCC achieved 84.33% of accuracy, 87.35% of MCC, and 76.94% of CSI, where the proposed EFS-MI achieved 98.87% of accuracy, 98.54% of MCC, and 97.96% of CSI on the apple dataset. The specificity and sensitivity of ensemble feature selection are 98.67% and 99.02%, where IFS achieved 97.62% of sensitivity and 98.13% of

specificity. The reason is that the ensemble feature selection works based on the MI for selecting the optimal features, where the PCC selects its optimal features based on only the correlation between the features.



# FEATURE SELECTION TECHNIQUES

Figure 9 Performance Comparison of Proposed EFS-MI on Apple Dataset.

Figure 10 describes the experimental results of the proposed ensemble feature selection with DNN on the Potato dataset.

EVALUATION METRICS (%) - 00 00 (%) - 07 00 (%) - 08 01 (%)	96.87 84.22 92.53 94.41 95.73	96.85 84.92 95.74 94.09 96.93 98.63	98.44 94.36 95.4 97.64 97.8	96.19 82.34 91.87 91.62 95.03	96.7 87.18 90.39 92.94 95.29 98.42	94.66 75.87 85.76 91.98 93.26 97.74
0 -	Accuracy	Sensitivity	Specificity	F1-Measure	МСС	CSI
ReliefF	96.87	96.85	98.44	96.19	96.7	94.66
PCC	84.22	84.92	94.36	82.34	87.18	75.87
F-score	92.53	95.74	95.4	91.87	90.39	85.76
TV	94.41	94.09	97.64	94.62	92.94	91.98
IFS	95.73	96.94	97.8	95.03	95.29	93.26
EFS-MI	98.77	98.63	98.84	96.62	98.42	97.74
LI		FEATUF	RE SELECTION T	ECHNIQUES		1

Figure10 Performance Comparison of Proposed EFS-MI with Existing Techniques on Potato Dataset.

While comparing with all single feature selection techniques, PCC achieved less performance, i.e., 84.22% of accuracy, 84.92% of sensitivity, 94.36% of specificity, 82.34% of F1-measure, 87.18% of MCC, and 75.87% of CSI. The reason is that PCC is sensitive to a linear relationship that provides poor performance than other feature selection techniques. The reliefF algorithm achieved 96.87% of accuracy, 98.44% of specificity, 96.85% of sensitivity, which is close to the performance of the proposed ensemble feature selection technique, i.e., 98.77% of accuracy, 98.84% of specificity, and 98.63% of sensitivity. However, the CSI of ReliefF is low (i.e.94.66%) than the CSI of the proposed ensemble technique (i.e.97.74%). This shows that the proposed ensemble feature selection technique achieved better performance on the potato dataset. **Figure 11** illustrates the performance of the proposed ensemble feature selection technique with DNN on the rice dataset.



Figure 11 Performance Comparison of Proposed EFS-MI with Existing Techniques on Rice Dataset.

From **Figure 11**, it is clearly shown that the proposed ensemble feature selection technique achieved higher performance than other single feature selection techniques on the rice dataset. For instance, the ensemble technique achieved 91% to 93% on sensitivity, specificity, MCC, F1-measure, and CSI, where the same method achieved only 89.67% of accuracy. While comparing with other single feature selection techniques, PCC achieved only 44% to 55% of CSI, sensitivity, F1-measure, and accuracy due to its weighting function on the extracted rice data features. The IFS and ReliefF feature selection technique achieved 89.67% accuracy and 92% F1-measure, where the proposed ensemble feature selection technique achieved 89.67% accuracy and 93.80% of F1-measure. However, the same proposed method achieved less performance on collected rice data than the Plant Village dataset of apple and potato, which is shown in Figure 9 and 10. The reason is that the collected dataset suffers from various lighting conditions, background blur, and different illumination conditions. However, the effective KFCM and multilevel Otsu thresholding are used in this research study; the characteristics of input collected images lead to low accuracy. **Figure 12** presents the experimental results of proposed ensemble feature selection techniques.



Figure 12 Performance Comparison of Proposed EFS-MI with Existing Techniques on Groundnut Dataset

The proposed ensemble feature selection technique achieved better performance on the groundnut dataset than the collected rice dataset. For instance, the same method achieved nearly 89% accuracy on the collected rice dataset, where it achieved nearly 96% accuracy on the groundnut dataset. The reason is that the rice dataset has multiple classes, but the groundnut dataset used has only two classes that lead to better performance using the ensemble feature selection technique on the groundnut dataset. Here, the single feature selection techniques also achieved higher performance in terms of all parameters. For instance, ReliefF, TV, F-score, PCC, and IFS achieved nearly 91% to 95% of accuracy, sensitivity, specificity, and F1-measure, then these techniques achieved nearly 87% to 92% of MCC and CSI. However, the ensemble feature selection technique achieved nearly 96% accuracy, sensitivity, and F1-measure and achieved nearly 93% of MCC and CSI.

# 5.2. Performance Analysis of Proposed DNN Classifier with Proposed Feature Selection Technique

International Journal of Modern Agriculture, Volume 9, No.4, 2020 ISSN: 2305 - 7246

In this section, the performance of the DNN classifier is tested with other classifiers, namely Multi-class Support Vector Machine (MSVM), Artificial Neural Network (ANN), K-Nearest Neighbour (KNN), and Random Forest (RF) with proposed feature selection techniques on both Plant Village dataset and collected dataset. These classifiers are also implemented with the proposed ensemble feature selection technique on the collected as well as on Plant Village datasets to test the performance of the DNN classifier. **Figure 13** illustrates the validated results of the proposed DNN classifier in terms of all parameters on the apple dataset. It compares the performance of (1) Proposed EFS-MI and MSVM (EFS-MI+MSVM), i.e., proposed EFS-MI for feature selection and MSVM for classification. (2) Proposed EFS-MI and ANN (EFS-MI+ANN). (3) Proposed EFS-MI and KNN (EFS-MI+KNN). (4) Proposed EFS-MI and RF (EFS-MI+RF). (5) Proposed EFS-MI and DNN (EFS-MI+DNN).



Figure 13 Performance Comparison of Proposed DNN classifier on Apple Dataset.

In the Plant Village Dataset for Apple, the performance of the proposed DNN classifier achieved better results, i.e., 98.87% of accuracy, 98.67% of sensitivity, 99.02% of specificity, 97.38% of F1-measure, 98.57% of MCC, and 97.96% of CSI. Among the other traditional classifiers, MSVM achieves better results, which is close to the proposed DNN classifier on all parameters. For instance, MSVM achieved 97.73% of accuracy, 98.07% of sensitivity, 98.51% of specificity, 97.29% of F1-measure, 97.75% of MCC, and 96.66% of CSI. But the ANN techniques achieve low performance than RF, KNN, and MSVM, i.e., 77.77% of CSI and 83.83% of F1-measure, where RF achieved 92.68% of CSI and 96.14% of F1-measure. The reason for achieving low performance by ANN is that it has a high computation burden, and it is sensitive to overfitting. These experiments proved that the DNN classifier achieved better performance by using the proposed ensemble feature selection technique on the apple dataset. **Figure 14** represents the validated results of the DNN classifier with other traditional classifiers on the potato dataset.



Figure 14 Performance Comparison of Proposed DNN classifier on Potato Dataset.

From **Figure 14**, it is evident that the traditional ANN and KNN achieved poor performance on the potato dataset than MSVM, DNN, and RF. For instance, KNN achieved 88% of CSI, 93.59% of accuracy, and 93.08% of F1-measure, where ANN achieved 76.81% of CSI, 85.05% of accuracy, and 83.38% of F1-measure. The reason is that KNN didn't learn the features from the training data; it works on the principle of nearest neighbours, and ANN has more computation burden that leads to poor performance than other classifiers. Even though the ensemble feature selection technique is used with all classifiers, the MSVM and RF achieved low performance than the proposed DNN classifier. For instance, the MSVM and RF achieved 97.03% and 96.60% of accuracy, where the proposed DNN classifier achieved 98.77% of accuracy on the potato dataset. **Figure 15** explains the validated results of the proposed DNN classifier on rice data.



Figure 15 Performance Comparison of Proposed DNN classifier on Rice Dataset.

The proposed DNN classifier has a low accuracy (89.67%) than other parameters, namely sensitivity (93.87%), F1-measure (93.80%), MCC (91.74%), and CSI (92.63%). The reason is that collected images on rice from the Nellore region, which is subjected to various challenges, include blurring the background and low lighting conditions. When compared with other classifiers, ANN achieved deficient performance on rice data, i.e., 59.30% of accuracy, 54% of sensitivity, 53.55% of F1-measure, 69.41% of MCC, and 48.89% of CSI. The performance of the RF is close to the performance of MSVM. For instance, these two methods achieved nearly 88% to 89% of accuracy and 91% to 93% of sensitivity, MCC, CSI, F1-measure, and specificity. The reason is that the use of the ensemble feature selection technique with MSVM and RF. But, the performance of the proposed DNN classifier is high than other classifiers on the collected rice dataset even. **Figure 16** provides the performance of the proposed DNN classifier with an ensemble feature selection technique on the groundnut dataset.



Figure 16 Performance Comparison of Proposed DNN classifier on Groundnut Dataset.

The experimental results proved that the proposed DNN achieved nearly 96% of accuracy, sensitivity, specificity, and F1-measure and nearly 93% of MCC and CSI on groundnut data. All the existing classifiers, namely M-SVM, ANN, KNN, and RF are also achieved nearly 91% to 95% of accuracy, sensitivity, specificity, and F1-measure, where these techniques achieved 89% to 92% of MCC and CSI. When compared with the rice dataset, the performance of the proposed DNN achieved better performance. This is because the groundnut has only two classes, i.e., Healthy and Affected. Therefore, the proposed DNN classifier achieved higher performance than other traditional classifiers on the groundnut dataset.

# 5.3. Comparative Analysis of Proposed Feature Selection with Classifier

In this section, the performance of DNN classifier with ensemble feature selection technique is compared with existing techniques, namely CCDF 14, GA with M-SVM 16, M-SVM 18, PDNet1+PDNet 2 20, and MODNN+MCNN 21 in terms of accuracy on Plant Village Dataset. **Table 3** provides the comparative results of proposed DNN with existing classifiers.

From the **Table 3**, it is clearly evident that the proposed DNN classifier achieved better classification accuracy on both apple and potato plant diseases than CNN and M-SVM. The existing M-SVM with GA 16 on the apple dataset achieved 98.00%, while the existing M-SVM without feature selection technique 18 achieved 97.80% of accuracy. However, the computation complexity was high in the M-SVM with GA, which requires deep learning techniques. Whereas, deep learning techniques, including the MCNN 21, and PDNet1+PDNet2 20, achieved only 93% and 96.84% of accuracy. The reason is that PDNet has no feature selection techniques, and MCNN has a single feature selection technique (SURF).

Author	Methodology	Plant Village Dataset	Accuracy (%)
Khan, et al., 14	CCDF	Apple	96.90
Khan, et al., 16	GA+M-SVM	Apple	98.00
Akram, et al., 18	M-SVM	Apple	97.80
Arsenovic, et al., 20	PDNet1+PDNet2	Apple	93.00
Al-bayati, et al., 21	MODNN+MCNN	Apple	96.84
		Potato	95.71
Proposed	Ensemble feature selection technique + DNN	Apple	98.87
		Potato	98.77

Table 3 Comparative Analysis of Proposed EFS-MI with DNN in terms of accuracy

In this proposed DNN classifier, ensemble feature selection techniques are used for plant disease classification, and hence it achieved 98.87% of accuracy on apple data and 98.77% of accuracy on potato data. This proves that the proposed DNN with ensemble feature selection technique achieved better performance in terms of accuracy on Plant Village Dataset (i.e., apple and potato data)

# Conclusion

In this research study, an optimized, automated computer-based method is developed for plant disease recognition. Five major steps are presented in this study that consists of pre-processing, image segmentation, hybrid feature extraction, hybrid ensemble feature selection, and deep learning classifier. In the first step, the input data are pre-processed by multi-scale retinex algorithm for image enhancement. Then, these enhanced images are given as input to the segmentation process. Two kinds of segmentations are carried out in this study: KFCM is used for background subtraction, and Multilevel Otsu Thresholding is used for segmenting the affected area of input leaves. Hybrid feature extraction techniques, namely GLCM, color features, and LTP, are used to extract the essential features of the segmented data. Then, the ensemble feature selection technique includes reliefF, TV, IFS, TV, and PCC are used to select the optimal subset of features for improving the classification accuracy. These input subset features are fed into the stack-autoencoder DNN for the final classification of plant diseases. The experiments are conducted on the Plant Village dataset (apple and potato) and collected dataset (rice and groundnut) in terms of accuracy, specificity, MCC, CSI, sensitivity, and F1measure. The results proved that the ensemble feature selection technique achieved 98.87% of accuracy on apple, 98.77% of accuracy on potato, 89.67% of accuracy on rice, and 96.60% of accuracy on groundnut data. In addition, the proposed DNN classifier achieved 96.60% of accuracy and 96.61% of F1-measure, where the RF achieved 95.33% of accuracy and 94.27% of F1-measure on groundnut data. However, the proposed ensemble feature selection and DNN achieved less accuracy (i.e., 89.67%) on collected rice data than collected groundnut data due to background blur and illumination conditions on rice data. Therefore, an effective image preprocessing and segmentation techniques are needed to solve the issues on collected rice data as future work.

# References

- 1. Cristin, R., Kumar, B.S., Priya, C. and Karthick, K., 2020. Deep neural network-based Rider-Cuckoo Search Algorithm for plant disease detection. *Artificial Intelligence Review*, pp.1-26.
- 2. Savary, S., Ficke, A., Aubertot, J.N. and Hollier, C., 2012. Crop losses due to diseases and their implications for global food production losses and food security.
- 3. Kazerouni, M.F., Saeed, N.T.M. and Kuhnert, K.D., 2019. Fully-automatic natural plant recognition system using deep neural network for dynamic outdoor environments. *SN Applied Sciences*, *1*(7), p.756.
- 4. Andrushia, A.D. and Patricia, A.T., 2020. Artificial bee colony optimization (ABC) for grape leaves disease detection. *Evolving Systems*, *11*(1), pp.105-117.
- Alafeef, M., Fraiwan, M., Alkhalaf, H. et al. Shannon entropy and fuzzy C-means weighting for AI-based diagnosis of vertebral column diseases. J Ambient Intell Human Comput 11, 2557–2566 (2020). https://doi.org/10.1007/s12652-019-01312-3.
- 6. Abu-Naser, S.S., Kashkash, K.A. and Fayyad, M., 2010. Developing an expert system for plant disease diagnosis.
- Jahanjoo, A., Naderan, M. & Rashti, M.J. Detection and multi-class classification of falling in elderly people by deep belief network algorithms. J Ambient Intell Human Comput 11, 4145–4165 (2020). https://doi.org/10.1007/s12652-020-01690-z
- Lee, S. Using entropy for similarity measures in collaborative filtering. J Ambient Intell Human Comput 11, 363–374 (2020). https://doi.org/10.1007/s12652-019-01226-0
- 9. Gebbers, R., and Adamchuk, V.I., 2010. Precision agriculture and food security. *Science*, 327 (5967), pp.828-831.
- Sharif, M., Khan, M.A., Iqbal, Z., Azam, M.F., Lali, M.I.U. and Javed, M.Y., 2018. Detection and classification of citrus diseases in agriculture based on optimized weighted segmentation and feature selection. *Computers and electronics in agriculture*, 150, pp.220-234.
- 11. Samajpati, B.J. and Degadwala, S.D., 2016, April. Hybrid approach for apple fruit diseases detection and classification using random forest classifier. In 2016 International Conference on Communication and Signal Processing (ICCSP) (pp. 1015-1019). IEEE.
- 12. Gavhale, K.R. and Gawande, U., 2014. An overview of the research on plant leaves disease detection using image processing techniques. *IOSR Journal of Computer Engineering (IOSR-JCE)*, *16*(1), pp.10-16.
- 13. Kumar, S., Sharma, B., Sharma, V.K., Sharma, H. and Bansal, J.C., 2018. Plant leaf disease identification using exponential spider monkey optimization. *Sustainable computing: Informatics and systems*.
- 14. Khan, M.A., Akram, T., Sharif, M., Awais, M., Javed, K., Ali, H. and Saba, T., 2018. CCDF: Automatic system for segmentation and recognition of fruit crops diseases based on correlation coefficient and deep CNN features. *Computers and electronics in agriculture*, 155, pp.220-236.
- 15. Sharma, P., Berwal, Y.P.S. and Ghai, W., 2019. Performance analysis of deep learning CNN models for disease detection in plants using image segmentation. *Information Processing in Agriculture*.
- Khan, M.A., Lali, M.I.U., Sharif, M., Javed, K., Aurangzeb, K., Haider, S.I., Altamrah, A.S. and Akram, T., 2019. An optimized method for segmentation and classification of apple diseases based on strong correlation and genetic algorithm based feature selection. *IEEE Access*, 7, pp.46261-46277.
- 17. Kamal, K.C., Yin, Z., Wu, M. and Wu, Z., 2019. Depthwise separable convolution architectures for plant disease classification. *Computers and Electronics in Agriculture*, *165*, p.104948.
- 18. Akram, T., Sharif, M. and Saba, T., 2020. Fruits diseases classification: exploiting a hierarchical framework for deep features fusion and selection. *Multimedia Tools and Applications*, pp.1-21.
- 19. Chouhan, S.S., Kaul, A., Singh, U.P. and Jain, S., 2018. Bacterial foraging optimization based radial basis function neural network (BRBFNN) for identification and classification of plant leaf diseases: An automatic approach towards plant pathology. *IEEE Access*, 6, pp.8852-8863.

- 20. Arsenovic, M., Karanovic, M., Sladojevic, S., Anderla, A. and Stefanovic, D., 2019. Solving current limitations of deep learning based approaches for plant disease detection. *Symmetry*, *11*(7), p.939.
- 21. Al-bayati, J.S.H. and Üstündağ, B.B., Early and Late Fusion of Deep Convolutional Neural Networks and Evolutionary feature optimization for Plant leaf illness recognition. *Journal of Xi'an University of Architecture & Technology* 12(2):1591-1610, 2020.
- 22. Petro, A.B., Sbert, C. and Morel, J.M., 2014. Multiscale retinex. Image Processing On Line, pp.71-88.
- 23. Kolkur, S., Kalbande, D., Shimpi, P., Bapat, C. and Jatakia, J., 2016, December. Human Skin Detection Using RGB, HSV and YCbCr Color Models. In *International Conference on Communication and Signal Processing 2016 (ICCASP 2016)*. Atlantis Press.
- Wu, C., Li, Y., Zhao, Z. et al. Research on image classification method of features of combinatorial convolution. J Ambient Intell Human Comput 11, 2913–2923 (2020). https://doi.org/10.1007/s12652-019-01433-9
- 25. Fachrurrozi, M., Dela, N.R., Mahyudin, Y. and Putra, H.K., 2019, March. Tongue Image Segmentation using Hybrid Multilevel Otsu Thresholding and Harmony Search Algorithm. In *Journal of Physics: Conference Series* (Vol. 1196, No. 1, p. 012072). IOP Publishing.
- 26. Ren, J., Jiang, X. and Yuan, J., 2013, September. Relaxed local ternary pattern for face recognition. In 2013 IEEE international conference on image processing (pp. 3680-3684). IEEE.
- 27. Tan, X. and Triggs, B., 2010. Enhanced local texture feature sets for face recognition under difficult lighting conditions. *IEEE transactions on image processing*, *19*(6), pp.1635-1650.
- 28. Chahi, A., Ruichek, Y. and Touahni, R., 2018. Local directional ternary pattern: A new texture descriptor for texture classification. *Computer vision and image understanding*, *169*, pp.14-27.
- 29. Patil, J.K. and Kumar, R., 2011. Color feature extraction of tomato leaf diseases. *International Journal of Engineering Trends and Technology*, 2(2), pp.72-74.
- 30. Banu, M.S. and Nallaperumal, K., 2010, December. Analysis of color feature extraction techniques for pathology image retrieval system. In 2010 IEEE International Conference on Computational Intelligence and Computing Research (pp. 1-7). IEEE.
- 31. Urbanowicz, R.J., Meeker, M., La Cava, W., Olson, R.S., and Moore, J.H., 2018. Relief-based feature selection: Introduction and review. *Journal of biomedical informatics*, 85, pp.189-203.
- 32. Saeys, Y., Abeel, T. and Van de Peer, Y., 2008, September. Robust feature selection using ensemble feature selection techniques. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases* (pp. 313-325). Springer, Berlin, Heidelberg.
- Polat, K. and Güneş, S., 2009. A new feature selection method on classification of medical datasets: Kernel F-score feature selection. *Expert Systems with Applications*, 36(7), pp.10367-10373.
- 34. Roffo, G., Melzi, S. and Cristani, M., 2015. Infinite feature selection. In *Proceedings of the IEEE International Conference on Computer Vision* (pp. 4202-4210).
- Yin, S. and Jiang, Z., 2009, March. A Variance–Mean Based Feature Selection in Text Classification. In 2009 First International Workshop on Education Technology and Computer Science (Vol. 3, pp. 519-522). IEEE.
- 36. Kraskov, A., Stogbauer, H. and Grassberger, P., 2004. Estimating mutual information. *Physical review E*, *69*(6), p.066138.