



A Method Based on Artificial Intelligence That Predicts Arabica Coffee Yield by Analyzing Abiotic Factors and the Prevalence of Coffee Leaf Rust

C S Santhosh¹, Umesh K K²

¹Assistant Professor, Department of Computer Applications, JSS Science and Technology University, Mysuru, Karnataka, India.

²Associate Professor, Department of Information Science and Engineering, JSS Science and technology University, Mysuru, Karnataka, India.

Email Id: sancs84@jssstuniv.in¹, umeshkatte@jssstuniv.in²

Abstract

Coffee is a perennial crop that harbors infections throughout the plant, which may worsen illness under favorable circumstances. Coffee leaf rust, a widespread coffee disease and Arabica is susceptible to leaf rust. Disease incidence and severity depend on abiotic variables. Cloudy and constant South-West monsoon weather (June–September) promotes coffee leaf rust growth. In this five-model study, an Extra Tree and Gradient Boosting regression model predicted coffee crop output in Chikamagaluru, Karnataka, with the least error utilizing biotic and abiotic factors. We investigated additional tree, gradient boosting, RF, Decision Tree, and KNN models using biotic and abiotic predictors. Used the independent testing dataset's MSE, MAE, RMSE Root mean square errors, and R-squared errors to compare model performance. The extra tree ($R^2=0.98$ kg/ha⁻¹ and RMSE = 7.96 kg/ha⁻¹) and gradient boosting ($R^2=0.96$ kg/ha⁻¹ and RMSE = 10.96 kg/ha⁻¹) regression models used Group 1 and 2 characteristics as predictor variables and different parameter fine tuning functions to estimate coffee yield most accurately. Compared to the less precise probabilistic models utilized in this work, such as Random Forest, decision tree, and KNN models, shown in the results section. The optimum weather parameter for coffee production forecasts and biotic-CLR incidence data outperformed random forest, decision tree, and K-Nearest Neighbor models.

Keywords: Abiotic Variables; Coffee; Coffee Leaf Rust (CLR); Stochastic Regression Models, Yield Prediction.

1. Introduction

The taste of coffee stimulates millions worldwide. Coffee is second for world commerce after petroleum. Over 80 countries grow the energizing beverage, some of which are its main producers. These 24 countries generate over 50,000 MT of coffee, including India. Coffee is cultivated largely in hilly regions of Karnataka, Kerala, Tamil Nadu, and the southern states, and to a lesser extent in non-traditional locations like Andhra Pradesh, Orissa, West Bengal, Maharashtra, and the north and east [1]. Arabica and Robusta are two of the most prominent coffee varieties farmed in India, out of 103 marketed worldwide [2]. Most coffee growers in Karnataka cultivate C.Arabica in Chikamagaluru and C.Robusta in Coorg. Productivity might fall by

10-20% by 2050 owing to agro-ecological conditions in these places [3]. Juliana Jaramillo et. al., (2009) have assessed the heat tolerance of Hypothenemus hampei, the most destructive coffee pest, and examined data from Ethiopia, Kenya, Tanzania, and Colombia to determine the consequences of climate change. This research looked at the bionomics of Helicobacter hampei under eight different temperature regimes: 15, 20, 23, 25, 27, 30, 33, and 35. We look at how different adaptation strategies have fared in the face of climate change and how it has affected the spread of H. hampei [4]. M. de C. Alves et. al., (2011) Using climate and crop distribution, geoinformation-based prediction models identified soybean rust, coffee



leaf rust, and black Sigatoka risk zones in Brazil. They classified three plant diseases by temperature and leaf wetness using a meteorological model [5]. Jaramillo, J et. al., (2011) produced maps of expected *H. hampei* distributions in East African coffee-producing areas to estimate hazards and prioritize management operations. For HADCM3, the CLIMEX model links insect distributions to present climate and estimates future climatic envelopes under scenarios A2A and B2B [6]. Cora B. Pérez-Arizac et. al., (2012) an agricultural case study using Bayesian networks (BNs) for coffee rust prediction is presented. An 8-year Brazilian experimental farm dataset was utilized. Pre-processing of the original dataset was influenced by preliminary data analysis [7]. Kutuywayo, D et. al., (2013) this study studies agricultural pest distribution under anticipated climatic scenarios, concentrating on Zimbabwe's African coffee white stem borer (CWB) [8]. Classen, A et. al., (2014) the authors conducted an experiment to examine how land-use intensification affects Mount Kilimanjaro coffee production pollination and pest control. The study tested ordinary home gardens, shady coffee plantations, and sun coffee plantations (total sample size $\frac{1}{4}$ 180 coffee bushes) for pollinators and vertebrates across a land-use gradient. Researchers found no significant ecological service decrease across land-use gradients [9]. Wang, N et. al., (2014) this research examined yields and Central, North, East, Southwest, and Northwest coffee-growing regions' 254 plots' production metrics. They employed boundary line analysis to quantify regional yield disparities and determine how production characteristics affect coffee output. They also examined how rainfall fluctuation affects coffee output using regression analysis [10]. Corrales et. al., (2015) this research forecasts disease and pest crops using supervised learning. This study reviews supervised learning methods for maize, rice, coffee, mango, peanut, and tomato pesto and disease detection. Finding the best agricultural algorithms is the objective here [11]. Hameed et. al., (2020) the study examined the correlation between agricultural and environmental parameters and final

cup quality features. A study found that agricultural and environmental variables significantly affect the physical and biochemical qualities of coffee [12]. Sudha, Met. al., (2020) this research examined the impact of meteorological factors on CLR incidence at the Central Coffee Research Institute (CCRI) in Chikkamagaluru District, Karnataka, India. To control for interspecific hybrid Sln.5B and Robusta cultivar C×R and used the leaf rust-resistant Arabica coffee selection Sln.3. CLR observed at CCRI farm in 2015-16, 2016-17, 2017-18, and 2018-19 seasons in Coffee Arabica L cultivars Sln.3 & Sln.5B and C. canephora C×R. CCRI's meteorological observatory recorded maximum and minimum temperature, relative humidity, and rainfall [13]. Yáñez-López et.al, (2012) this re-view provides a summary of current climate change studies. This study analyzed climate change studies on plant diseases due to its possible influence on agricultural disorders. This review examines how climate change affects plant diseases, agricultural growth, development, and productivity [14]. Suresh, N et. al., (2012) this research examines the interaction between the host, environment, and pathogen in the coffee-leaf rust disease complex. Climate change affects variables affecting climate, but not disease development, since spore germination needs lower temperatures (15-20oC) and dim light. To ensure sustainable coffee production, breeding should prioritize the host plant's capacity to withstand pathogen attacks [15]. Cerda et. al., (2017) aimed to assess coffee pests and diseases' main and secondary yield losses and their drivers. We created a full-sun coffee package with six treatments, each including a distinct pesticide application sequence. The three-year study (2013-2015) evaluated yield components, dead productive branches, and foliar pests and diseases as yield predictors [16]. Abreu Junior et. al., (2022) This study used multispectral images and machine learning to predict coffee crop productivity. Sentinel 2 images, Random Forest (RF), Support Vector Machine (SVM), Neural Network (NN), and Linear Regression (LR) provided data from the same study location in 2017, 2018, and 2019. Using 85% training and 15%



validation data, the NN algorithm calculated yield best with 23% RMSE, 20% MAPE, and R2 0.82. NN predicted yields better (27% RMSE) [17]. De Leijst-er et. al., (2021) this research explores the relationships between ecosystem services in coffee systems in Colombia, examining their trajectories after agroforestry transition and the underlying variables. Research was conducted to examine the chronology of agro-forestry coffee plantations from 1-40 years after installing shade trees. This study found that agroforestry restores ecosystem services [18]. Bebber et. al., (2016) Using climatic reanalyzes, researchers examined the idea that climate change caused Coffee Leaf Rust in Colombia from 2008 to 2011. Experiments have shown that germination and infection are Weibull functions with different temperature optimums [19]. Berihun, G et. al., (2022) this review discusses Ethiopian CLR disease. The paper studies climate change and CLR out-breaks and suggests disease control strategies [20]. Fanelli Carvalho et. al., (2020) the goal was to evaluate Arabica coffee genome selection by accounting for biennially impacts on yield prediction accuracy. Low (2005 and 2007) and high (2006 and 2008) yield years assessed the GBS high-density genotyped population ($n = 586$). Genotypic selection outperformed prediction methods due to shorter breeding cycles [21]. Santana et. al., (2022) systematically studying soil variability and plant effects improves the field. Scientists at the PC provide improved coffee manufacturing management and security. Designed to reduce pesticide use and soil nutrients for sustainable coffee production [22]. Tadesse et. al., (2021) researchers have sited that main meteorological factors affecting CLR are temperature, moisture, and wind. Understanding climatic and meteorological links to CLR outbreaks may help farmers anticipate and control the disease, reducing crop losses. Future study aims to identify realistic climate change adaption solutions for smallholder growers [23]. Avelino et. al., (2015) to better understand the 2008–2013 outbreaks, experts examined Mesoamerica's most severe epidemics

from 1987 to 2003, particularly in Central America and Colombia. Following these recent outbreaks and the projected climate change that would impact Mesoamerica in the near future, they have proposed various strategies for enhanced disease management [24]. Tadesse et. al., (2020) This study examined coffee production constraints and possibilities in significant coffee-growing districts (Wereda) in four South Nation Nationalities and Peoples Region zones (Sidama, Gedeo, Gamo Goffa, and Wolayta). They interviewed 161 houses for qualitative and quantitative data and analyses it using SPSS and descriptive statistics. Diseases, insect pests, weeds, and vertebrates are the biggest biotic factors. Drought, frost, fluctuating rainfall, high humidity, temperature, low moisture, hail, storm, wind, and soil fertility may reduce coffee yield by 70% [25]. To estimate crop yields at the Coffee Research Station (CCRI) in Chikamgaluru, Karnataka, India, we analyzed several machine learning algorithms and included biotic and abiotic data. Many experts' claim data mining and AI are ready to construct effective forecasting models. Machine learning methods were explored to forecast coffee output from biotic and abiotic characteristics. We extensively discussed data sets and techniques here. Detailed findings explanation provided. Detailed analysis of five probabilistic regression models for predicting coffee production using various parameter fine tuning strategies.

2. Dataset And Methods

2.1. Study Area and Dataset

Located in the Balehonnur district of the Karnataka State of India, the Central Coffee Research Institute (CCRI) was the source of the data gathered from 2015 to 2022. Seven input characteristics derived from one thousand samples make up the dataset for the model used in this study, which includes both biotic and abiotic components. Table 1 shows the datasets and describes the biotic (Coffee Leaf Rust incidence, or CLR) and abiotic (parameters affecting yield) variables in detail. We performed research at Central Coffee Research Institute (CCRI), Coffee Research Station, Chikkamagaluru District, Karnataka, India, to examine how

meteorological factors affect CLR occurrence. Recorded CLR incidence at fortnightly intervals from Coffee Arabica L cultivar Sln.3 at CCRI farm during 2015-2022. The meteorological observatory at CCRI collected weather data, including maximum and lowest temperatures, relative humidity, and rainfall amounts.

Table 1 Biotic and Abiotic Parameters

Parameters Name	Description
Year	2015-2022
Temperature – Minimum in Degree’s	12.4 to 22.8
Temperature – Maximum in Degree’s	20.6 to 29.08
Coffee Leaf Rust (CLR) Incidence	Two Fort-Night’s in Each month – incidence data

- Day-1 to Day-15: First Fort Night.	were recorded in Percentage. (0% to 79.69%)
- Day-16 to Last day of respective month: Second Fort Night.	
Relative Humidity – in Percentage	32% to 100%
Rainfall in Centimeters	186.1 cm to 329.5 cm
Coffee Type	SLN3 (Arabica)
Yield in Kg/Acre	210 to 400

2.2. Proposed Methodology

To make and predict coffee production according to biotic and abiotic criteria, we provide probabilistic machine learning methods in this section. Figure 1 depicts the recommended technique in square form.

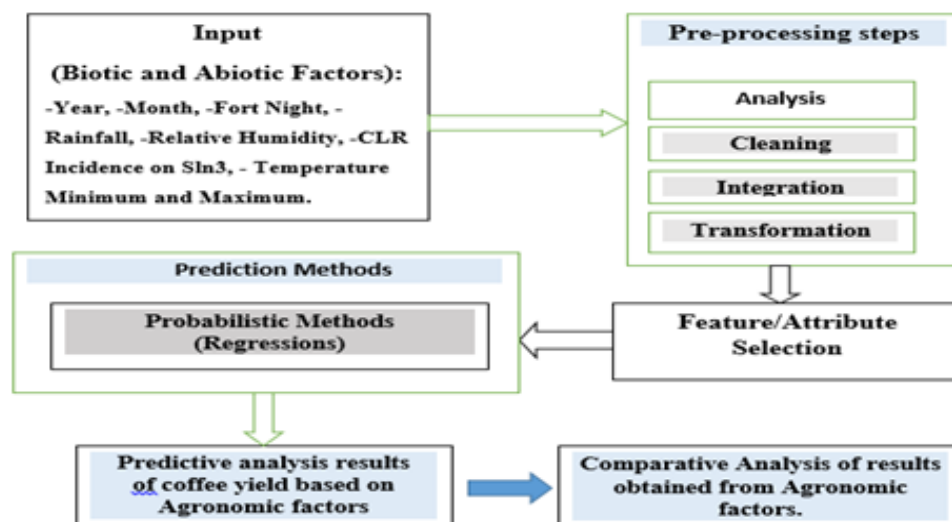


Figure 1 Proposed Methodology Block Diagram for Coffee Yield Prediction System Based on Biotic – CLR Incidence and Abiotic Factors

The pre-processing part of our suggested method just required a few steps. The datasets didn't seem to have many gaps at first look. Mathematical means, medians, and binning algorithms are among the statistical methods used to fill in the gaps in the data. Next, we may use the climate-related component standard critical values that the CCRI has set forth [26]. By applying a boxplot to the biotic and abiotic factors, we were able to eliminate any outliers [27]. By eliminating outliers, we may decrease the initial

1000 data points to 794. Then, we used skewness to check whether the data are regularly distributed [28].

2.3. Methods

We discussed probabilistic machine learning methods in this section that are being investigated for the purpose of coffee yield prediction using biotic (CLR Incidence) and abiotic factors.

2.3.1. Extra Tree Regression

The "Extra Tree Regression" RF model is a novel



method. Unedited judgments or regression trees result from top-down ETR. Random Forest regression bags and bootstraps. Starting with random training datasets, bootstrapping creates decision trees. After ensemble creation, two-step bagging separates decision tree nodes. Bagging begins with random selection of training data groupings. Deciding on the top subgroup and its value completes decision-making [29].

Instructions for using the Extra Tree Regression Algorithm for numerical attributes are provided below [30]:

Step-1: Dividing a node (A)

The learning subset A for the neighboring node that has to be divided is the input.

A node split [x xy] or a zero split is the output.

Return 0 if Stop split (A) equals.

Anyhow, from all non-Consistent (in A) Applicant characteristics, choose

B attributes (c1... cb);

Describe the locations of the B divides e1... eb. ei = choose a random number

Splice (A, xi), v j = 1, and then B;

To ensure that Count (d*, A) = maxi=1... B Count (ei, A), return a split e*.

Step-2: Take a haphazard split. (A,x)

Inputs include an attribute x and x subclass A.

Results: x split

Let xAmax and xAmin denote the maximum and insignificant estimates of m in A, respectively.

Illustrate any boundary ac in accordance with [xAmin, xAmax].

Send the split ([x xy]) back.

Step-3: Reverse split (A)

Enter: x subclass A binary output for x return TRUE if |A| xmin; Return TRUE if A's properties are all consistent.

Return TRUE if the result is in accordance with A; FALSE otherwise.

The steps above describe the Extra-Trees splitting approach for numerical features. It contains two parameters: ymin, the lowest sample size for splitting a node, and B, the number of characteristics randomly picked at each node. It is combined with the (whole) original multiple times. Creating an

ensemble model by learning a sample. The final forecast is produced by combining the tree predictions. In classification issues, use the majority vote, while in regression issues, and use the arithmetic average [30].

2.3.2. Gradient Boosting Regression

In ensemble learning, gradient boosting regression tree approaches employ weak learner regression trees (decision trees) to create reliable forecasting models. Poorly trained models (repressors or classifiers) experience fewer errors. Poorly learned models show a bias toward the training dataset, little volatility, and minimal regularization compared to random guesses [31]. Boosting uses additive models, weak learners, and loss functions. Gradient boosting machines establish weak model constraints using gradients. Iteratively connecting decision trees using an additive model and decreasing the loss function with gradient descent reduces prediction errors by linking base learners.

The gradient boosting tree, also known as the $F_n(x)$ algorithm, is the accumulation of n regression trees.: $F_n(x) = \sum_{i=1}^n f_i(x)$ ----- (i)

Every $f_i(x)$ is a decision tree or regression. The equation estimates the new decision tree $f_{n+1}(x)$ to form the ensemble of trees:

$$\operatorname{argmin}_{\sum_{i=1}^n f_i(x)} [L(y, F_n(x) + f_{n+1}(x))] \text{ ---}$$

(ii) Where the loss function $L(\cdot)$ is differentiable. This study employed 0.2 learning rate and 100 estimator value. When learning rate is smaller, stopping before over fitting is easier.

2.3.3. Random Forest Regression

Random Forest (RF) ensemble-based data mining makes accurate predictions without over fitting. They use model aggregation-based learning [32]. Random forests use bootstrapped learning samples, Binary decision trees, and a random selection of explanatory factors. Validation or "out-of-bag" predictions create up to 2000 random trees using data for a third of samples in the RF technique. Each tree is bootstrapped. Random predictors branch the tree at each node, and the result is the average of all trees. Out-of-bag samples are used by Random Forest to test independent tree data model error. Need no cross-validation [32]. Sequentially, the



Random Forest algorithm:

Step-1: Select and replace N training set data cases randomly. Growing original trees is training.

Step-2: At each node, Random Forest randomly picks m variables from M inputs (or predictors). The best split on these m variables separates the node while keeping m's value as the forest Grows.

Step-3: The Random Forest method maximizes tree size without cutting its structure.

Step-4: If regression is a problem, pooling n trees' predictions gives a mean value for predicting Incoming data.

The out-of-bag error rate estimate may be accurate if enough Random Forest trees duplicate particular data. The out-of-bag error rate estimate may be accurate if enough Random Forest trees duplicate particular data.

2.3.4. Decision Tree Regression

Machine learning and data mining employ classification and regression decision trees. Are the anticipated outcomes the observed or given variable's class. Classification tree analysis. If the predicted value is actual, use regression tree analysis. Regression decision trees include non-leaf nodes for binary attribute tests, branches for test results, and leaf nodes for projected values or labels. The highest tree node is the root. Recursive partitioning or sub-division develops tree branches and predicts using binary questions for each feature value. Classification or regression trees need feature vectors or observed and known values [33]. To

create a tree, hierarchically subdividing the space takes three steps: splitting nodes, identifying terminal nodes, and adding labels or anticipated values to terminal nodes. The majority or weighted votes determine class labels or projected values, with some being more probable or costlier than others [34].

2.3.5. K-Nearest Neighbour Regression

The k-Nearest Neighbor (kNN) approach classifies a new item based on the k nearest training dataset points, allocating it to the class with the most points. Regression calculations use a weighted sum of answers from k neighbors, with weight inversely proportional to the distance (normalized Euclidean) from the input record. The simplest variant uses k = 1. This creates an unstable prediction model with high volatility and data sensitivity. Increasing k decreases variance but may increase bias. The method is sensitive to selecting k properly. There is no need for optimization or training beyond selecting k and the distance measure. Additionally, the technique utilizes local information to create nonlinear, adaptive decision limits. However, the approach is popular for its simplicity and the specified features [35].

2.4. Model Performance Evaluation Metrics

This research evaluated five stochastic regression models for coffee yield prediction by comparing measured yield data with test phase yield data. Analysis of R2, MAE, MSE, and root MSE. Table 2 basic performance metrics predict production:

Table 2 Standard Performance Metrics

S.No	Performance Metrics	Formula
1	R-Squared (R2)	$R^2 \text{ Squared} = 1 - \frac{SSr}{SSm}$ SSr – Squared regression line sum error SSm - Squared mean line sum error.
2	Mean Absolute Error (MAE)	$MAE = \frac{1}{N} \sum_{i=1}^{i=N} (y_i - \hat{y}_i)$ There are N anticipated values.. The ith data's real true value is represented by yi. \hat{y}_i is the i-th data's anticipated value.

3	MSE, or Mean Squared Error	$MSE = 1/N \sum (y_i - \hat{y}_i)^2$
4	RMSE, or Root Mean Square Error	$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$

3. Results And Discussion

3.1. Results

Scatterplots (fig. 2) shown the degree of agreement between actual yield and expected yield data during testing to condense all characteristics to a single scale without modifying the probabilistic models' range of values. This research analyzed probabilistic model outcomes using normal parameter, ordinary scalar, and principal component analysis fine tuning functions. Additionally, we have developed our proposed models in three groups like: Group-1: All Input Parameters vs. Yield Using Normal Parameter Fine-tuning function, Group-2: All Input Parameters vs. Yield Using Ordinary Scalar Fine-tuning function and Group-3: All Input parameters

vs. Yield using Principal Component Analysis Functions respectively. Figure-2 shows performance metrics using scatterplots, including R-Square (R^2), RMS Error (RMSE), actual and predicted yield, and model error rate.

3.1.1. Extra Tree Regression

Extremely Randomized Trees (Extra Trees) is a regression like Random Forest. Nodes are built using all preparation set data. Mold the root or any other node before checking a sqrt subset of randomly produced characteristics for the appropriate split. Feature divides randomly. An Extra tree regression model was developed utilizing 100 DTs from various relevant categories from the test and train datasets show in Table 3.

Table 3 Extra Tree Regression Model Was Used to Forecast Coffee Yield with Variable Biotic and Abiotic Factors for Three Groups. We Looked at R2, Mae, Mses, And Rmse Ratios for Different Split Ratios in The Testing Phase. The Best Bargains Are Boldfaced

Group-1 Extra Tree Model: All Parameters vs. Yield Using Normal Fine-tuning function.				
Quantity Shared	R-Square kg/ha	Mean Absolute Error kg per ha	Mean Square Error kg per ha	Root Mean Square Error kg per ha
90:10	0.98	5.58	68.91	8.30
80:20	0.98	5.87	63.58	7.97
70:30	0.98	6.01	63.42	7.96
60:40	0.97	6.81	89.23	9.45
50:50	0.94	8.85	148.04	12.17
Group-2 Extra Tree Model: All Parameters vs. Yield Using Ordinary Scalar Fine-tuning function.				
90:10	0.98	5.38	70.56	8.40
80:20	0.98	6.01	64.77	8.05
70:30	0.98	6.25	64.81	8.05
60:40	0.97	6.33	79.93	8.94
50:50	0.95	8.60	145.34	12.06
Group-3 Extra Tree Model: All Parameters vs. Yield Using Principal Component Analysis function.				
90:10	0.53	30.03	1596.61	39.96
80:20	0.54	29.32	1498.90	38.72
70:30	0.50	28.40	1519.10	38.98
60:40	0.50	25.80	1333.85	36.52
50:50	0.49	27.61	1380.61	37.16



Considering all input parameters vs. yield, Group-1 fine tuning function model yields the most promising results. The additional tree regression model had the best coefficient of determination (R-square = 0.98 kg per ha and Root Mean Square Error = 7.96 kg per ha) (figure-2). Comparing measured yield to expected yield for 70:30 split ratio showed this. Performance indicators showed the extra tree regression model performed less error rate for groups 2 and 3. Gradient boosting regression, random forest regression, KNN regression, and Decision Tree regression models fail to extract

predictive features from multi-parameter data, as does the Extra Tree regression model using the same inputs and group-1 parameters. Table-3 displays R², MAE, MSE, and RMSE performance for several splitting models during testing.

3.1.2. Gradient Boosting Regression

Boosting uses 100 weak learners and a random forest base_estimator. Every stage adds weak learners to make up for existing weak learners. Gradients show the merged model's weaknesses is show in Table 4.

Table 4 Gradient Boosting Regression Model Was Used to Forecast Coffee Yield with Variable Biotic and Abiotic Factors for Three Groups. We Looked at R², Mae, Mse, And Rmse Ratios for Different Split Ratios in The Testing Phase. The Best Bargains Are Boldfaced

Group-1 Gradient Boosting Model: All Parameters vs. Yield Using Normal Fine-tuning function.				
Quantity Shared	R-Square kg/ha	Mean Absolute Error kg per ha	Mean Square Error kg per ha	Root Mean Square Error kg per ha
90:10	0.98	6.29	79.68	8.93
80:20	0.97	6.68	91.94	9.59
70:30	0.96	8.20	121.90	11.04
60:40	0.97	6.80	87.61	9.36
50:50	0.93	9.35	186.44	13.65
Group-2 Gradient Boosting Model: All Parameters vs. Yield Using Ordinary Scalar Fine-tuning function.				
90:10	0.98	6.07	77.77	8.82
80:20	0.97	6.58	91.26	9.55
70:30	0.96	8.14	120.18	10.96
60:40	0.97	6.86	89.95	9.48
50:50	0.93	9.21	183.79	13.56
Group-3 Gradient Boosting Model: All Parameters vs. Yield Using Principal Component Analysis function.				
90:10	0.17	36.13	2821.00	53.11
80:20	0.72	21.01	914.35	30.24
70:30	0.48	25.95	1577.82	39.72
60:40	0.49	24.25	1379.67	37.14
50:50	0.09	34.01	2434.70	49.34



Considering all input parameters vs. yield, Group-2 fine tuning function model yields the most promising results. The gradient boosting regression model had the best coefficient of determination (R-square = 0.96 kg per ha and Root Mean Square Error = 10.96 kg per ha) (figure-2). Comparing measured yield to expected yield for 70:30 split ratio showed this. Performance indicators showed the gradient boosting regression model performed less error rate for groups 1 and 3. Extra Tree regression, random forest regression, KNN regression, and Decision Tree regression models fail to extract predictive

features from multi-parameter data, as does the gradient boosting regression model using the same inputs and group-2 parameters. Table-4 displays R², MAE, MSE, and RMSE performance for several splitting models during testing.

3.1.3. Random Forest Regression

We tested the random forest model with tens of thousands of samples to determine whether it improved reliability over regression models. Table 5 indicates significant parameter group differences in splitting.

Table 5 Random Forest Regression Model Was Used to Forecast Coffee Yield with Variable Biotic and Abiotic Factors for Three Groups. We Looked at R², Mae, Mse, And Rmse Ratios for Different Split Ratios in The Testing Phase. The Best Bargains Are Boldfaced

Group-1 Random Forest Regression Model: All Parameters vs. Yield Using Normal Fine-tuning function.				
Quantity Shared	R-Square kg/ha	Mean Absolute Error kg per ha	Mean Square Error kg per ha	Root Mean Square Error kg per ha
90:10	0.97	5.79	107.46	10.37
80:20	0.97	6.65	104.90	10.24
70:30	0.96	7.93	114.65	10.71
60:40	0.96	6.69	99.02	9.95
50:50	0.92	9.24	206.53	14.37
Group-2 Random Forest Regression Model: All Parameters vs. Yield Using Ordinary Scalar Fine-tuning function.				
90:10	0.97	5.99	107.68	10.38
80:20	0.97	6.84	109.44	10.46
70:30	0.96	7.92	115.77	10.76
60:40	0.96	6.85	97.62	9.88
50:50	0.92	9.38	208.79	14.45
Group-3 Random Forest Regression Model: All Parameters vs. Yield Using Principal Component Analysis function.				
90:10	0.55	27.90	1541.31	39.26
80:20	0.53	28.28	1522.03	39.01
70:30	0.57	26.12	1313.37	36.24
60:40	0.57	24.62	1162.09	34.09
50:50	0.47	28.02	1433.07	37.86



Considering all input parameters vs. yield, Group-1 fine tuning function model yields the most promising results. The random forest regression model had the best coefficient of determination (R-square = 0.96 kg per ha and Root Mean Square Error = 10.71 kg per ha) (Figure-2). Comparing measured yield to expected yield for 70:30 split ratio showed this. Performance indicators showed the random forest regression model performed less error rate for groups 2 and 3. Gradient boosting regression, extra tree regression, KNN regression, and Decision Tree regression models fail to extract predictive features

from multi-parameter data, as does the random forest regression model using the same inputs and group-1 parameters. Table 5 displays R², MAE, MSE, and RMSE performance for several splitting models during testing.

3.1.4. Decision Tree Regression

The decision tree regression model took into account 100 DTs from different subgroups by combining the test and train datasets. In splitting, Table 6 shows substantial parameter group disparities.

Table 6 Decision Tree Regression Model Was Used to Forecast Coffee Yield with Variable Biotic and Abiotic Factors for Three Groups. We Looked at R², Mae, Mse, And Rmse Ratios for Different Split Ratios in The Testing Phase. The Best Bargains Are Boldfaced

Group-1 Decision Tree Regression Model: All Parameters vs. Yield Using Normal Fine-tuning function.				
Quantity Shared	R-Square kg/ha	Mean Absolute Error kg per ha	Mean Square Error kg per ha	Root Mean Square Error kg per ha
90:10	0.97	4.55	90.65	9.52
80:20	0.93	8.59	233.77	15.29
70:30	0.92	7.43	253.91	15.93
60:40	0.94	7.34	162.01	12.73
50:50	0.93	8.58	191.50	13.84
Group-2 Decision Tree Regression Model: All Parameters vs. Yield Using Ordinary Scalar Fine-tuning function.				
90:10	0.97	4.55	90.65	9.52
80:20	0.93	8.59	233.77	15.29
70:30	0.92	7.43	253.91	15.93
60:40	0.94	7.34	162.01	12.73
50:50	0.93	8.58	191.50	13.84
Group-3 Decision Tree Regression Model: All Parameters vs. Yield Using Principal Component Analysis function.				
90:10	0.28	26.10	2467.20	49.67
80:20	0.33	23.62	2165.77	46.54
70:30	0.25	26.79	2276.69	47.71
60:40	0.17	24.79	2222.32	47.14
50:50	-0.03	30.06	2777.50	52.70



Considering all input parameters vs. yield, Group-1 and 2 fine tuning function model yields the most promising results. The Decision tree regression model had the best coefficient of determination (R-square = 0.92 kg per ha and Root Mean Square Error = 15.93 kg per ha) (figure-2). Comparing measured yield to expected yield for 70:30 split ratio showed this. Performance indicators showed the decision tree regression model performed less error rate for group-3. Gradient boosting regression, extra tree regression, KNN regression, and random forest regression models fail to extract predictive features

from multi-parameter data, as does the random forest regression model using the same inputs and group-1 and 2 parameters. Table-6 displays R², MAE, MSE, and RMSE performance for several splitting models during testing.

3.1.5. K-Nearest Neighbour Regression

With a weight that is inversely proportionate to the distance (normalized Euclidean) from the input record, regression computations employ a weighted sum of replies from k neighbours. The most basic version used here is k = 1.

Table 7 K-Nearest Neighbour Regression Model Was Used to Forecast Coffee Yield with Variable Biotic and Abiotic Factors for Three Groups. We Looked at R², Mae, Mse, And Rmse Ratios for Different Split Ratios in The Testing Phase. The Best Bargains Are Boldfaced

Group-1 K-Nearest Neighbour Regression Model: All Parameters vs. Yield Using Normal Fine-tuning function.				
Quantity Shared	R-Square kg/ha	Mean Absolute Error kg per ha	Mean Square Error kg per ha	Root Mean Square Error kg per ha
90:10	0.52	26.40	1636.70	40.46
80:20	0.58	24.05	1370.51	37.02
70:30	0.46	28.05	1647.94	40.59
60:40	0.51	24.23	1329.98	36.47
50:50	0.23	30.13	2060.15	45.39
Group-2 K-Nearest Neighbour Regression Model: All Parameters vs. Yield Using Ordinary Scalar Fine-tuning function.				
90:10	0.05	41.25	3250.15	57.01
80:20	0.23	37.63	2495.51	49.96
70:30	0.06	41.31	2853.76	53.42
60:40	0.10	36.86	2421.78	49.21
50:50	-0.07	41.35	2864.95	53.53
Group-3 K-Nearest Neighbour Regression Model: All Parameters vs. Yield Using Principal Component Analysis function.				
90:10	0.28	34.38	2449.49	49.49
80:20	0.45	28.19	1799.20	42.42
70:30	0.35	31.96	1966.56	44.35
60:40	0.41	27.58	1582.38	39.78
50:50	0.17	31.95	2241.58	47.35

Considering all input parameters vs. yield, Group-1 fine tuning function model yields the most promising results. The K-Nearest Neighbour regression model had the best coefficient of determination ($R^2 = 0.46$ kg per ha and Root Mean Square Error = 40.59 kg per ha) (Figure-2). Comparing measured yield to expected yield for 70:30 split ratio showed this. Performance indicators showed the K-Nearest Neighbour

regression model performed less error rate for groups-2 and 3. Gradient boosting regression, extra tree regression, decision tree regression, and random forest regression models fail to extract predictive features from multi-parameter data, as does the random forest regression model using the same inputs and group-1 parameters. Table-7 displays R^2 , MAE, MSE, and RMSE performance for several splitting models during testing.

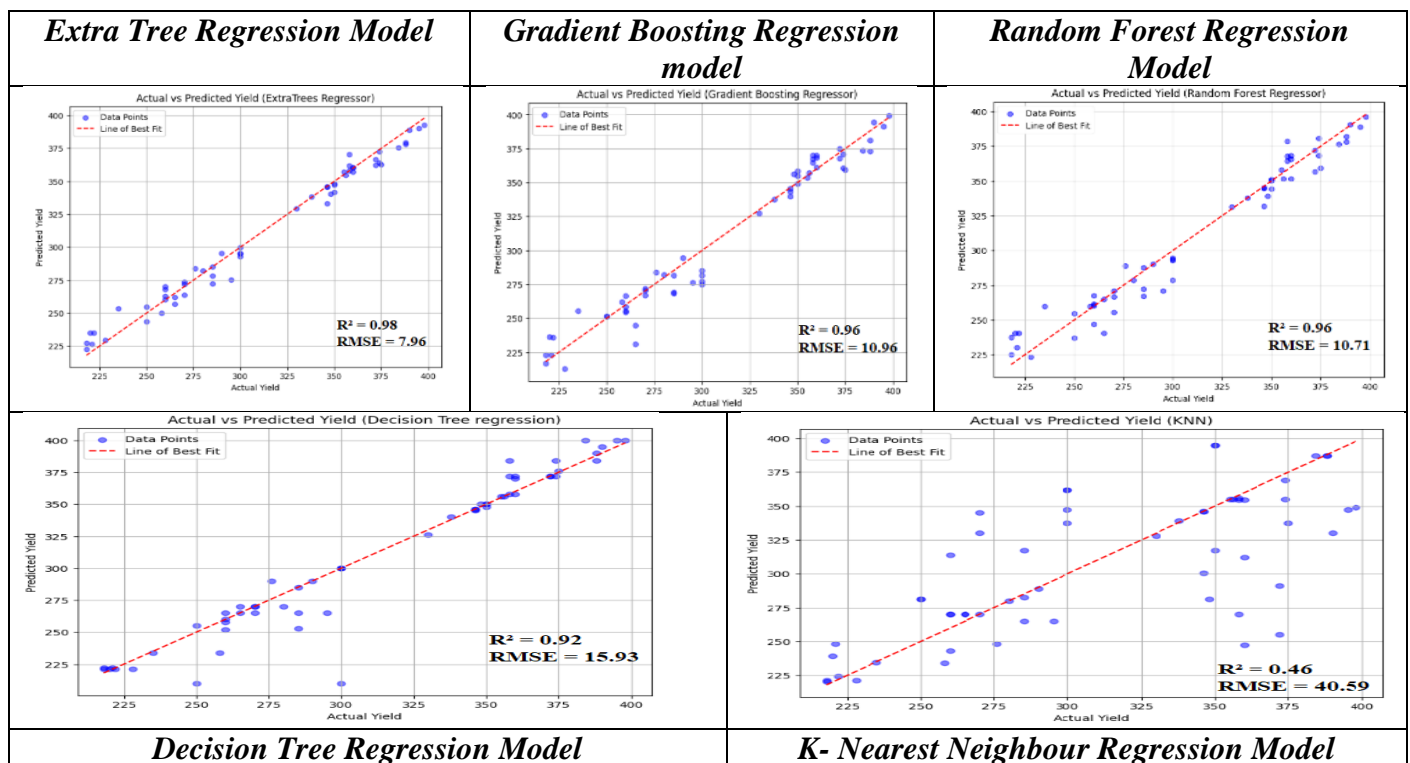


Figure 2 Group 1 & 2 Served as Input Parameters for The Proposed Models Shown in Above Scatterplots Based on Probabilistic Models Based on the Extra Tree, Gradient Boosting, Random Forest, Decision Tree And K-Nearest Neighbor Regressions. Above Are Scatterplots Depicting the Highest Expected and Actual Coffee Yields During the Experimentation Period (70:30 Splits)

3.2. Discussion

There are a lot of factors that affect coffee crop productivity, including the amount of rainfall, temperature, daylight hours, vapour, dew points, and relative humidity, as well as the area produced in. Changes in these factors affect the monthly coffee harvest yield. Prediction is necessary for accurate production estimation and meeting customer demand. In the past, neither farmers nor entrepreneurs had a good grasp of the monthly

coffee crop production potential nor the factors that may impact it. This research used five different probabilistic models to assess and predict the coffee crop output based on data on biotic and abiotic parameters, such as coffee leaf rust incidence, from 2015 to 2022. Using data intelligence and statistical analytics, we may try to predict which abiotic and biotic components, when combined, would lead to the greatest yield in crop development simulation models, which is no easy feat. Despite various



challenges, Indian coffee producers must increase crop nutrition and resource use efficiency to maintain output. Our technique identified yield-estimating biotic and abiotic variables. Even if non-selected abiotic factors continue to impact coffee growth, our findings may guide future research to enhance decision-making on smallholder coffee farms. We did not take into account any other variables that can affect the stated yields in coffee production. Our research found that biotic and abiotic data-based models may predict coffee output and its key components. Other studies indicated that soil parameters were the biggest yield confounders. So far, Tables 3, 4, 5, 6, and 7 under the model generated for group 1, 2, and 3 parameters utilizing fine tuning functions are correct. (Fig. 2) R-squared compares model fit to horizontal straight Line. Throughout our analysis, we found no instances where the selected model fit the data worse than a horizontal line. Although the extra tree and gradient boosting regression models improved, they could do better. Biophysical modelling of coffee production at larger regional scales (such as Karnataka's Coorg area) requires additional empirical correlations using data from several coffee-producing provinces. Comparing and comparing smallholder and big farm practices is important. We performed research at Central Coffee Research Institute (CCRI), Coffee Research Station, Chikkamagaluru District, Karnataka, India, to examine how meteorological factors affect CLR occurrence. Recorded CLR incidence at fortnightly intervals from Coffee Arabica L cultivar Sln.3 at CCRI farm during 2015-2022. The meteorological observatory at CCRI collected weather data, including maximum and lowest temperatures, relative humidity, and rainfall amounts. The study found that rainfall distribution over the period September to November differed from monthly rainfall amounts. Coffee leaf rust disease depends on time and temperature. This research found that linkages matter more than individual meanings for biotic and abiotic components. Now that land management, soil, and environmental conditions need various chemical, biological, and physical indicators, studies must

concentrate on more than one or two environmental parameters. We divided each parameter's relevance by its greatest value to get its index weight. Therefore, the model with the highest R² and lowest RMSE is best.

Conclusion

Using biotic – CLR incidence data and abiotic variables measured at the Central Coffee Research Institute Station in Karnataka's Chikkamagaluru coffee-growing area, this study compares actual and expected coffee yield. The coffee research station in the Chikkamagaluru area aimed to improve coffee crop yield, therefore they evaluated the usefulness of probabilistic models, a data-driven approach for examining biotic and abiotic variable data for predictive traits. In order to create the probabilistic models under consideration, a unique machine learning technique for handling complicated and ill-defined situations was used. The goal variable was the coffee yield (Y), and the seven biotic and abiotic characteristics were divided into three groups using different parameter fine tuning functions. The predictor variables included year, temperatures (both minimum and maximum), rainfall, relative humidity (both minimum and maximum), and CLR incidence, month). Our results showed that compared to Gradient Boost, Decision Tree, KNN, Random Forest, and Gradient Boost, Extra Tree was the most successful regression model. Using a variety of variables, regression models are able to more accurately predict coffee yields by extracting characteristics from interactions between biotic and abiotic factors. This study shown the possible benefits of combining biophysical crop models with AI algorithms in precision agriculture decision-support systems by using a set of properly screened data for biotic and abiotic traits to intentionally increase productivity in smallholder farms. Analysing these machine learning methods may help us build better models for analysis and predictions.

Acknowledgements

The Authors Would Like to Express Their Gratitude to The Director of Research at The Central Coffee Research Institute (CCRI), As Well As the Coffee



Board, For Providing and Encouraging Us to Carry Out the Experiment. Additionally, We Would Like to Express Their Gratitude to The Division of Agricultural Chemistry and The Division of Entomology at The CCRI For Supplying the Meteorological and CLR Incidence Data.

References

- [1]. Santhosh C S., & Umesh K K. (2023) An abiotic factors-based compendium probabilistic forthcoming methodology for prediction of coffee crop yield. International Journal of Science, Mathematics and Technology Learning, ISSN: 2327-7971 (print) ISSN: 2327-915X (online), Volume: 31 No.22023.<https://doi.org/10.5281/zenodo.8388471>.
- [2]. Santhosh, C. S., & Umesh, K. K., (2022) A Compendium Probabilistic Prospective for Predicting Coffee Crop Yield Based on Agronomical Factors. In 2022 Fourth International Conference on Emerging Research in Electronics, Computer Science and Technology (ICERECT), (pp. 1-8). IEEE.
- [3]. Santhosh, C. S., & Umesh, K. K. (2022) An ensemble approach for coffee crop yield prediction based on agronomic factors. ASEAN Engineering Journal (SCOPUS), 13(3), 29-38. <https://doi.org/10.11113/aej.v13.18846>.
- [4]. Jaramillo, J., Chabi-Olaye, A., Kamonjo, C., Jaramillo, A., Vega, F. E., Poehling, H. M., & Borgemeister, C., (2009) Thermal tolerance of the coffee berry borer *Hypothenemus hampei*: predictions of climate change impact on a tropical insect pest. PLoS one, 4(8), e6487,
- [5]. Alves, M. D. C., de Carvalho, L. G., Pozza, E. A., Sanches, L., & Maia, J. D. S., (2011). Ecological zoning of soybean rust, coffee rust and banana black sigatoka based on Brazilian climate changes. Procedia Environmental Sciences, 6, 35-49.
- [6]. Jaramillo, J., Muchugu, E., Vega, F. E., Davis, A., Borgemeister, C., & Chabi-Olaye, A., (2011). Some like it hot: the influence and implications of climate change on coffee berry borer (*Hypothenemus hampei*) and coffee production in East Africa. PLoS one, 6(9), e24528.
- [7]. Pérez-Ariza, C. B., Nicholson, A. E., & Flores, M. J., (2012). Prediction of coffee rust disease using Bayesian networks. In Proceedings of the Sixth European Workshop on Probabilistic Graphical Models, (Vol. 6, pp. 259-266).
- [8]. Kutwayo, D., Chemura, A., Kusena, W., Chidoko, P., & Mahoya, C., (2013). The impact of climate change on the potential distribution of agricultural pests: the case of the coffee white stem borer (*Monochamus leuconotus* P.) in Zimbabwe. PLoS One, 8(8), e73432.
- [9]. Classen, A., Peters, M. K., Ferger, S. W., Helbig-Bonitz, M., Schmack, J. M., Maassen, G., & Steffan-Dewenter, I. (2014). Complementary ecosystem services provided by pest predators and pollinators increase quantity and quality of coffee yields. Proceedings of the Royal Society B: Biological Sciences, 281(1779), 20133148.
- [10]. Wang, N., Jassogne, L., van Asten, P. J., Mukasa, D., Wanyama, I., Kagezi, G., & Giller, K. E., (2015) Evaluating coffee yield gaps and important biotic, abiotic, and management factors limiting coffee production in Uganda. European Journal of Agronomy, 63, 1-11.
- [11]. Corrales, D. C., Corrales, J. C., & Figueroa-Casas, A., (2015) towards detecting crop diseases and pest by supervised learning. Ingeniería y Universidad, 19(1), 207-228.
- [12]. Hameed, A., Hussain, S. A., & Suleria, H. A. R. (2015) "Coffee Bean-Related" agroecological factors affecting the coffee. Co-evolution of secondary metabolites, 641-705.
- [13]. Sudha, M., Machenahalli, S., Giri, M. S., Ranjini, A. P., & Daivasikamani, S. (2020).



- Influence of abiotic factors on coffee leaf rust disease caused by the fungus *Hemileia vastatrix* Berk. & Br. under changing climate. *Journal of Agro meteorology*, 22(3), 367-371.
- [14]. Yáñez-López, R., Torres-Pacheco, I., Guevara-González, R. G., Hernández-Zul, M. I., Quijano-Carranza, J. A., & Rico-García, E. (2012). The effect of climate change on plant diseases. *African Journal of Biotechnology*, 11(10), 2417-2428.
- [15]. Suresh, N., Santa Ram, A., & Shivanna, M. B. (2012). Coffee leaf rust (CLR) and disease triangle: A case study. *International Journal of Food, Agriculture and Veterinary Sciences*, 2(2), 50-55.
- [16]. Cerda, R., Avelino, J., Gary, C., Tixier, P., Lechevallier, E., & Allinne, C. (2017). Primary and secondary yield losses caused by pests and diseases: Assessment and modelling in coffee. *PloS one*, 12(1), e0169133.
- [17]. Abreu Júnior, C. A. M. D., Martins, G. D., Xavier, L. C. M., Vieira, B. S., Gallis, R. B. D. A., Fraga Junior, E. F., ... & Lima, J. V. D. N. (2022) Estimating Coffee Plant Yield Based on Multispectral Images and Machine Learning Models. *Agronomy*, 12(12), 3195.
- [18]. De Leijster, V., Santos, M. J., Wassen, M. W., García, J. C., Fernandez, I. L., Verkuil, L., ... & Verweij, P. A. (2021). Ecosystem services trajectories in coffee agroforestry in Colombia over 40 years. *Ecosystem Services*, 48, 101246.
- [19]. Bebbler, D. P., Castillo, Á. D., & Gurr, S. J. (2016) Modelling coffee leaf rust risk in Colombia with climate reanalysis data. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 371(1709), 20150458.
- [20]. Berihun, G., & Alemu, K. (2022). Status of coffee leaf rust (*Hemileia vastatrix*) and its management in Ethiopia: a review. *Archives of Phytopathology and Plant Protection*, 55(20), 2283-2300.
- [21]. Fanelli Carvalho, H., Galli, G., Ventorim Ferrão, L. F., Vieira Almeida Nonato, J., Padilha, L., Perez Maluf, M., ... & Fritsche-Neto, R. (2020). The effect of bienniality on genomic prediction of yield in Arabica coffee. *Euphytica*, 216(6), 101.
- [22]. Santana, L. S., Ferraz, G. A. E. S., dos Santos, S. A., & Dias, J. E. L. (2022). Precision coffee growing: a review.
- [23]. Tadesse, Y., Amare, D., & Kesho, A., (2021). Coffee leaf rust disease and climate change. *World Journal of Agricultural Science*, 17(5), 418-429.
- [24]. Avelino, J., Cristancho, M., Georgiou, S., Imbach, P., Aguilar, L., Bornemann, G., & Morales, C. (2015). The coffee rust crises in Colombia and Central America (2008–2013): impacts, plausible causes and proposed solutions. *Food security*, 7, 303-321.
- [25]. Tadesse, T., Tesfaye, B., & Abera, G. (2020) Coffee production constraints and opportunities at major growing districts of southern Ethiopia. *Cogent Food & Agriculture*, 6(1), 1741982.
- [26]. Coffee guide book-a manual of coffee Cultivation, Central coffee research Institute (Ministry of commerce and Industry, Govt. of India) (2014).
- [27]. Schwertman, N. C., Owens, M. A., and Adnan, R., (2014) A simple more general boxplot Method for identifying outliers. *Computational statistics and data analysis*, 47(1), 165-174.
- [28]. Joanes, D. N., and Gill, C. A., (1998) Comparing measures of sample skewness and kurtosis. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 47(1), 183-189.
- [29]. Hamed, M. M., AlOmar, M. K., Khaleel, F., & Al-Ansari, N., (2021) An extra tree regression model for discharge coefficient prediction: novel, practical applications in the hydraulic sector and future research directions. *Mathematical problems in*



- engineering, 1-19.
- [30]. Geurts, P., Ernst, D., & Wehenkel, L., (2006).
- [31]. Extremely randomized trees. *Machine learning*, 63, 3-42.
- [32]. Singh, U., Rizwan, M., Alaraj, M., & Alsaidan, I., (2021) A machine learning-based gradient boosting regression approach for wind power production forecasting: A step towards smart grid environments. *Energies*, 14(16), 5196.
- [33]. Louis Kouadioa, Ravinesh C. Deob, Vivekananda Byrareddy, Jan F. Adamowskic, Shahbaz Mushtaq, Van Phuong Nguyend (2018) “Artificial intelligence approach for the prediction of Robusta Coffee yield using soil fertility properties “*Computers and Electronics in Agriculture*, 155 324–338.
- [34]. Chowdhury, D., Sarkar, M., Haider, M. Z., & Alam, T. (2018) Zone wise hourly load prediction using regression decision tree model. In 2018 International Conference on Innovation in Engineering and Technology (ICIET), (pp. 1-6). IEEE.
- [35]. S. B. Gelfand, C. S. Ravishankar and E. J. Delp, (1991) “An Iterative Growing and Pruning Algorithm for Classification Tree Design”, *IEEE Trans. On Pattern Analysis AND Machine Intelligence*, vol. 13, no. 2, pp. 163-174, Feb.
- [36]. Gonzalez-Sanchez, A., Frausto-Solis, J., & Ojeda-Bustamante, W. (2014) Predictive ability of machine learning methods for massive crop yield prediction.
- [37]. Sudhakar, S.B., Kiran kumar, K.C., Daivasikamani, S., Hanumantha, B.T., Prakash, N.S. and Jayarama.. (2014) Diseases of coffee. In: *Diseases of Plantation crops* (Eds.) Chowdappa, P., Sharma, P., Anandaraj, M. and Khetrapal, R. K., Indian Psychopathological Society, New Delhi, Pp.55-109.