*Article*

# Comparing Regression and Classification Models to Estimate Leaf Spot Disease in Peanut (*Arachis hypogaea* L.) for Implementation in Breeding Selection

Ivan Chapu [1], Abhilash Chandel [2], Emmanuel Kofi Sie [3], David Kalule Okello [4], Richard Oteng-Frimpong [3], Robert Cyrus Ongom Okello [1], David Hoisington [5] and Maria Balota [6],*

[1] College of Agricultural and Environmental Sciences, Makerere University, Kampala P.O. Box 7062, Uganda; chapuivan@gmail.com (I.C.); cyrusongom@gmail.com (R.C.O.O.)
[2] Department of Biological Systems Engineering, Virginia Tech Tidewater Agricultural Research and Extension Center, Suffolk, VA 23437, USA; abhilashchandel@vt.edu
[3] CSIR-Savanna Agricultural Research Institute, Nyankpala P.O. Box TL 52, Ghana; sieemmanuel@gmail.com (E.K.S.); kotengfrimpong@gmail.com (R.O.-F.)
[4] National Semi-Arid Resources Research Institute (NaSARRI), Soroti P.O. Box 56, Uganda; kod143@gmail.com
[5] Peanut Innovation Lab, University of Georgia, Athens, GA 30601, USA; davehois@uga.edu
[6] School of Plant and Environmental Sciences, Virginia Tech Tidewater Agricultural Research and Extension Center, Suffolk, VA 23437, USA
* Correspondence: mbalota@vt.edu

**Abstract:** Late leaf spot (LLS) is an important disease of peanut, causing global yield losses. Developing resistant varieties through breeding is crucial for yield stability, especially for smallholder farmers. However, traditional phenotyping methods used for resistance selection are laborious and subjective. Remote sensing offers an accurate, objective, and efficient alternative for phenotyping for resistance. The objectives of this study were to compare between regression and classification for breeding, and to identify the best models and indices to be used for selection. We evaluated 223 genotypes in three environments: Serere in 2020, and Nakabango and Nyankpala in 2021. Phenotypic data were collected using visual scores and two handheld sensors: a red–green–blue (RGB) camera and GreenSeeker. RGB indices derived from the images, along with the normalized difference vegetation index (NDVI), were used to model LLS resistance using statistical and machine learning methods. Both regression and classification methods were also evaluated for selection. Random Forest (RF), the artificial neural network (ANN), and k-nearest neighbors (KNNs) were the top-performing algorithms for both regression and classification. The ANN ($R^2$: 0.81, RMSE: 22%) was the best regression algorithm, while the RF was the best classification algorithm for both binary (90%) and multiclass (78% and 73% accuracy) classification. The classification accuracy of the models decreased with the increase in classification classes. NDVI, crop senescence index (CSI), hue, and greenness index were strongly associated with LLS and useful for selection. Our study demonstrates that the integration of remote sensing and machine learning can enhance selection for LLS-resistant genotypes, aiding plant breeders in managing large populations effectively.

**Keywords:** peanut breeding; late leaf spot; machine learning; remote sensing; genotype resistance classification

## 1. Introduction

Peanut (*Arachis hypogaea* L.), also known as groundnut, is recognized as a major source of vegetable oil, protein, and income, thereby improving food security and livelihoods [1]. It is grown in over 100 countries, on over 32 million hectares worldwide, and has an average annual production of 53.4 MT [2]. In Africa, over 18 million hectares is under peanut production, and it has become an important source of income for the smallholder farmers who are the most abundant on the continent [1,3]. Uganda and Ghana are among the most

important producers in Eastern and Western Africa, respectively. In Uganda, peanut is one of the most important legumes, only second to common bean (*Phaseolus vulgaris* L.), while in Ghana, it is the most important grain legume in terms of area under cultivation, production, and consumption [4]. The average productivity in both countries is lower compared to that in most developed countries. For example, the average productivity of peanut in the USA was 4550 Kg/ha in 2022, while that in Uganda and Ghana was 2014 Kg/ha [2]. The low productivity in these countries is attributed to low agricultural inputs, eroded soils, abiotic stress, and the high pressure of biotic stress [5].

Late leaf spot [LLS; caused by Northopassalora personata (Berk. and Curt)] is one of the most damaging foliar diseases of peanut in Uganda and Ghana alike. It can cause up to 50% yield losses when it occurs alone, 70% when it occurs together with early leaf spot [6], and up to 100% yield losses when it occurs together with groundnut rosette disease (GRD). The yield losses are dependent on the environmental conditions, cultivar, and disease severity [7]. The regular application of fungicides is effective in controlling LLS [8,9], and it is fairly economical on large farms [10]. However, fungicide use is not widely adopted by smallholder farmers in Africa, because of its high cost. Furthermore, most farmers are not aware of LLS symptoms such as the yellowing of leaves and defoliation and often mistake it for a sign of maturity [11,12]. Therefore, the development of LLS-resistant cultivars is viewed as the most cost-effective and environmentally friendly method of controlling LLS and maintaining stable yields among smallholder farmers [10].

Selection for disease-resistant and high-yielding genotypes during breeding involves the phenotyping of large populations across several breeding locations [13]. Phenotyping and selection for resistance in many breeding programs are mainly based on destructive sampling [14] and visual ratings using predefined scales, which are the traditional phenotyping methods. Examples of such scales include the 1–9 severity scale proposed by [15] and the Florida 1–10 scale [16]. These visual ratings are easy to use and have successfully been applied to release several varieties in the past. However, the traditional measurements are subjective [17], labor-intensive, destructive, time-consuming, and, in the long run, expensive for breeding [13] because breeding populations are large and evaluated in several target locations.

The adoption of remote sensing technologies and high-throughput phenotyping (HTP) platforms enables breeders to collect spectral data on plants in different growth stages [18] and has the potential to overcome the shortfalls of traditional phenotyping and accelerate genetic gain. HTP measurements are non-destructive, have repeatability, are fast over large trials, and are less expensive in comparison with direct evaluations [19]. Plant spectral properties are genotype-specific, dependent on the morphology and physiology of the plant [20], and can thus be used to screen and select for traits of interest including disease resistance and yield. These methods generate large multidimensional datasets that linear models are limited in analyzing [21,22]. Instead, non-linear models and machine learning (ML) may be better suited for data with complex characteristics of non-linearity and outliers [23]. Several ML algorithms such as support vector machine (SVM; [24]), Random Forest (RF; [25]), and artificial neural networks (ANNs; [26]) have the capacity to adapt to complex data while constantly looping in search for the best parameters and models [21]. Since biotic and abiotic stresses often express similar symptoms, i.e., leaf wilting, defoliation, and senescence, optimization of the ML algorithms is critical for the identification of specific symptoms associated with particular stresses [22]. Several studies have compared various algorithms based on their accuracy to identify the best algorithm for a problem [27]. The RF, SVM, ANN, and k-nearest neighbors (KNNs) seem to be the most accurate algorithms in plant science thus far [28]. None of these algorithms possess significant advantages over the other in terms of accuracy [27]; however, several studies have suggested that SVM presents higher accuracy in classification problems. Successful applications of HTP and ML methods in plant breeding for the prediction of important agronomic traits including disease identification have been reported in several crops including tomato (*Solanum lycopersicum* L.) [29], maize (*Zea mays* L.) [30], radish

(*Raphanus sativus* L.) [31], and sugar beet (*Beta vulgaris*) [32]. In peanut breeding, HTP and ML methods have been applied for agronomic traits such as plant height [33], leaf area [34], and pod maturity [35], but these ML methods have barely been applied for selection for LLS resistance.

Previous efforts have demonstrated that remote sensing techniques can be successfully used for indirect selection for LLS resistance in Uganda [36] and Ghana [37]. RGB color space indices and vegetation indices that are strongly associated with LLS severity were identified and used in models to predict LLS severity. Although not always supported by the results [27], the application of ML has the potential to improve the accuracy of LLS severity prediction. Therefore, the objectives of this study were to (i) compare regression and classification prediction for accuracy and usability in breeding programs, and (ii) identify the best indices and models to be used by the breeding programs during LLS screening.

## 2. Materials and Methods

### 2.1. Genetic Material

A total of two hundred twenty-three (223) peanut genotypes derived from the African core collection were used in this study. The African core collection comprises a total of 300 genotypes assembled from nine different countries from West, East, and South Africa (Table 1). This population consists of different market types, i.e., Spanish (*A. hypogaea* sub. *vulgaris*), Valencia (*A. hypogaea* sub. *fastigiata*), Virginia (*A. hypogaea* sub. *hypogaea*), and the hybrid (combination between subspecies). The collection is believed to represent the total diversity available within the nominating breeding programs on the African continent and a detailed description of this population is found in [38,39]. Of the 223 lines, a total of 97 were planted in only Uganda, 43 in only Ghana, and 98 in both countries. Genotype selection was based on seed availability in these counties.

**Table 1.** Summary of the 223 genotypes of the African mini-core collection used in this study showing the market types and the country of origin. The hybrid represents genotypes from crossing different market types.

| Country of Origin | Market Types | | | | |
|---|---|---|---|---|---|
| | **Hybrid** | **Spanish** | **Valencia** | **Virginia** | **Total** |
| Ghana | 3 | 16 | 2 | 15 | 36 |
| Malawi | 9 | 9 | 1 | 8 | 27 |
| Mali | 1 | 17 | - | 3 | 21 |
| Mozambique | - | 13 | 1 | 4 | 18 |
| Niger | 9 | 24 | 1 | - | 34 |
| Senegal | 1 | 7 | - | 9 | 17 |
| Togo | 2 | 6 | 2 | 4 | 14 |
| Uganda | 3 | 15 | 4 | 20 | 42 |
| Zambia | 4 | 3 | | 7 | 14 |
| Total | 32 | 110 | 11 | 70 | 223 |

### 2.2. Site Description

The genotypes were evaluated in Ghana and Uganda. In Ghana, the experiment was set up at the Savanna Agricultural Research Institute research field located at Nyankpala (09°25′41″ N, 00°58′42″ W) in the northern region of Ghana during the May–October rainy season. In Uganda, two experiments were set up at Nakabango (0°31′ N, 33°12′ E) in Jinja during the September–December rainy season, and on station at the National Semi-Arid Resources Research Institute (NaSARRI; 1°35′ N; 33°35′ E) in Serere District during the March–July rainy season. Generally, the rainfall patterns during the growing seasons were similar across the three locations (Table 2). Moderate to high rainfall ranging between 83 mm and 260 mm was received during the growing months. In Serere, the highest precipitation was received in April (242 mm), 260 mm in August at Nyankpala, and

186 mm in November at Nakabango. The average monthly relative humidity was equally high across the three locations ranging between 67 and 83%. However, the average monthly temperature in Nyankpala was higher than at Serere and Nakabango. Average monthly temperatures at Nyankpala during the growing season ranged between 26 and 29 °C, while the average temperatures at Serere and Nakabango ranged between 21 and 24 °C, 5 °C less.

**Table 2.** Summary of the monthly weather data for Serere and Nakabango, in Uganda, and Nyankpala in Ghana, Africa.

| | Jan | Feb | Mar | Apr | May | Jun | Jul | Aug | Sep | Oct | Nov | Dec |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Weather Parameters | | | | | | Serere, 2020 | | | | | | |
| Av. Temperature (°C) | 24.3 | 25.1 | 24.6 | 23.9 | 22.8 | 21.5 | 21 | 22.2 | 23.1 | 22.8 | 23.4 | 23.8 |
| Relative humidity (%) | 67.9 | 64.9 | 75.1 | 76.7 | 82.8 | 83.8 | 82.1 | 79.4 | 77.6 | 80.4 | 75.1 | 63.2 |
| Wind speed (m/s) | 1.9 | 2.1 | 1.5 | 1.5 | 1.3 | 1.6 | 1.6 | 1.2 | 1.3 | 1.2 | 1.8 | 2.1 |
| Precipitation (mm) | 63.3 | 73.8 | 189.8 | 242.6 | 195.1 | 142.4 | 152.9 | 116 | 342.8 | 189.8 | 94.9 | 31.6 |
| | | | | | | Nakabango, 2021 | | | | | | |
| Av. Temperature (°C) | 22.6 | 22.9 | 23.4 | 22 | 21.9 | 21.5 | 21.4 | 22.8 | 22.4 | 23.4 | 22.8 | 23.4 |
| Relative humidity (%) | 73.6 | 74 | 74.4 | 83.7 | 83 | 79.5 | 73.8 | 71.2 | 75.4 | 73.6 | 77.1 | 72.9 |
| Wind speed (m/s) | 1.5 | 1.6 | 1.7 | 1.5 | 1.4 | 2.0 | 2.4 | 1.7 | 1.5 | 1.4 | 1.4 | 1.5 |
| Precipitation (mm) | 163.5 | 42.2 | 100.2 | 227.1 | 98.8 | 23.5 | 12.4 | 25 | 96.5 | 119.3 | 186.9 | 100.2 |
| | | | | | | Nyankapala, 2021 | | | | | | |
| Av. Temperature (°C) | 27.6 | 29.3 | 30.6 | 31.3 | 29.8 | 28.4 | 26.8 | 26 | 26.4 | 27.3 | 27.8 | 26.1 |
| Relative humidity (%) | 35.7 | 29.4 | 52.6 | 55.9 | 67.6 | 71.7 | 78.8 | 83.4 | 82.7 | 80.1 | 68.9 | 49.3 |
| Wind speed (m/s) | 2.7 | 2.6 | 2.4 | 2.6 | 2.2 | 2.0 | 2.1 | 1.6 | 1.3 | 1.3 | 1.3 | 2.4 |
| Precipitation (mm) | 0 | 0 | 15.8 | 55.3 | 101 | 83.3 | 161.9 | 261 | 173.1 | 126.9 | 14.3 | 0.1 |

*2.3. Experimental Set Up*

The trials were set up in an alpha lattice design and replicated thrice in Ghana and twice in Nakabango and Serere. In Uganda, each genotype was planted in 2 plots across the two replications each measuring 1 m × 0.45 m, with a spacing of 0.45 m between rows and 0.15 m within rows. The plots were separated by 0.90 m alleys and the replicates were separated by a 1.5 m alley. In Ghana, the genotypes were planted in two-meter single rows with 0.60 m between the rows.

*2.4. Data Collection*

2.4.1. Visual Ground Rating

LLS severity was visually scored using the 1–9 modified severity scale [15] (Subrahmanyam et al., 1995) starting from 4 weeks after planting and every 4 weeks until harvest at 14 weeks after planting. Remote sensing data were collected on the same days with visual ratings.

2.4.2. Normalized Difference Vegetation Index (NDVI)

A handheld spectroradiometer (GreenSeeker; Trimble Navigation, Sunnyvale, CA, USA) was used to measure and record canopy normalized difference vegetation index (NDVI) values. The GreenSeeker was held 0.60 m above the plant canopy and dragged through the entire length of each row to obtain the average NDVI value for each plot. NDVI was calculated from the equation.

$$NDVI = (NIR - R)/(NIR + R) \tag{1}$$

where R is the reflectance in the red band (660 nm) and NIR is the reflectance in the near-infrared band (780 nm).

### 2.4.3. Red–Green–Blue (RGB) Imaging

A Sony $\alpha$-6000 camera was used to take RGB images of the experimental plots. The camera was set to Auto mode so that the lens adjusts to the best sharpness and brightness based on the available light. The 58 mm camera lens was used and zero zoom was used for acquiring all the images. The camera was held 0.90 m above the plant canopy and images of the entire plot were taken. The RGB images were saved in JPEG format with a resolution of 350 dpi. RGB color space indices were then extracted (Table 3) from the images using the BreedPix 0.2 option of the CIMMYT maize scanner 1.6 plugin [open software] (GitHub–george-haddad/CIMMYT: CIMMYT MaizeScanner); Copyright 2015 Shawn Carlisle Kefauver, University of Barcelona [40]; produced as part of Image J/Fiji (http://fiji.sc/Fiji) (open source software (Fiji: ImageJ, with "Batteries Included") [41,42]).

**Table 3.** Red–green–blue (RGB) color space indices derived from the RGB images using BreedPix and indices derived from the combinations of various indices.

| RGB Indices | Basis of Derivation | Reference |
|---|---|---|
| Hue | Color description in form of angles [0–360° (0°—red; 60°—yellow; 120°—green; 240°—blue)] | [43] |
| a* | Green (−a*)–red (+a*) component in CIE-Lab color space | [43] |
| b* | Blue (−b*)–yellow (+b*) component in CIE-Lab color space | [43] |
| u* | Green (−u*)–red (+u*) component in CIE-Luv color space | [43] |
| v* | Blue (−v*)–yellow (+v*) component in CIE-Luv color space | [43] |
| Green Area (GA) | Percentage of Pixels from 60–120° of the hue angle | [44] |
| Greener Area (GGA) | Percentage of Pixels from 80–120° of the hue angle | [44] |
| Crop Senescence Index (CSI) | $100 \times (GA - GGA)/GA$ | [45] |
| Greenness index (GI) | GA/GGA | [46] |
| Greenness Product Index (GPI) | GA × GGA | [46] |
| Normalized Greenness Product Index (NGPI) | (GA − GGA)/(GA + GGA) | [46] |

### 2.5. Data Analysis

In this study, both regression and classification models were evaluated to support selection for LLS resistance. The traditional statistical models—stepwise linear regression (SLR) and partial least-squares regression (PLSR)—in addition to machine learning algorithms —support vector machine (SVM), Random Forest (RF), k-nearest neighbor (KNN), and the artificial neural network (ANN)—were used for regression modeling. All those models in addition to Naïve Bayes (NB) and linear discriminant analysis (LDA) were also used for the classification of LLS resistance. All these models were formulated in R version 4.3.1 [47] with the caret package framework [48]. A brief explanation of all models is as follows.

SVM is a supervised ML algorithm that works by finding the hyperplane that best separates the classes in the feature space. The hyperplane is chosen such that it maximizes the space between the nearest data points of the different classes [24]. The RF is an ensemble learning method that uses multiple decision trees during training and the output is a mode of the classes (classification) or the mean prediction (regression) of the individual trees [25]. The KNN is a non-parametric supervised learning algorithm that is not necessarily trained to produce a model. Instead, the unknown class is compared with the rest of the data and assigned to the most common class. It assigns the labels based on majority votes or averages of the "K" number of nearest neighbors. The performance depends on the value of K (number of neighbors to consider). The smaller the K, the more sensitive the model is to noise in the data; however, larger K values might cause the model to miss the patterns in the data due to overfitting [49]. The ANN is a computational model that consists of interconnected nodes called neurons arranged into different layers [26]. The layers consist of the input layer, hidden layer, and output layer. Each neuron in the network processes input data and passes it to the next layer. During the training process, the network learns through a process called back propagation, adjusting the weights to minimize the error between the predicted output and the actual output. The Naïve Bayes (BN) is a

probabilistic algorithm that is based on the Bayes theorem and assumes that features are conditionally independent given the class label [50]. The probability of each class is calculated given the input variables and the class with the highest probability is chosen as the prediction. Stepwise regression is a technique used for feature selection in logistic regression models. It involves adding or removing predictors from the model based on their statistical significance. The algorithm starts with no predictors and iteratively adds or removes predictors based on significance or model fit. This aims to find the subset that best describes the variation while minimizing overfitting. Linear discriminant analysis (LDA) is a supervised classification technique used to find a linear combination of features that best separates different classes in the data. LDA assumes that the features are normally distributed and then calculates linear discriminants, which are axes that maximize the separation between classes while minimizing the variation within each class. During prediction, LDA assigns new data points to the class with the highest posterior probability based on the linear discriminant functions [51].

For the regression models, the disease severity scores (1–9) were considered as continuous variables while RGB color indices and NDVI were used as predictor variables. The data were randomly split into training and testing datasets in a 70:30% ratio. The "train" function of the caret package was used to train six different regression models: SLR, PLSR, SVM, RF, KNN, and ANN. The linear kernel was used in the SVM and a cost value of 5 was applied to control for the overfitting of the model. For the RF, a total of 400 decision trees were used in the ensemble and 5 variables were considered at each split. The ANN comprised four layers: the input layer, two hidden layers, and the output layer. The hidden layers consisted of eight neurons: five in the first hidden layer and three in the second. For the KNN model, a 'K' value of 11 nearest neighbors was used.

The trained models were validated on the testing dataset (30%) using $R^2$ and the root-mean-squared error (RMSE) between actual LLS and estimated LLS ratings.

RF, SVM, KNN, ANN, NB, and LDA were also evaluated for classifying genotype selection based on LLS ratings. These models were evaluated for three different classification modes: (1) two-class or binary classification where LLS scores of 1–6 were designated as acceptable and LLS scores of 7–9 as unacceptable; (2) three-class classification where LLS scores of 1–3 were designated as resistant, 4–6 as moderately resistant/tolerant, and 7–9 as susceptible; and (3) four-class classification where LLS scores of 1–2 were designated as highly resistant, 3–4 as resistant, 5–6 as moderately resistant/tolerant, and 7–9 as susceptible. The "train" function of the caret package [48] was used to train the classification models using 70% of the data using the same inputs as in regression models, i.e., RGB color indices and NDVI as predictors while resistance classes as observations. A fivefold cross-validation with three repeats was used for training the selected algorithms and their performance was evaluated using accuracy, specificity, sensitivity, Recall, and F1-score from confusion matrices.

$$\text{Accuracy} = (TP + TN)/(TP + TN + FP + FN) \tag{2}$$

$$\text{Sensitivity (Recall)} = TP/(TP + FN) \tag{3}$$

$$\text{Specificity} = TN/(TN + FP) \tag{4}$$

$$\text{Precision} = TP/(TP + FP) \tag{5}$$

$$\text{F1} = 2 \times (\text{Precision} \times \text{Recall})/(\text{Precision} + \text{Recall}) \tag{6}$$

where TP is true positives, TN is true negatives, FP is false positives, and FN is false negatives [52].

## 3. Results

### 3.1. Disease Distribution across the Three Locations

The LLS severity scores showed high disease pressure in all test locations (Figure 1). The medians in Figure 1 indicate that the disease was higher in Serere compared to Naka-

bango and Nyankpala. In addition, the Virginia market types generally had the lowest severity and disease spread across all locations. Within Nyankpala, the Valencia market types showed the lowest LLS severity, while in Nakabango, the hybrid type experienced the highest LLS severity compared to the Spanish, Valencia, and Virginia types which showed similar severities. In Nyankpala and Serere, the hybrid and Spanish types appeared to have been affected more by LLS compared to subspecies Virginia and Spanish.
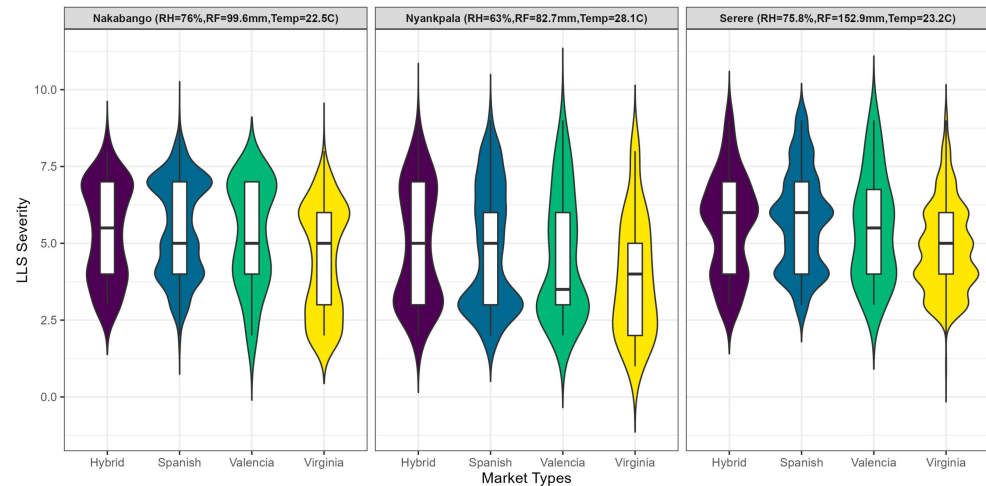


**Figure 1.** Distribution of LLS severity scores at three experimental sites for the different peanut market types (Spanish, Valencia, Virginia, and Hybrid) used in this study. The locations used were Serere and Nakabango in Uganda and Nyankpala in Ghana.

### 3.2. Late Leaf Spot Estimation Using Machine Learning

Overall, ML models performed better than the traditional statistical models (Figure 2). The ANN was the best-performing model ($R^2$: 0.81, RMSE: ~22%), followed by RF ($R^2$: 0.75, RMSE: ~25%) and KNN ($R^2$: 0.73, RMSE: ~26%). Traditional statistical models PLSR and SLR had lower accuracies ($R^2$: 0.56, RMSE: ~33%) compared to the ML models, except SVM which was the least performing ($R^2$: 0.52, RMSE: ~35%). The RF model identified NDVI, hue, CSI, b*, and v* as the most important predictor variables while GGA, NDGI, and GI were identified as the least important.
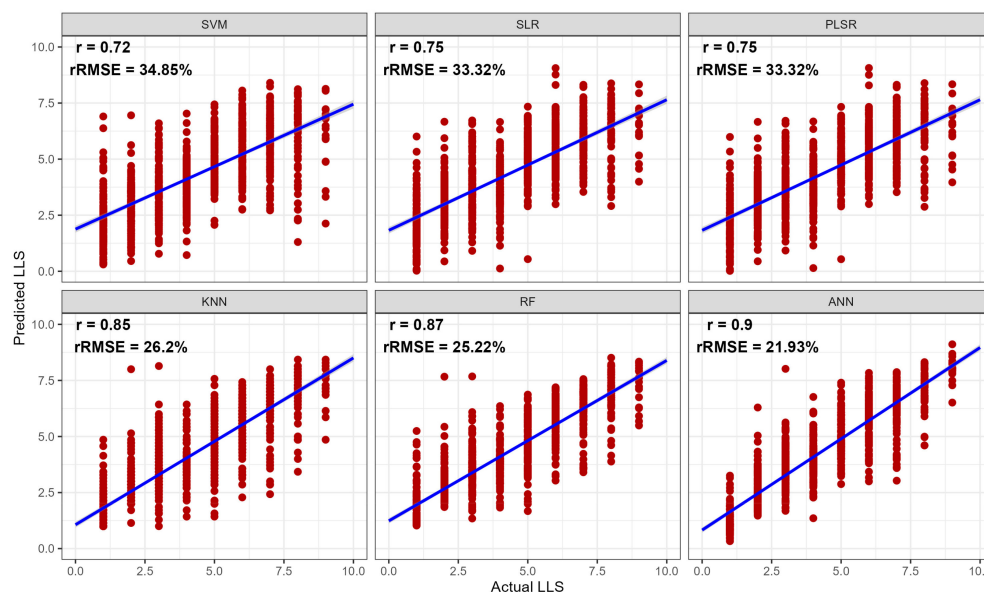
**Figure 2.** Associations between the actual and predicted LLS scores of the test dataset using six models derived from the testing dataset. The models were trained on 70% of the data and tested on 30% of the testing data. ANN—artificial neural network; KNN—K-nearest neighbors; PLSR—partial least-squares regression; RF—Random Forest; SLR—stepwise logistic regression; SVM—support vector machine.

### 3.3. LLS Severity Classification Using Machine Learning

#### 3.3.1. Binary Classification Models

The RF was the best performer with an accuracy of 90%, followed by KNN with 89.5% accuracy, and ANN with 88.8% accuracy (Table 4). NB and LDA had lower classification accuracies of 83% and 85.9%, respectively. Traditional classification models, PLSC and SLC, had classification accuracies of 85.9% and 87.6%, respectively. Nonetheless, all the models were able to effectively differentiate between the acceptable and non-acceptable classes of LLS resistance. The RF model also had the highest Kappa of 66%, followed by ANN with 64% and KNN with 63.7%, while PLSC and LDA had the highest sensitivities of 96% and 94.6% and the lowest specificities of 41% and 47%, respectively. The ANN and RF had the highest specificities of 74% and 70.5%, respectively.

**Table 4.** Performance of binary classification models in predicting late leaf spot (LLS) resistance of peanut genotypes.

| Metric (%) | ANN | KNN | LDA | PLSC | NB | RF | SLC | SVM |
|---|---|---|---|---|---|---|---|---|
| Accuracy | 88.85 | 89.59 | 85.90 | 85.99 | 83.41 | 90.05 | 87.65 | 86.73 |
| Kappa | 64.09 | 63.75 | 47.06 | 44.37 | 50.57 | 66.25 | 56.75 | 52.56 |
| Sensitivity | 92.2 | 94.92 | 94.69 | 96.16 | 86.44 | 94.46 | 93.90 | 93.90 |
| Precision | 94.01 | 92.51 | 88.77 | 87.82 | 92.73 | 93.41 | 91.22 | 90.23 |
| Recall | 92.2 | 94.92 | 94.69 | 96.16 | 86.44 | 94.46 | 93.90 | 93.90 |
| F1 | 93.1 | 93.70 | 91.63 | 91.80 | 89.47 | 93.93 | 92.54 | 92.03 |
| Specificity | 74.00 | 66.00 | 47.00 | 41.00 | 70.00 | 71.00 | 60.00 | 55.00 |

ANN—artificial neural network; KNN—K-nearest neighbors; LDA—linear discriminant analysis; NB—Naïve Bayes; PLSC—partial least-squares classification; RF—Random Forest; SLC—stepwise logistic classification; SVM—support vector machine.

#### 3.3.2. Multiclass Resistance Classification

The multiclass resistance classification using ML models apparently performed better than the traditional statistical models but worse compared to the binary classification (Table 5, Figures 3 and 4). Moreover, classification accuracies were higher for the three-class resistance classification compared to the four-class resistance classification (Table 5, Figures 3 and 4).

The RF model was the best-performing model with 78% accuracy and a Kappa score of 65 for the 3-class classification, with 73% accuracy and a Kappa score of 64 for the 4-class classification. The KNN and ANN were the second- and third-best performing models for both the 3- and 4-class multiclassification, while NB and SVM were the least-performing models with accuracies of 63% and 68%, respectively, for the 3-class classification, and 58% and 67%, respectively, for the 4-class classification. Traditional statistical classification models, SLC and PLSC, had accuracies of 73% and 69%, respectively, for the 3-class classification, and 70% and 65%, respectively, for the 4-class classification. These statistical models performed better than both NB and SVM models for the 4-class classification.

**Table 5.** Performance of multiclass algorithms in categorizing genotypes into respective predefined late leaf spot (LLS) resistance classes.

| Model | 3 Classes | | 4 Classes | |
|---|---|---|---|---|
| | Accuracy | Kappa | Accuracy | Kappa |
| ANN | 76.10 | 62.22 | 74.45 | 65.89 |
| KNN | 77.42 | 64.25 | 73.36 | 64.19 |
| LDA | 72.07 | 55.25 | 67.10 | 55.56 |
| NB | 62.86 | 42.21 | 57.51 | 43.01 |
| PLSC | 68.66 | 49.57 | 65.90 | 54.04 |
| RF | 77.60 | 64.75 | 72.81 | 63.51 |
| SLC | 72.72 | 56.85 | 69.86 | 59.54 |
| SVM | 68.20 | 49.43 | 63.04 | 50.27 |

ANN—artificial neural network; KNN—K-nearest neighbors; LDA—linear discriminant analysis; NB—Naïve Bayes; PLSC—partial least-squares classification; RF—Random Forest; SLC—stepwise logistic classification; SVM—support vector machine.
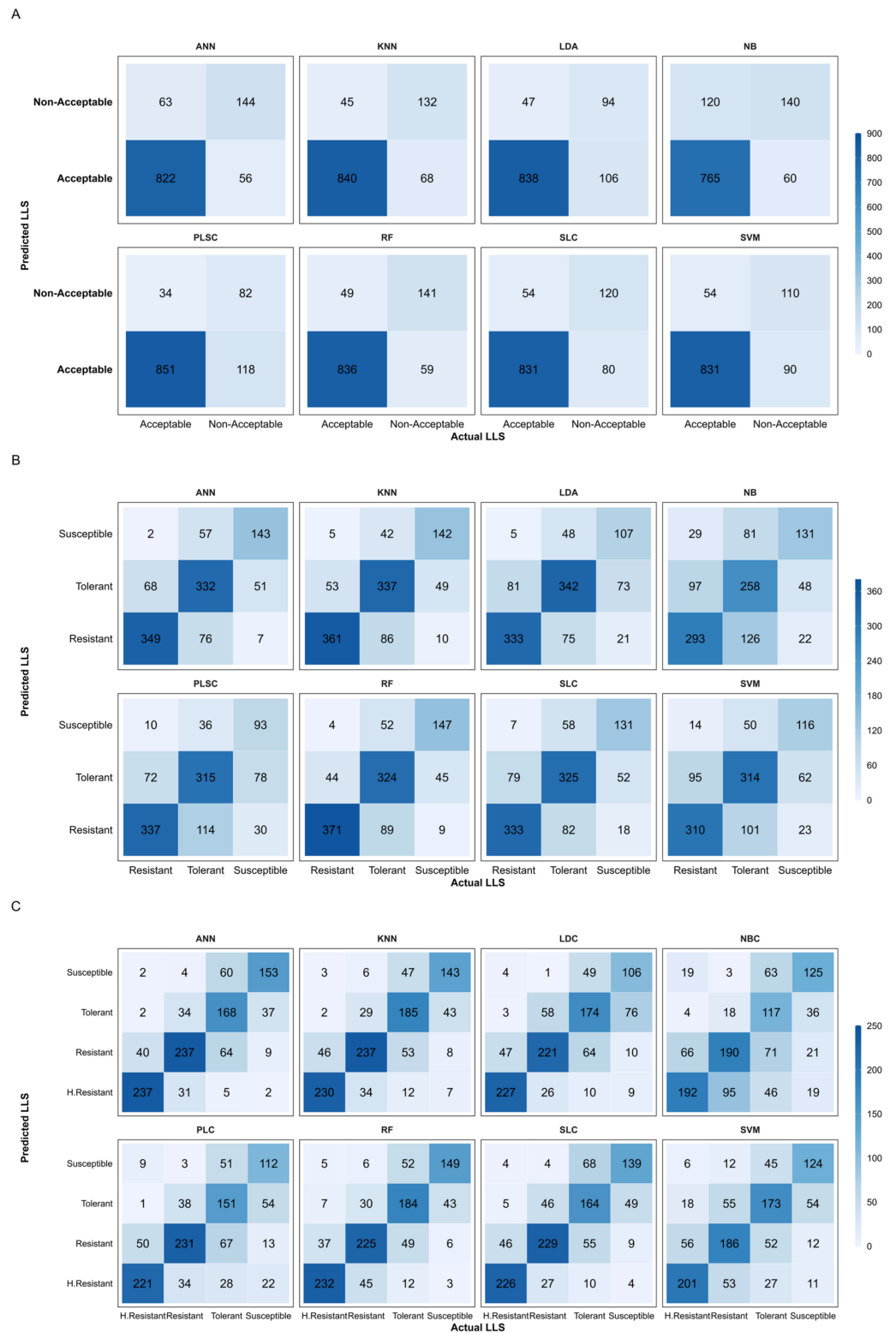
**Figure 3.** Confusion matrices of predicted and actual LLS categories of (**A**) binary classification, (**B**) 3-class multiclassification, and (**C**) 4-class multiclassification. The confusion matrices show the correctly classified and misclassified plots by the different methods. Each number and color placed in each box represents the number of plots classified by the different methods, and the diagonals in B and C represent classes correctly classified by the models. ANN—artificial neural network; KNN—K-nearest neighbors; LDA—linear discriminant analysis; NB—Naïve Bayes; PLSC—partial least-squares classification; RF—Random Forest; SLC—stepwise logistic classification; SVM—support vector machine.
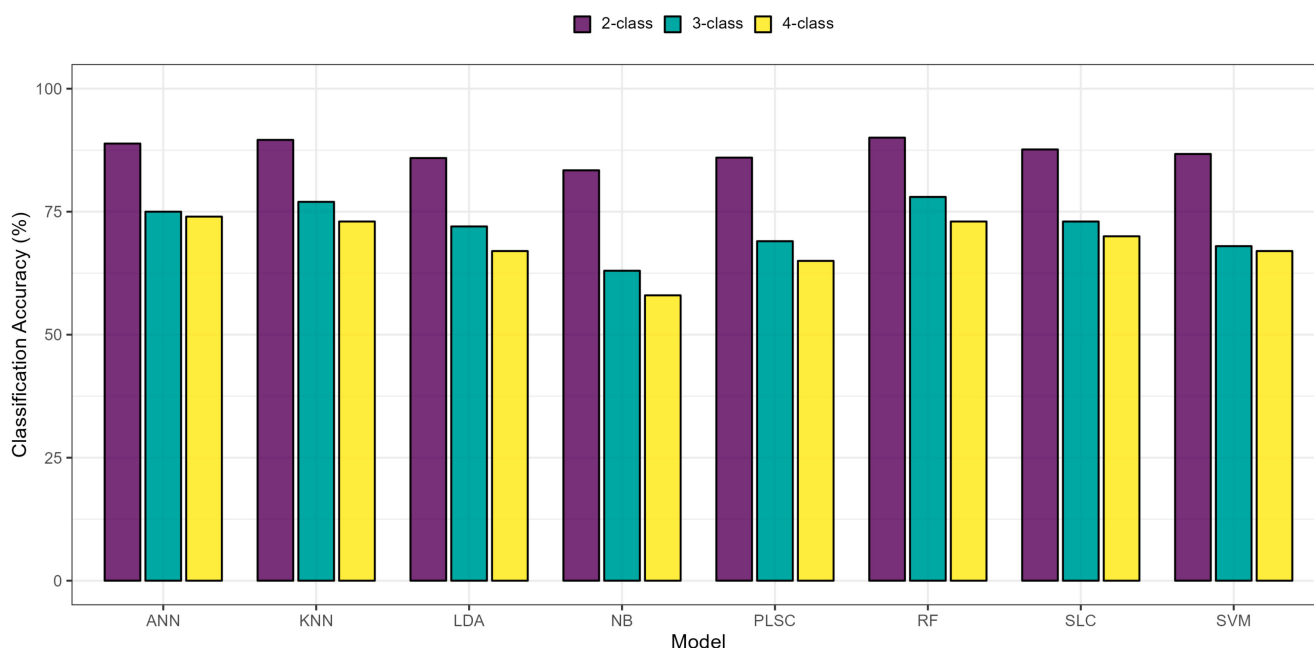
**Figure 4.** Comparison of the accuracy of the different classification models. (a) binary classification, (b) 3-class multiclassification, and (c) 4-class multiclassification using the different algorithms. ANN—artificial neural network; KNN—K-nearest neighbors; LDA—linear discriminant analysis; NB—Naïve Bayes; PLSC—partial least-squares classification; RF—Random Forest; SLC—stepwise logistic classification; SVM—support vector machine.

## 4. Discussion

Late leaf spot (LLS) is an important fungal disease that is favored by high relative humidity and moderate temperature [53,54]. This was evident from the results: LLS severities were highest in Serere followed by Nakabango and Nyankpala. All these locations experienced high relative humidity and moderate temperatures that could have led to high disease pressure with severity medians of 5 and maximum values up to 9 (Figure 1). Butler (1990) indicated earlier that temperatures between 18 °C and 27 °C favored the early development of LLS lesions. Earlier studies have also identified Nakabango and Serere as LLS hotspots in Uganda [1].

Currently, the selection for LLS-resistant genotypes is dependent on the visual rating of the disease severity on a scale of 1–9 or 1–10 [15]. These ratings are based on the presence of lesions on the leaves and leaf defoliation. A scale of 1–9 is the most used in the peanut breeding programs in Uganda and Ghana and is categorized into resistant (1–3), tolerant to moderately resistant (4–6), and susceptible (7–9), and several varieties have been selected and released using this criterion, for example, Serenut 5R [55] and Naronut 1R [56] in Uganda, and the Babile series [57] in Ethiopia. Selection for the most resistant (1–3) lines is the most ideal situation; however, because selection for several other traits is important for production, i.e., yield, maturity, resistance to other diseases, etc., they play important roles in variety release decisions. To accommodate other traits, lines with high (1–3) and medium (4–6) LLS resistance are commonly retained for advancement and potential release by breeders. In addition, the use of a visual rating is subjective to human error and depends on the individual's expertise and ability to visually discern LLS from other lesions on the leaves, the time of rating, and the peanut market type, among a few. It is relatively easy to identify the most resistant (1) and the most susceptible (9) genotypes for LLS, but differentiating middle scores such as between 3 and 4 and between 6 and 7 can be difficult and can have consequences for the selection process.

For the LLS severity estimation, the ANN model depicted the highest performance ($R^2$: 0.81, RMSE: ~22%), indicating that the predicted values were in high agreement

with the actual observed values (Figure 2). This is the advantage of the ANN, combining multiple layers of neurons, in learning complex data patterns from multidimensional data including remote sensing and making accurate predictions [58,59]. ANNs have been widely applied in various plant breeding programs including for yield modeling and prediction in soybean [60], and root regeneration [61] and drought tolerance screening [62] in wheat. The ANN has been applied for the identification and prediction of several crop diseases such as late blight (*Phytophthora infestans*) in tomato with a prediction accuracy of 66% [63], *Botrytis cinerea* of eggplant (*Solanum melongena* L.) with a prediction accuracy of 70% [64], and seedling diseases of orchids (*Phalaenopsis* spp.) with a prediction accuracy of 89.6% [65].

The classification of LLS severity into two classes (acceptable and non-acceptable) exhibited robust discriminatory performance (83 to 90%) for all models (Table 4). However, high discriminatory accuracy does not necessarily translate into high selection accuracy. The number of genotypes selected using this criterion is large, often between 760 and 850 (Figure 3A), and yet the breeding target is to reduce large populations to a small, manageable number [66]. Increasing the number of classification classes to three (resistant, tolerant, and susceptible) reduced the classification accuracy of the models (63 to 78%) but increased the precision of selection. This reduction in accuracy is due to the increased complexity of the model due to increased possibilities and decision boundaries, and the similarity of features between the classes leading to increased misclassifications [67]. Increasing to four classes (highly resistant, resistant, tolerant, and susceptible) further reduced classification accuracies (58 to 74%) but also improved the potential applicability of this classification model in a breeding program. This is because four classes allowed the reduction in selected genotypes to manageable numbers. RF, for example, predicted 184 cases of highly resistant genotypes in agreement with the visual rating when four classes were used versus 371 cases when the three-class classification was applied (Figure 3). In summary, the classification accuracy of the models decreased with an increase in the classification classes (Figure 4), but the practicality for peanut breeding seems to have increased with the increase in the number of classes.

For resistance classification, RF was consistently the best classifier for the binary, 3-class, and 4-class condition (Tables 4 and 5). RF follows an efficient training phase and therefore obtains a high generalization accuracy. This is because RF trains multiple decision trees using multiple random subsamples of the original dataset where the generalization error reduces as more trees are added, thereby also reducing the model overfitting [25]. This supported accurate identification of the true positives and true negatives and minimized the possibility of identifying false positives and false negatives (Figure 3). This is particularly important in breeding where false positives and false negatives can impede genetic gain [68]. RF has been successfully applied in studies of various diseases and crops such as the identification of *Alternaria* diseases of rape oil seed [69] with a discrimination accuracy of 82.6%, the identification of several tomato diseases with a classification accuracy of 95.2% [70], and the differentiation between wheat rust (*Puccinia recondita* f. sp. *tritici*) and rye rust (*Puccinia recondita* f. sp. *recondita*) with classification accuracies of 96.6% and 91.7% when using spectral wavelengths and vegetation indices, respectively [71]. Although, in this work, all models had a misclassification of several classes, there was a relatively clear distinction of the genotypes into the respective resistance classes (Figure 3), which is the most important purpose of the models.

The SVM, which is one of the most used algorithms in plant science, was also used in this study. However, its accuracy was low compared to the other ML algorithms applied in this study and those reported in previous studies. This could be due to the fact that the data classes were imbalanced, and there were close margins between the classification classes, all leading to overfitting of the model [24].

Accurate selection is important for genetic gain [19] and making wrong selections such as classifying a susceptible genotype as a resistant genotype can be detrimental to the breeding goals. Visualization of model performance using a confusion matrix helps show the number of correct and incorrect classifications made by the model as well as

the misclassification of various classes (Figure 3). This enables the calculation of other performance metrics such as sensitivity, precision, recall, and specificity which correct for false positives and negatives [72].

The indices identified by RF to be highly important included NDVI, CSI, hue, GI, b*, and v*, which have also been identified as critical in earlier studies for LLS detection [36,73]. LLS severity is associated with leaf defoliation, which reduces the leaf area index of the canopy and the NDVI as well [74]. The CSI, GI, b*, and v* which represent the senescent fraction of the canopy are related to the senescence and canopy yellowing associated with LLS severity [7]. The increase in phenotyping accuracy of LLS severity in breeding is key for genetic gain improvement [75]. The use of remote sensing HTP methods has the potential to remove the bias and subjectivity associated with traditional phenotyping and reduce the time spent on data collection. These methods are repeatable [76], for which the selection response for the trait of interest can be improved over traditional methods [18].

Although initial investments in remote sensing sensors are high, the cost of phenotyping is eventually reduced in the long run [19], thereby reducing the cost of developing new varieties. The tools used in this study were handheld and have demonstrated a high repeatability of measurements [36,76,77]. However, they do not allow the throughput required in breeding. Therefore, the adoption of unmanned aerial vehicles for data collection is being proposed to increase the amount and throughput of data collected and reduce the cost per data point collected [78]. It is also important to note that the performance of ML models tends to be crop- and environment-dependent and therefore often obtain inconsistent and non-generalizable results [79]. This could be eliminated by several approaches, primarily by either gathering data from multiple agroclimatic conditions (crop, varieties, weather, and soil, among others) and later testing ensemble models, i.e., stacking multiple ML algorithms together as one [31,80–82]. A recent study by [81] utilized an ensemble of RF, KNN, and SVR models to predict alfalfa yield that yielded a much better accuracy than individual models.

## 5. Conclusions

Late leaf spot is an important constraint to peanut production worldwide. The accurate phenotyping of the disease is important for selection for resistance in breeding. The adoption of remote sensing methods together with machine learning is considered an alternative to make selection for LLS resistance faster and accurate. The objectives of this study were to compare between regression and classification for breeding, and to identify the best models and indices to be used for selection. The results of our study indicate that the ANN, RF, and KNN are the best performing algorithms for both regression and classification methods. The classification accuracy of all algorithms decreased with the increase in the classification classes. Of the three different modes of classification tested, the four-class classification was the most practical for selection although the classification accuracies were lower compared to the other two modes. NDVI, and RGB indices CSI, hue, GI, b*, and v*, were identified as the most important indices for selection for LLS resistance. Our study demonstrated the efficacy of machine learning methods for the selection for LLS resistance in peanut breeding. While our results demonstrate promising outcomes for the utilization of machine learning for selection, there were some limitations in this study. The data acquisition was limited to handheld sensors which generated a small set of data and yet the machine learning models needed a large amount of data for efficient training of the models. Another limitation faced was the quality of the visual scores. Since the data were collected from different countries, it was challenging to control the influence of human error and ascertain the consistence of the data collected. This could have affected the accuracy of the models developed. Therefore, for future work, there is a need to develop standardized data collection protocols and incorporate unmanned aerial vehicles (UAVs) equipped with multiple sensors for the fast and objective collection of large amounts of spectral data for model training. There is also a need to include weather parameters such as temperature and relative humidity in the models, and explore more robust analytical methodologies

such as deep learning to enhance classification accuracy and selection for LLS resistance in peanut breeding.

## References

1. Okello, D.K.; Biruma, M.; Deom, C.M. Overview of groundnuts research in Uganda: Past, present and future. *Afr. J. Biotechnol.* **2010**, *9*, 6448–6459.

2. FAOSTAT. *Food and Agriculture Organization of the United Nations*; FAOSTAT Statistical Database: Rome, Italy, 2023; Available online: https://www.fao.org/faostat/en/#data/QCL (accessed on 22 March 2024).

3. Deom, C.M.; Okello, D.K. *Developing Improved Varieties of Groundnut*; Sivasankar, S., Ed.; Burleigh Dodds Science Publishing: Cambridge, UK, 2018; pp. 145–176.

4. Oteng-Frimpong, R.; Konlan, S.P.; Denwar, N.N. Evaluation of Selected Groundnut (*Arachis hypogaea* L.) Lines for Yield and Haulm Nutritive Quality Traits. *Int. J. Agron.* **2017**, *2017*, 7479309. [CrossRef]

5. Abady, S.; Shimelis, H.; Janila, P.; Mashilo, J. Groundnut (*Arachis hypogaea* L.) improvement in sub-Saharan Africa: A review. *Acta Agric. Scand. Sect. B-Soil Plant Sci.* **2019**, *69*, 528–545. [CrossRef]

6. Waliyar, F.; Kumar, P.L.; Ntare, B.R.; Monyo, E.; Nigam, S.N.; Reddy, A.S.; Osiru, M.; Diallo, A.T. *A Century of Research on Groundnut Rosette Disease and Its Management*; Information Bulletin No. 75; International Crops Research Institute for the Semi-Arid Tropics: Patancheru, India, 2007.

7. Singh, M.P.; Erickson, J.E.; Boote, K.J.; Tillman, B.L.; van Bruggen, A.H.C.; Jones, J.W. Photosynthetic consequences of late leaf spot differ between two peanut cultivars with variable levels of resistance. *Crop Sci.* **2011**, *51*, 2741–2748. [CrossRef]

8. Culbreath, A.K.; Stevenson, K.L.; Brenneman, T.B. Management of Late Leaf Spot of Peanut with Benomyl and Chlorothalonil: A Study in Preserving Fungicide Utility. *Plant Dis.* **2002**, *86*, 349–355. [CrossRef] [PubMed]

9. Shokes, F.; Gorbert, D.W.; Jackson, L.F. Control of Early and Late Leafspot on two peanut cultivars. *Peanut Sci.* **1983**, *10*, 17–21. [CrossRef]

10. Lamon, S.; Chu, Y.; Guimaraes, L.A.; Bertioli, D.J.; Leal-bertioli, S.C.M.; Santos, J.F.; Godoy, I.J.; Culbreath, A.K.; Holbrook, C.C.; Ozias-akins, P. Characterization of peanut lines with interspecific introgressions conferring late leaf spot resistance. *Crop Sci.* **2021**, *61*, 1724–1738. [CrossRef]

11. Kalule Okello, D.; Monyo, E.; Michael, D.C.; Jane, I.; Herbert Kefa, O. *Groundnut Production Guide for Uganda: Recommended Practices for Farmers*; National Agricultural Research Organisation: Entebbe, Uganda, 2013.

12. Natsugah, S.K.; Abudulai, M.; Oti-Boateng, C.; Brandenburg, R.; Jordan, D. Management of Leaf Spot Diseases of Peanut with Fungicides and Local Detergents in Ghana. *Plant Pathol. J.* **2007**, *6*, 248–253.

13. Araus, J.L.; Cairns, J.E. Field high-throughput phenotyping: The new crop breeding frontier. *Trends Plant Sci.* **2014**, *19*, 52–61. [CrossRef]

14. Foster, D.J.; Wynne, J.C.; Beute, M.K. Evaluation of detached leaf culture for screening peanuts for leafspot resistance. *Peanut Sci.* **1980**, *7*, 98–100. [CrossRef]

15. Subrahmanyam, P.; McDonald, D.; Waliayar, F.; Reddy, L.J.; Nigam, S.N.; Gibbons, R.W.; Rao, V.R.; Singh, A.K.; Pande, S.; Reddy, P.M.; et al. *Screening Methods and Sources of Resistance to Rust and Late Leaf Spot of Groundnut*; Information Bulletin No. 47; International Crops Research Institute for the Semi-Arid Tropics: Patancheruvu, India, 1995; p. 21.

16. Chiteka, Z.A.; Gorbet, D.W.; Shokes, F.M.; Kucharek, T.A.; The, F. Components of Resistance to Late Leafspot in Peanut. I. Levels and Variability—Implications for Selection. *Peanut Sci.* **1988**, *15*, 25–30. [CrossRef]

17. Milberg, P.; Bergstedt, J.; Fridman, J.; Odell, G.; Westerberg, L. Observer bias and random variation in vegetation monitoring data. *J. Veg. Sci.* **2008**, *19*, 633–644. [CrossRef]

18. Rutkoski, J.; Poland, J.; Mondal, S.; Autrique, E.; Pérez, L.G.; Crossa, J.; Reynolds, M.; Singh, R. Canopy Temperature and Vegetation Indices from High-Throughput Phenotyping Improve Accuracy of Pedigree and Genomic Selection for Grain Yield in Wheat. *G3 Genes Genomes Genet.* **2016**, *6*, 2799–2808. [CrossRef] [PubMed]

19. Araus, J.L.; Kefauver, S.C.; Zaman-allah, M.; Olsen, M.S.; Cairns, J.E. Translating High-Throughput Phenotyping into Genetic Gain. *Trends Plant Sci.* **2018**, *23*, 451–466. [CrossRef] [PubMed]

20. Schweiger, A.K.; Cavender-Bares, J.; Townsend, P.A.; Hobbie, S.E.; Madritch, M.D.; Wang, R.; Tilman, D.; Gamon, J.A. Plant spectral diversity integrates functional and phylogenetic components of biodiversity and predicts ecosystem function. *Nat. Ecol. Evol.* **2018**, *2*, 976–982. [CrossRef] [PubMed]

21. Behmann, J.; Mahlein, A.; Rumpf, T.; Ro, C.; Plu, L. A review of advanced machine learning methods for the detection of biotic stress in precision crop protection. *Precis. Agric.* **2014**, *16*, 239–260. [CrossRef]

22. Watt, M.; Fiorani, F.; Usadel, B.; Rascher, U.; Muller, O.; Schurr, U. Phenotyping: New Windows into the Plant for Breeders. *Annu. Rev. Plant Biol.* **2020**, *71*, 689–712. [CrossRef] [PubMed]

23. White, J.W.; Andrade-Sanchez, P.; Gore, M.A.; Bronson, K.F.; Coffelt, T.A.; Conley, M.M.; Feldmann, K.A.; French, A.N.; Heun, J.T.; Hunsaker, D.J.; et al. Field-based phenomics for plant genetics research. *Field Crops Res.* **2012**, *133*, 101–112. [CrossRef]

24. Vapnik, V.N. The Nature of Statistical Learning Theory. In *Statistics for Engineering and Information Science*, 2nd ed.; Springer: New York, NY, USA, 2000.

25. Breiman, L. *Statistics*; Department University of California: Berkeley, CA, USA, 2001.

26. Pal, S.K.; Member, S.; Mitra, S.; Member, S. Multilayer Perceptron, Fuzzy Sets, and Classification. *IEEE Trans. Neural Netw.* **1992**, *3*, 683–697. [CrossRef]

27. Maxwell, A.E.; Warner, T.A.; Fang, F. Implementation of machine-learning classification in remote sensing: An applied review. *Int. J. Remote Sens.* **2018**, *39*, 2784–2817. [CrossRef]

28. Singh, A.; Ganapathysubramanian, B.; Singh, A.K.; Sarkar, S. Machine Learning for High-Throughput Stress Phenotyping in Plants. *Trends Plant Sci.* **2016**, *21*, 110–124. [CrossRef] [PubMed]

29. Raza, S.; Prince, G.; Clarkson, J.P.; Rajpoot, N.M. Automatic Detection of Diseased Tomato Plants Using Thermal and Stereo Visible Light Images. *PLoS ONE* **2015**, *10*, e0123262. [CrossRef] [PubMed]

30. Wiesner-Hanks, T.; Wu, H.; Stewart, E.; DeChant, C.; Kaczmar, N.; Lipson, H.; Gore, M.A.; Nelson, R.J. Millimeter-Level Plant Disease Detection From Aerial Photographs via Deep Learning and Crowdsourced Data. *Front. Plant Sci.* **2019**, *10*, 1550. [CrossRef] [PubMed]

31. Dang, L.M.; Ibrahim Hassan, S.; Suhyeon, I.; Kumar Sangaiah, A.; Mehmood, I.; Rho, S.; Seo, S.; Moon, H. UAV based wilt detection system via convolutional neural networks. *Sustain. Comput. Inform. Syst.* **2020**, *28*, 100250. [CrossRef]

32. Rumpf, T.; Mahlein, A.K.; Steiner, U.; Oerke, E.C.; Dehne, H.W.; Plümer, L. Early detection and classification of plant diseases with Support Vector Machines based on hyperspectral reflectance. *Comput. Electron. Agric.* **2010**, *74*, 91–99. [CrossRef]

33. Sarkar, S.; Thomason, W.; Cazenave, A.; Abbot, L.; Balota, M.; Oakes, J.; Mccall, D. High-throughput measurement of peanut canopy height using digital surface models. *Plant Phenome J.* **2020**, *3*, e20003. [CrossRef]

34. Sarkar, S.; Cazenave, A.B.; Oakes, J.; McCall, D.; Thomason, W.; Abbott, L.; Balota, M. Aerial high-throughput phenotyping of peanut leaf area index and lateral growth. *Sci. Rep.* **2021**, *11*, 21661. [CrossRef] [PubMed]

35. Brunno, J.; Souza, C.; Luns, S.; De Almeida, H.; De Oliveira, M.F. Integrating Satellite and UAV Data to Predict Peanut Maturity upon Artificial Neural Networks. *Agronomy* **2022**, *12*, 1512. [CrossRef]

36. Chapu, I.; Okello, D.K.; Okello, R.C.O.; Odong, T.L.; Sarkar, S.; Balota, M. Exploration of alternative approaches to phenotyping of late leaf spot and groundnut rosette virus disease for groundnut breeding. *Front. Plant Sci.* **2022**, *13*, 912332. [CrossRef]

37. Sie, E.K.; Oteng-Frimpong, R.; Kassim, Y.B.; Puozaa, D.K.; Adjebeng-Danquah, J.; Masawudu, A.R.; Ofori, K.; Danquah, A.; Cazenave, A.B.; Hoisington, D.; et al. RGB-image method enables indirect selection for leaf spot resistance and yield estimation in a groundnut breeding program in Western Africa. *Front. Plant Sci.* **2022**, *13*, 957061. [CrossRef]

38. Achola, E.; Wasswa, P.; Fonceka, D.; Paul, J.; Prasad, C.; Ozias, P.; Jean, A.; Rami, F.; Michael, C.; David, D.; et al. Genome-wide association studies reveal novel loci for resistance to groundnut rosette disease in the African core groundnut collection. *Theor. Appl. Genet.* **2023**, *136*, 35. [CrossRef] [PubMed]

39. Conde, S.; Rami, J.-F.; Okello, D.K.; Sambou, A.; Muitia, A.; Oteng-Frimpong, R.; Makweti, L.; Sako, D.; Faye, I.; Chintu, J.; et al. The groundnut improvement network for Africa (GINA) germplasm collection: A unique genetic resource for breeding and gene discovery. *G3 Genes Genomes Genet.* **2024**, *14*, jkad244. [CrossRef] [PubMed]

40. Kefauver, S.C.; Romero, A.G.; Buchaillot, M.L.; Vergara-Diaz, O.; Fernandez-Gallego, J.A.; El-Haddad, G.; Akl, A.; Araus, J.L. Open-Source Software for Crop Physiological Assessments Using High Resolution RGB Images. In Proceedings of the IGARSS 2020-2020 IEEE International Geoscience and Remote Sensing Symposium, Waikoloa, HI, USA, 26 September–2 October 2020; pp. 4359–4362. [CrossRef]

41. Rueden, C.T.; Schindelin, J.; Hiner, M.C.; DeZonia, B.E.; Walter, A.E.; Arena, E.T.; Eliceiri, K.W. ImageJ2: ImageJ for the next generation of scientific image data. *BMC Bioinform.* **2017**, *18*, 1–26. [CrossRef] [PubMed]

42. Schindelin, J.; Arganda-Carreras, I.; Frise, E.; Kaynig, V.; Longair, M.; Pietzsch, T.; Preibisch, S.; Rueden, C.; Saalfeld, S.; Schmid, B.; et al. Fiji: An open-source platform for biological-image analysis. *Nat. Methods* **2012**, *9*, 676–682. [CrossRef] [PubMed]

43. Cheng, H.D.; Jiang, X.H.; Sun, Y.; Wang, J. Color image segmentation: Advances and prospects. *Pattern Recognit.* **2001**, *34*, 2259–2281. [CrossRef]

44.  Casadesus, J.; Kaya, Y.; Bort, J.; Nachit, M.M.; Araus, J.L.; Amor, S.; Ferrazzano, G.; Maalouf, F. Using vegetation indices derived from conventional digital cameras as selection criteria for wheat breeding in water-limited environments. *Ann. Appl. Biol.* **2007**, *150*, 227–236. [CrossRef]

45.  Zaman-Allah, M.; Vergara, O.; Araus, J.L.; Tarekegne, A.; Magorokosho, C.; Tejada, P.J.Z.; Hornero, A. Unmanned aerial platform-based multi-spectral imaging for field phenotyping of maize. *Plant Methods* **2015**, *11*, 35. [CrossRef] [PubMed]

46.  Sarkar, S.; Oakes, J.; Cazenave, A.-B.; Burow, M.D.; Bennett, R.S.; Chamberlin, K.D.; Wang, N.; White, M.; Payton, P.; Mahan, J.; et al. Evaluation of the US peanut germplasm mini-core collection in the Virginia-Carolina region using traditional and new high-throughput methods. *Agronomy* **2022**, *12*, 1945. [CrossRef]

47.  R Core Team. *R: A Language and Environment for Statistical Computing*; R Foundation for Statistical Computing: Vienna, Austria, 2023.

48.  Max, A.; Wing, J.; Weston, S.; Williams, A.; Keefer, C.; Engelhardt, A.; Cooper, T.; Mayer, Z.; Ziem, A.; Scrucca, L.; et al. *Package 'caret' R Topics Documented*. 2021. Available online: https://cran.r-project.org/web/packages/caret/caret.pdf (accessed on 22 March 2024).

49.  Taunk, K.; De, S.; Verma, S.; Swetapadma, A. A brief review of nearest neighbor algorithm for learning and classification. In Proceedings of the 2019 International Conference on Intelligent Computing and Control Systems (ICCS), Madurai, India, 15–17 May 2019; pp. 1255–1260.

50.  Chen, S.; Webb, G.I.; Liu, L.; Ma, X. A novel selective naïve Bayes algorithm. *Knowl.-Based Syst.* **2020**, *192*, 105361. [CrossRef]

51.  Xanthopoulos, P.; Pardalos, P.M.; Trafalis, T.B.; Xanthopoulos, P.; Pardalos, P.M.; Trafalis, T.B. Linear discriminant analysis. In *Robust Data Mining*; Springer: Berlin/Heidelberg, Germany, 2013; pp. 27–33. [CrossRef]

52.  Jeppesen, J.H.; Jacobsen, R.H.; Inceoglu, F.; Toftegaard, T.S. A cloud detection algorithm for satellite imagery based on deep learning. *Remote Sens. Environ.* **2019**, *229*, 247–259. [CrossRef]

53.  Butler, D. Weather Requirements for Infection by Late Leaf Spot in Groundnut. International Crops Research Institute for the Semi-Arid Tropics: Patancheruvu, India, 1990.

54.  Wadia, K.; Butler, D. Relationship between temperature and latent periods of rust and leaf-spot diseases of groundnut. *Plant Pathol.* **1994**, *43*, 121–129. [CrossRef]

55.  Okello, D.K.; Deom, C.M.; Puppala, N.; Monyo, E.; Bravo-Ureta, B. Registration of 'Serenut 5R' Groundnut. *J. Plant Regist.* **2016**, *10*, 115–118. [CrossRef]

56.  Okello, D.K.; Deom, C.M.; Puppala, N. Registration of 'Naronut 1R' groundnut. *J. Plant Regist.* **2023**, *17*, 40–46. [CrossRef]

57.  Amare, K.; Seltene, A.; Daniel, E.; Jemal, A.; Addisu, G.; Aliyi, R.; Yohanese, P. East African Journal of Sciences (2017) Registration of 'Babile-1', 'Babile-2', and 'Babile-3' Groundnut Varieties 2. Agronomic and Morphological. *East Afr. J. Sci.* **2017**, *11*, 59–64.

58.  Kavzoglu, T. Environmental Modelling & Software Increasing the accuracy of neural network classification using refined training data. *Environ. Model. Softw.* **2009**, *24*, 850–858. [CrossRef]

59.  Ghorbani, M.A.; Khatibi, R.; Hosseini, B.; Bilgili, M. Relative importance of parameters affecting wind speed prediction using artificial neural networks. *Theor. Appl. Climatol.* **2013**, *114*, 107–114. [CrossRef]

60.  Yoosefzadeh-najafabadi, M.; Earl, H.J.; Tulpan, D.; Sulik, J.; Eskandari, M. Application of Machine Learning Algorithms in Plant Breeding: Predicting Yield from Hyperspectral Reflectance in Soybean. *Front. Plant Sci.* **2021**, *11*, 624273. [CrossRef] [PubMed]

61.  Hesami, M.; Condori-apfata, J.A.; Valencia, M.V.; Mohammadi, M. Application of Artificial Neural Network for Modeling and Studying In Vitro Genotype-Independent Shoot Regeneration in Wheat. *Appl. Sci.* **2020**, *10*, 5370. [CrossRef]

62.  Etminan, A.; Pour-Aboughadareh, A.; Mohammadi, R.; Shooshtari, L.; Yousefiazarkhanian, M.; Moradkhani, H. Determining the best drought tolerance indices using artificial neural network (ANN): Insight into application of intelligent agriculture in agronomy and plant breeding. *Cereal Res. Commun.* **2019**, *47*, 170–181. [CrossRef]

63.  Wang, X.; Zhang, M.; Zhu, J.; Geng, S. Spectral prediction of Phytophthora infestans infection on tomatoes using artificial neural network (ANN). *Int. J. Remote Sens.* **2008**, *29*, 1693–1706. [CrossRef]

64.  Wu, D.; Feng, L.; Zhang, C.; He, Y. Early detection of Botrytis cinerea on eggplant leaves based on visible and near-infrared spectroscopy. *Trans. ASABE* **2008**, *51*, 1133–1139. [CrossRef]

65.  Huang, K.Y. Application of artificial neural network for detecting Phalaenopsis seedling diseases using color and texture features. *Comput. Electron. Agric.* **2007**, *57*, 3–11. [CrossRef]

66.  Cabrera-Bosquet, L.; Crossa, J.; Von Zitzewitz, J.; Serret, M.D.; Araus, L. High-throughput Phenotyping and Genomic Selection: The Frontiers of Crop Breeding Converge Genomic Selection: A Step Forward from. *J. Integr. Plant Biol.* **2012**, *54*, 312–320. [CrossRef] [PubMed]

67.  Gann, D.; Richards, J. Scaling of classification systems—Effects of class precision on detection accuracy from medium resolution multispectral data. *Landsc. Ecol.* **2023**, *38*, 659–687. [CrossRef]

68.  Dwivedi, S.L.; Goldman, I.; Ceccarelli, S.; Ortiz, R. *Advanced Analytics, Phenomics and Biotechnology Approaches to Enhance Genetic Gains in Plant Breeding*, 1st ed.; Elsevier Inc.: Amsterdam, The Netherlands, 2020; Volume 162.

69.  Baranowski, P.; Jedryczka, M.; Mazurek, W.; Babula-Skowronska, D.; Siedliska, A.; Kaczmarek, J. Hyperspectral and thermal imaging of oilseed rape (Brassica napus) response to fungal species of the genus Alternaria. *PLoS ONE* **2015**, *10*, e0122913. [CrossRef]

70.  Govardhan, M.; Veena, M.B. Diagnosis of Tomato Plant Diseases using Random Forest. In Proceedings of the 2019 Global Conference for Advancement in Technology (GCAT), Bangaluru, India, 18–20 October 2019; pp. 1–5. [CrossRef]

71. Wójtowicz, A.; Piekarczyk, J.; Czernecki, B.; Ratajkiewicz, H. A random forest model for the classification of wheat and rye leaf rust symptoms based on pure spectra at leaf scale. *J. Photochem. Photobiol. B Biol.* **2021**, *223*, 112278. [CrossRef] [PubMed]
72. Caelen, O. A Bayesian interpretation of the confusion matrix. *Ann. Math. Artif. Intell.* **2017**, *81*, 429–450. [CrossRef]
73. Oteng-Frimpong, R.; Karikari, B.; Sie, E.K.; Kassim, Y.B.; Puozaa, D.K.; Rasheed, M.A.; Fonceka, D.; Okello, D.K.; Balota, M.; Burow, M.; et al. Multi-locus genome-wide association studies reveal genomic regions and putative candidate genes associated with leaf spot diseases in African groundnut (*Arachis hypogaea* L.) germplasm. *Front. Plant Sci.* **2023**, *13*, 1076744. [CrossRef]
74. Liu, H.Q.; Huete, A. Feedback based modification of the NDVI to minimize canopy background and atmospheric noise. *IEEE Trans. Geosci. Remote Sens.* **1995**, *33*, 457–465. [CrossRef]
75. Kassim, Y.B.; Oteng-Frimpong, R.; Puozaa, D.K.; Sie, E.K.; Abdul Rasheed, M.; Abdul Rashid, I.; Danquah, A.; Akogo, D.A.; Rhoads, J.; Hoisington, D.; et al. High-Throughput Plant Phenotyping (HTPP) in Resource-Constrained Research Programs: A Working Example in Ghana. *Agronomy* **2022**, *12*, 2733. [CrossRef]
76. Crain, J.L.; Wei, Y.; Barker, J.; Thompson, S.M.; Alderman, P.D.; Reynolds, M.; Zhang, N.; Poland, J. Development and Deployment of a Portable Field Phenotyping Platform. *Crop Sci.* **2016**, *56*, 965. [CrossRef]
77. Andrade-Sanchez, P.; Gore, M.A.; Heun, J.T.; Thorp, K.R.; Carmo-Silva, A.E.; French, A.N.; Salvucci, M.E.; White, J.W. Development and evaluation of a field-based high-throughput phenotyping platform. *Funct. Plant Biol.* **2014**, *41*, 68–79. [CrossRef] [PubMed]
78. Lane, H.M.; Murray, S.C. High throughput can produce better decisions than high accuracy when phenotyping plant populations. *Crop Sci.* **2021**, *61*, 3301–3313. [CrossRef]
79. Bali, N.; Singla, A. Emerging Trends in Machine Learning to Predict Crop Yield and Study Its Influential Factors: A Survey. *Arch. Comput. Methods Eng.* **2022**, *29*, 95–112. [CrossRef]
80. Zhou, Z.-H. *Ensemble Learning BT—Encyclopedia of Biometrics*; Li, S.Z., Jain, A., Eds.; Springer: Boston, MA, USA, 2009; pp. 270–273.
81. Feng, L.; Zhang, Z.; Ma, Y.; Du, Q.; Williams, P.; Drewry, J.; Luck, B. Alfalfa yield prediction using UAV-based hyperspectral imagery and ensemble learning. *Remote Sens.* **2020**, *12*, 2028. [CrossRef]
82. Vasilakos, C.; Kavroudakis, D.; Georganta, A. Machine Learning Classification Ensemble of Multitemporal Sentinel-2 Images: The Case of a Mixed Mediterranean Ecosystem. *Remote Sens.* **2020**, *12*, 2005. [CrossRef]