# Precise Sweetness Grading of Mangoes (*Mangifera indica* L.) Based on Random Forest Technique With Low-Cost Multispectral Sensors

**CHANH-NGHIEM NGUYEN**[1], **(Member, IEEE), QUOC-THANG PHAN**[1], **NHUT-THANH TRAN**[1,2], **MASAYUKI FUKUZAWA**[2], **(Member, IEEE), PHUOC-LOC NGUYEN**[1,3], **AND CHI-NGON NGUYEN**[1]

[1]Department of Automation Technology, Can Tho University, Can Tho 94000, Vietnam
[2]Graduate School of Science and Technology, Kyoto Institute of Technology, Kyoto 606-8585, Japan
[3]Faculty of Electric-Electronics and Computer, Kien Giang Vocational College, Rach Gia 91000, Vietnam

Corresponding author: Chanh-Nghiem Nguyen (ncnghiem@ctu.edu.vn)

**ABSTRACT** Before being exported, mangoes generally undergo rigorous external and internal quality inspection processes in which near-infrared (NIR) spectral approaches are favorable for grading purposes. A successful NIR-based grading system depends largely on high-quality spectral sensors and the reliability of the classifier. Motivated by the high economic impact of Cat Hoa Loc mangoes (Mangifera indica L.), we demonstrated that the sweetness of that mangoes could be precisely graded based on a random forest (RF) classifier in a three-phase approach with a low-cost Visible-Near infrared (VIS-NIR) multispectral sensor chipset. This approach is so-called RPR because RF, Partial Least Squares regression, and RF were respectively applied to consecutively determine the significant VIS-NIR responses, the good features as input variables, and the reliable RF classifier via our formulated discriminant index (DI). The experimental results confirmed that higher classification accuracy was achieved by using the extracted latent features rather than the raw VIS-NIR data. The DI was effectively used as a reliability measure to select the optimal classifier among those of identical training and testing accuracies of 100% and 82.1%, respectively. Performance comparison between the optimal RF classifier with a Support Vector Machines classifier and a multinomial logistic regression showed that the developed RF classifier was superior in various performance indices. Therefore, it is promising to extend the proposed approach to more complicated fruit grading problems with sufficient VIS-NIR datasets that are acquired from low-cost multispectral sensors.

**INDEX TERMS** Cost-effective, sweetness grading of Cat Hoa Loc mango, spectral response selection, VIS-NIR features extraction, discriminant index for random forest classifier.

## I. INTRODUCTION

Historically, Mangifera indica L. cultivations have been widely planted in tropical areas of India, Africa, Asia, and Central America. Due to its richness in nutrients and minerals, and its good taste and aroma, many Mangifera indica L. cultivars have been developed in favorable subtropical climates and adapted soils [1]. With an increasing global import demand, mangoes are becoming an export fruit of very high economic value for many countries in Africa and Southeast Asia [2]. Therefore, many techniques have been developed for non-invasive inspection and grading of high-quality mangoes.

One of the simple approaches was to utilize the visual or extrinsic features of the mango such as color, pattern, size, and shapes, etc. Based on the fact that the sugar content increased as the color changed during the maturation process, the hue feature extracted from an RGB image of a Chokanan mango (Mangifera indica) was used to grade the sweetness of Chokanan mangoes with up to an average success rate of 95.67% [3]. In a similar approach, multivariate discriminant analysis was applied to obtain a maturity classification model with a classification rate of 90% for ''Manila'' mangoes using more color information (i.e., a*, b*, S, and H color

The associate editor coordinating the review of this manuscript and approving it for publication was Juan Wang.

coordinates) [4]. Although good discriminant results were reported, it should be noted that the skin color feature of a mango fruit might not be robust indicators of its maturity as different environments might cause the skin color to change [5], [6]. More sophisticated data fusion methods were also reported for maturity and ripeness classification. These methods applied data fusion using various sensory information such as odor, acoustic response [7]; physical, mechanical, and optical properties of the mango [8]. Classification accuracy could also be enhanced with machine learning-based techniques using convolutional neural networks with sophisticated multispectral information fusion [9], [10]. The neural network approach was reported to have a high-accuracy prediction of the sweetness of a mango using its physical features including its mass, length, width, and volume [11]. However, this might not be a stable approach because little internal features were utilized. Moreover, convolutional neural networks require high computation power for real-time implementation. It is also complicated to implement various kinds of sensors and data fusion methods for an automated fruit sorting system.

To provide more internal quality information of the fruit, spectral data were utilized. With the ease of implementation, few studies took advantage of the spectrum within the range of 400-867 nm [12], [13]. However, the application of near-infrared spectroscopy (NIRS) to the horticultural field has attracted considerable attention. According to the recent review, 80% of the NIRS-based horticultural applications were for investigations and studies on fruit, among which about 13% were carried out on mangoes which were the fruit of the second most interest just after apples [14].

Because near-infrared (NIR) spectra contain information about the major C–H, O–H, and N–H bonds, NIR spectra might reveal valuable information of many organic materials that could be used for internal fruit quality analysis [15]. Many studies have reported the use of NIR spectra to estimate dry matter (DM) content to determine the mango maturity [6], [16]. Classification and prediction problems related to the mango sweetness were also tackled. Therefore, the shortwave infrared (SWIR) spectrum in the range of 900-2400 nm was investigated [17]–[19] because there are major absorption bands of sugar within that spectrum [20]. For cost-efficient applications, the visible and near-infrared (VIS-NIR) spectrum of 500-1100 nm was also adopted [6], [16], [21]–[24]. The utilization of this spectral region could be mainly due to various C–H and O–H vibrations reported for glucose and sucrose; and the availability of commercially available, somewhat low-cost, miniaturized NIR spectrometers for this specific spectral zone [21], [23], [25].

It should be noted that the ripening process is not only characterized by the increase of sugar content but also by the degradation of green chlorophylls and the accumulation of colored pigments such as carotenoids and anthocyanins which account for the yellowish color of the flesh and the reddish color of the peel, respectively [26], [27]. It is thus very probable that the visible spectrum also contains useful information for grading mangoes based on their sweetness because a correlation might exist between the contents of sugar and those colored pigments. Specifically, the low visible spectral zone of 400-500 nm is the broad absorption band of isolated yellow carotenoids. Within this spectral zone, isolated chlorophyll (Chl) *a* and Chl *b* also exhibit the narrow absorption bands (maxima) near 428 and 453 nm, respectively [28]. Since the ripening process has a complex effect on the NIR spectra, the machine learning-based classifier using NIR spectra features could be favorable to improve the classification accuracy. It would be a more stable approach than the neural network approaches using visual features mentioned above because it directly learns spectral data that contains huge information on the internal quality.

Recently, the VIS spectrum can be obtained easily by low-cost sensors, and some low-cost devices to obtain the NIR spectrum such as integrated multispectral sensors became available. A low-cost multispectral chipset could be used to develop a portable spectrometric system to predict the soluble solids content of *Kiou* apples [29]. With a simple optical setup, the proposed system showed much potential for practical applications in terms of manufacturability and reproducibility. It also suggested that cost-effective multispectral sensors could provide useful features from a few significant wavebands for the assessment of internal fruit quality. It is thus worthwhile to examine the potential of such low-cost multispectral sensors in the noninvasive assessment of fruit quality because successes in applying such sensors can facilitate the development of commercial grading applications not only for mango exporting companies but also for mango consumers. In this case, an effective grading technique is very important to make full use of the VIS-NIR data acquired from such sensors.

As an initiative for noninvasive mango grading with low-cost VIS-NIR multispectral sensors, this study aimed to develop an effective classification technique for precise sweetness grading of mangoes (Mangifera indica L. cv. Cat Hoa Loc) based on the spectral data acquired from a low-cost multispectral sensor chipset with a range from 410 nm to 940 nm. To achieve such a goal, a systematic approach for developing a reliable random forest classifier was proposed with specific contributions as follows.

1) Significant information (i.e., spectral responses) could be extracted from the spectral data.
2) Classification accuracy could be enhanced by using the latent features extracted from significant spectral responses.
3) The proposed Discriminant Index (DI) was a simple and effective performance measure for determining the optimal random forest (RF) classifier.

## II. MATERIAL AND METHODS
### A. INSTRUMENTATION
In this study, AS7265x Smart Spectral Sensor (ams AG) chipset was utilized. This evaluation board includes three
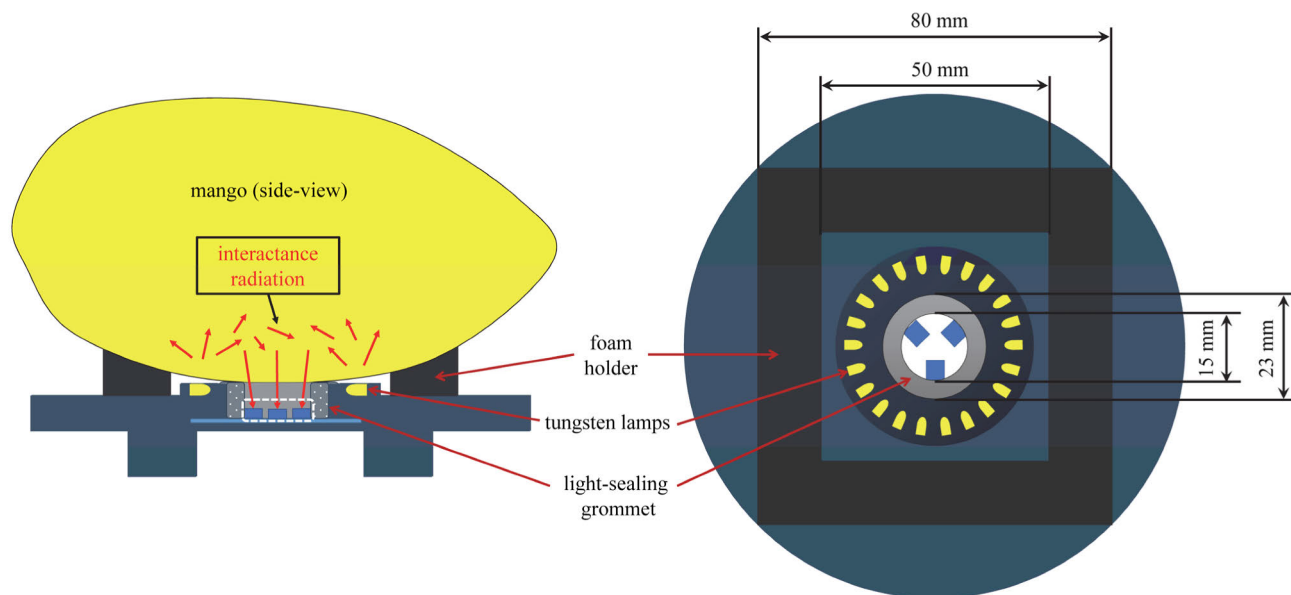
**FIGURE 1.** Design of adapter for AS7265x kit for interaction mode application.

6-channel sensor devices, delivering 18 VIS and NIR channels from 410 nm to 940 nm, each with 20 nm Full Width at Half Maximum (FWHM) [30]. For most NIRS applications, the measurement mode plays a very important role. The transmittance mode is very well known for detecting both external and internal qualities of a sample. Nevertheless, transmission measurements are generally difficult to obtain. Thus, the interactance mode was applied because it provided a compromise between the reflectance and transmittance modes [31] and it was a convenient method to obtain internal mango qualities in various studies [6], [16], [18], [22]–[24], [32].

Since the interactance mode was utilized, a simple optical setup was designed to utilize the characteristics of the AS7265x chipset such that appropriate lighting could be provided without affecting the measurement area. As illustrated in Fig. 1, the AS7265x kit was mounted at the center of the adapter, surrounded by a grommet to block the light from a ring of 22 tungsten lamps (Mineshima P-23, average dissipation power of 0.5 W/bulb) that ensured the illumination over the entire spectral region of the sensors. The outer area of the adapter had a foam holder to provide a soft contact with the mango and to prevent the illuminated area from the ambient light.

### B. MANGO DATA ACQUISITION

As mangoes (Mangifera indica L. cv. Cat Hoa Loc) can be harvested throughout the year, a reliable model for grading "Cat Hoa Loc" mangoes should be developed with consideration of the variability due to the harvest time. Therefore, samples of cv. "Cat Hoa Loc" were collected at different harvest time at the local fruit stores. All mangoes originated from three provinces in Vietnam, namely

**TABLE 1.** Mango dataset information.

| Month of harvest | Total samples | Number of samples by grade | | | | | |
| | | Training dataset | | | Testing dataset | | |
| | | I | II | III | I | II | III |
|---|---|---|---|---|---|---|---|
| Jan | 28 | 1 | 12 | 15 | 0 | 0 | 0 |
| Mar | 12 | 0 | 0 | 6 | 0 | 0 | 6 |
| April | 10 | 0 | 2 | 3 | 0 | 2 | 3 |
| Jun | 10 | 4 | 1 | 0 | 3 | 2 | 0 |
| Sep | 18 | 0 | 6 | 3 | 0 | 5 | 4 |
| Oct | 16 | 0 | 3 | 5 | 0 | 3 | 5 |
| Dec | 12 | 0 | 0 | 6 | 0 | 0 | 6 |
| Subtotal | | 5 | 24 | 38 | 3 | 12 | 24 |
| Total | 106 | 67 (63.2%) | | | 39 (36.8%) | | |

Can Tho, Tien Giang, and Dong Thap provinces. To help maintain the representative and robustness of the model, the samples were split into training and testing datasets such that most samples, for each harvest period, were present in both training and testing datasets as summarized in Table 1.

Because sugars constitute the majority of total soluble solids (TSS) in many fruits [33], the mango sweetness was determined based on the "degrees Brix," a conventional measure of the TSS present in the fruit.

According to the quality guide for Cat Hoa Loc mangoes [5], they could be classified into three categories of highly acceptable, acceptable, and unacceptable based on their Brix values. Based on this suggestion, all mango samples were graded based on their Brix values as

$$grade = \begin{cases} \text{I}, & \text{if } Brix > 24 \\ \text{II}, & \text{if } 20 \leq Brix \leq 24 \\ \text{III}, & \text{if otherwise} \end{cases} \quad (1)$$

### 1) MULTISPECTRAL DATA ACQUISITION

To ensure correct measurement, preliminary calibration for AS7265x sensors had been performed. First, the sensors were completely covered and spectral intensities of all channels were confirmed to be zeros, indicating that the sensors did not include any bias. Next, a white calibration plate (X-rite Inc.) was applied on the adapter to acquire the $i$-th reflectance spectral response $I_{i0}$ so that the normalized interactance response could be calculated as

$$\hat{I}_i = \frac{I_i}{I_{i0}} \qquad (2)$$

where $I_i$ is the interactance response of a mango sample for the $i$-th spectral channel with the respective wavelength $\lambda_i$ of 410, 435, 460, 485, 510, 535, 560, 585, 610, 645, 680, 705, 730, 760, 810, 860, 900, 940 nm; and $i = 1, 2, \cdots, 18$.
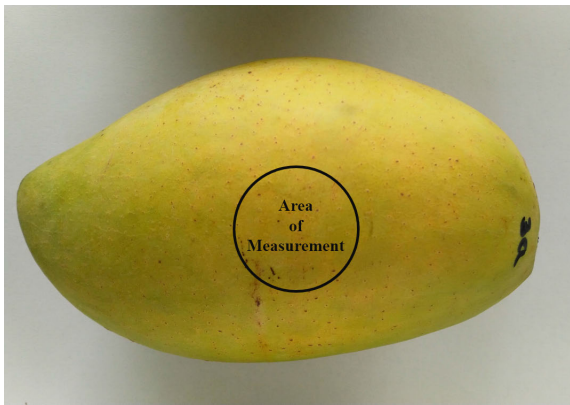


**FIGURE 2.** Area of measuring VIS-NIR interactance response on a mango.

Fig. 2 shows the actual area of measurement on one side of a sample mango where multispectral data were acquired. Both sides of the mangoes were subjected to multispectral data acquisition to increase the data size. To minimize the measurement noise, the averaged measurement $\bar{I}_i$ from five consecutive measurements of $\hat{I}_i$ was used for developing the classification model.

### 2) BRIX MEASUREMENT

After acquiring multispectral data of a mango sample, a piece of the mango flesh of about 10 millimeters deep was pulled out from the area of measurement. After it had been peeled, its juice was extracted, filtered, and subjected to MA871 digital Brix refractometer (Milwaukee Instruments) for Brix measurement. A mango sample could be graded from the Brix measure of the flesh by using (1). A number of 63.2% sample data were used in all phases of developing the classification model whereas the remaining were solely used for testing the model (Table 1). The statistics of Brix measurements are shown in Table 2. Few samples were found as "grade I" because high-quality mangoes were normally selected for export rather than for local consumption.

**TABLE 2.** Statistics of Brix measurements.

| Statistic | Grade I | Grade II | Grade III |
|---|---|---|---|
| Number of samples | 8 | 36 | 62 |
| Minimum (°Brix) | 24.2 | 20.0 | 12.4 |
| Maximum (°Brix) | 27.2 | 24.0 | 19.9 |
| Mean (°Brix) | 25.4 | 21.4 | 17.1 |
| Standard deviation (°Brix) | 1.1 | 1.2 | 2.0 |

### C. DEVELOPMENT OF CLASSIFIER FOR MANGO GRADING

In the process of developing a reliable classifier or a regressor, it is very important to select significant input variables and features so that overfitting is prevented, and reasonable generalization is allowed. Since all spectral responses in the investigated range of wavelength might not contribute equally to the classifier, only significant responses were selected so that discriminant features that highly correlated with the target grades were extracted. A grid search was then performed for $N$ iterations to generate $N$ high-performance random forest (RF) classifiers from which the optimal one could be determined based on a discriminant index. Thus, this approach included three phases in which RF, partial least squares (PLS), and RF were respectively applied in a consecutive order to determine the significant spectral responses, extract the good features, and identify the optimal RF classifier. This so-called RPR approach is summarized in Fig. 3 and described in detail as follows.

### 1) VARIABLE SELECTION BASED ON RANDOM FOREST

Random forest algorithm has been extremely successful for general-purpose classification and regression. It is a kind of ensemble learning technique that generated several randomized decision trees and aggregates their learning results by averaging [34]. Random forests can be used to rank the importance of variables in both regression or classification problems via two measures of variable importance, which are Mean Decrease Impurity (MDI) and Mean Decrease Accuracy (MDA). A detailed explanation and theoretical background of these measures were described elsewhere [34], [35]. In this study, the normalized MDI measure, calculated by using scikit-learn module [36], was used to determine the classification impact of the input variables so that significant variables could be selected for mango grading.

Let define $MDI_i$ the normalized MDI of the interactance response $\bar{I}_i$ at wavelength $\lambda_i$ where $i$ denotes the index of VIS and NIR channels that AS7265x smart spectral sensor can deliver and $i = 1, 2, \cdots, 18$. Let $v_i$ be the vote value for an input variable (i.e., interactance response $\bar{I}_i$) calculated from a certain RF model. An input variable would be considered important if it receives a vote value of 1; otherwise 0. The vote $v_i$ was determined from the corresponding $MDI_i$ as

$$v_i = \begin{cases} 1, & \text{if } MDI_i > \frac{1}{18} \\ 0, & \text{if otherwise} \end{cases} \qquad (3)$$
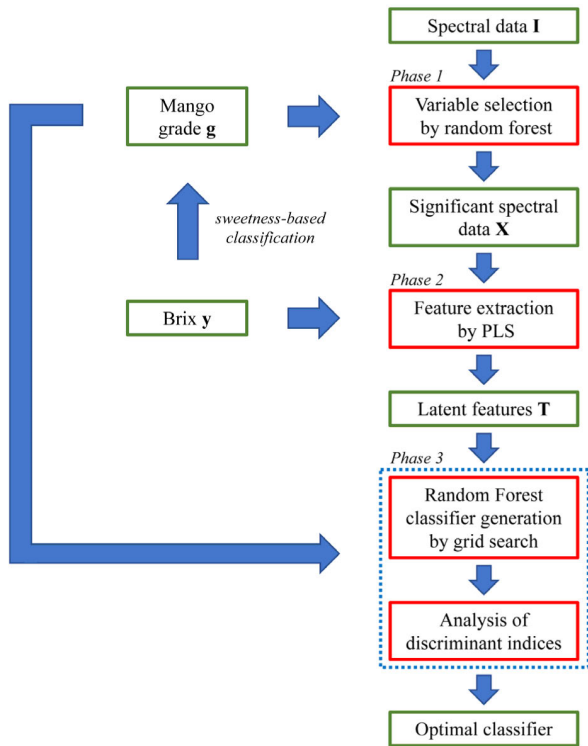
**FIGURE 3.** Illustration of the proposed RPR approach.

where 1/18 was the mean of normalized MDI values of all 18 interactance responses.

To obtain a confident decision on the significance of $\bar{I}_i$, a number of $N = 100$ iterations of grid search were carried out to generate 100 random forests. Thus, the total number of votes each variable $\bar{I}_i$ received was calculated as

$$V_i = \sum_{j=1}^{N} v_i^j, \tag{4}$$

where $v_i^j$ is the vote value of the $j$-th RF model for input variable $\bar{I}_i$. An input variable $\bar{I}_i$ was considered important for mango grading if it met the following criteria

$$V_i > \frac{N}{2} = 50. \tag{5}$$

### 2) FEATURE EXTRACTION BASED ON PARTIAL LEAST SQUARE REGRESSION

Partial least square (PLS) regression is a powerful tool for developing a regression model. The goal of PLS analysis is to search for a set of components (called *latent vectors*) from **X** for the best prediction of **Y** by performing a simultaneous decomposition of **X** and **Y** with the constraint that the components explain as much as possible the covariance between **X** and **Y**.

The independent variables are decomposed as

$$\mathbf{X} = \mathbf{TP}^T \quad \text{with } \mathbf{T}^T\mathbf{T} = \mathbf{I}, \tag{6}$$

where **T** and **P** are the score matrix and the loading matrix, respectively. The columns of **T** are the latent vectors. **Y** is estimated as

$$\hat{\mathbf{Y}} = \mathbf{TBC}^T, \tag{7}$$

where **B** is the diagonal matrix whose diagonal elements are the "regression weights," and **C** is the "weight matrix" of the dependent variables. The main target of PLS regression is to specify **T** by obtaining iteratively all pairs of vectors

$$\mathbf{t} = \mathbf{Xw} \quad \text{and } \mathbf{u} = \mathbf{Yc}, \tag{8}$$

with the constraint that $\mathbf{w}^T\mathbf{w} = 1$, $\mathbf{t}^T\mathbf{t} = 1$, and importantly their covariance $\mathbf{t}^T\mathbf{u}$ is maximal [37]. Further explanation of PLS regression could be found elsewhere [37]–[39].

The latent components are then used for prediction in place of the original variables. However, it is necessary to determine the optimal number of latent variables to keep for building the PLS model. Typically, cross-validation is the most popular method for PLS model selection and data overfitting can be detected if increasing the number of latent variables leads to a decrease in the prediction accuracy.

Let $\bar{\mathbf{I}}$ ($n$ x 18) be the spectral data of the 18 interactance responses. Let $\mathbf{X}(n$ x $m)$ be the extracted version of $\bar{\mathbf{I}}$ so that it consisted of the responses of $m$ spectral interaction signals that were optimally selected based on their impact using (5). Let $k \leq m$ be the optimal number of latent variables that were obtained by PLS regression from **X** and corresponding Brix **y** ($n$ x 1) of the mango samples in the study. These latent variables (LV) should have a maximal covariance with the univariate response **y**.

To determine the optimal LVs, leave-one-out cross-validation was performed with different numbers of LVs, as suggested by many studies [40], [41]. Because the coefficient of determination indicates the goodness of fit for the observations, the optimal LVs were associated with the largest coefficient of determination after leave-one-out cross-validation. Using scikit-learn module [36], the coefficient of determination, usually denoted as $R^2$ was computed as [42]

$$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n} (y_i - \bar{y})^2}, \tag{9}$$

where

$$\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i, \tag{10}$$

and $\hat{y}_i$ is the predicted value of the $i$-th sample; $y_i$ is the corresponding true value for the total $n$ samples.

The latent features were used as inputs for developing the mango classifier. The latent-feature data $\mathbf{T}(n$ x $k$ matrix) could be calculated as

$$\mathbf{T} = \mathbf{XW}, \tag{11}$$

where **W** is the $m$ x $k$ transformation matrix that could be obtained by horizontally concatenating the corresponding column vector **w** that had been derived from **t** using (8).
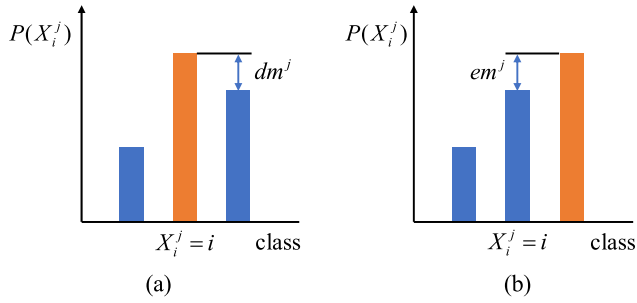
**FIGURE 4.** Illustration of (a) probability discriminant margin for a correct classification case and (b) probability error margin for a wrong classification case. Orange color denotes the predicted result.

### 3) DEVELOPMENT OF RANDOM FOREST CLASSIFIER

To obtain a good classifier within a certain parameter space, a grid search with cross-validation was iterated 100 times, each generated a candidate model. Then, all candidate models were evaluated with the testing dataset such that a smaller set of these candidate models with high classification accuracy could be shortlisted. The classification accuracy of a model was calculated as

$$\text{accuracy}\left(\hat{\mathbf{g}}, \mathbf{g}\right) = \frac{1}{n} \sum_{l=0}^{n-1} 1\left(\hat{g}_l = g_l\right), \tag{12}$$

where $1(x)$ is the indicator function; $\hat{g}_l$ and $g_l$ are the predicted and the true grade of the $l$-th sample, respectively; $n$ is the number of samples of the dataset.

Although the reliability of different classifiers could be analyzed statistically with two-dimensional reliability diagrams [43], a simple approach was proposed in this study based on the class probabilities of a random forest classifier that could be inherently obtained.

Let denote $X_i^j$ the predicted class for sample $j$ which belongs to class $i$. Let denote $P(X_i^j = k)$ the probability of predicting sample $j$ into class $k$. It is obvious that

$$\sum_k P(X_i^j = k) = 1. \tag{13}$$

Sample $j$ is correctly predicted when $X_i^j = i$, which means

$$P(X_i^j = i) = \max_k(P(X_i^j = k)). \tag{14}$$

For the correct classification of sample $j$, let denote $dm^j \in (0, 1]$ the probability discriminant margin which is calculated as

$$dm^j = P(X_i^j = i) - \max_{k \neq i}(P(X_i^j = k)). \tag{15}$$

As shown in Fig. 4a, a greater value of $dm^j$ indicates greater certainty of the $j$-th classification result, and the most confident result is obtained when

$$\begin{cases} P(X_i^j = i) = 1 \\ \underset{i \neq k}{P}(X_i^j = k) = 0 \\ dm^j = 1. \end{cases} \tag{16}$$

For the wrong classification of sample $j$, let calculate the probability error margin $em^j \in (0, 1]$ as

$$em^j = \max(P(X_i^j = k)) - P(X_i^j = i). \tag{17}$$

A greater value of $em^j$ indicates a greater undesired certainty of the false $j$-th classification result as shown in Fig. 4b. For a good classification model, it is desirable to obtain a large probability discriminant margin and a small probability error margin for all samples. These margins can be incorporated into a discriminant index (DI), defined as

$$DI = \frac{1}{N_c + N_w}\left(\sum_{j=1}^{N_c} dm^j - \sum_{j=1}^{N_w} de^j\right), \tag{18}$$

where $N_c$ and $N_w$ are the numbers of correct and wrong classification cases, respectively. DI can also be written as

$$DI = \frac{N_c}{N_c + N_w} \frac{1}{N_c} \sum_{j=1}^{N_c} dm^j - \frac{N_w}{N_c + N_w} \frac{1}{N_w} \sum_{j=1}^{N_w} de^j, \tag{19}$$

or

$$DI = \frac{N_c}{N_c + N_w} ADM - \frac{N_w}{N_c + N_w} AEM, \tag{20}$$

where *ADM* and *AEM* are respectively the average probability discriminant margin and the average probability error margin. Therefore, DI can be rewritten based on the accuracy of the classification model as

$$DI = accuracy * ADM - (1 - accuracy) * AEM. \tag{21}$$

It is obvious that DI falls in the range of $[-1, 1]$; however, DI does not probably receive the theoretical minimum or maximum. The formulation of DI showed that DI was a more effective measure for relative comparison between RF classifiers than any metrics solely based on the performance of an RF classifier. Because $P(X_i^j = k)$ could be obtained as the mean probability estimate across the trees in the forest $j$ when predicting the class for the sample $i$, the training and testing DI could be calculated with training and testing dataset, respectively. Therefore, these DI values of all model candidates could give some guidelines for the determination of the best RF classifier. The code for DI was prepared and available for use as-is on GitHub [44].

### D. COMPARISON WITH OTHER CLASSIFIERS

To demonstrate the effectiveness of the proposed RF classifier, a Support Vector Machines (SVM) classifier and a multinomial logistic regression classifier were developed using the same training dataset. The testing accuracies of the developed RF classifier and these classifiers were also compared based on their performance on the same testing dataset.

### 1) DEVELOPING AN SVM CLASSIFIER

Although SVMs were originally designed for binary classification, SVM classifiers have effectively been extended for multiclass classification [45]. SVM is a well-known

non-linear learning algorithm that uses a kernel function to transform the data into higher-dimensional spaces where a hyperplane can be constructed to separate the data with the maximal margin between the desired classes. They have also been intensively coupled with NIR spectroscopy in food analysis [46].

In this study, we developed the SVM classifier using the radial basis function, which is expressed as

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\gamma \left\| \mathbf{x}_i - \mathbf{x}_j \right\|^2\right), \quad (22)$$

where $\mathbf{x}_i$, $\mathbf{x}_j$ are the training vectors; and $\gamma$ is the kernel parameter. SVM training is a constrained optimization problem of

$$\min \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{m}\xi_i, \quad \text{subject to} \quad (23)$$
$$y_i(\langle \mathbf{w}\cdot\mathbf{x}\rangle + b) \geq 1 - \xi_i, \quad (24)$$
$$\xi_i \geq 0, \quad i = 1, 2, \cdots, m, \quad (25)$$

where $\mathbf{w}$ is the weight vector, b is the bias, $\xi_i$ measures the degree of misclassification of datapoint $i$, $y_i$ is the class label for the $i$-th sample, and $C$ is the trade-off parameter between the margin and error [47], [48].

Because an RBF kernel was used for SVM training, the penalty parameter $C$ and the kernel parameter $\gamma$ should be identified beforehand. Thus, we performed a grid-search on $C$ and $\gamma$ using cross-validation following the practical guide recommended in [49].

### 2) DEVELOPING A MULTINOMIAL LOGISTIC REGRESSION CLASSIFIER

Logistic regression (LG) is a statistical method that has been widely used for classification purposes. To provide more precise categorical solutions for multi-class classification problems, multinomial logistic regression (MLR) and many of its extensions have been proposed. A brief explanation and summary of the recent development of MLR could be found in [50].

In this study, we also developed an MLR classifier to compare the effectiveness of the proposed RF classifier. Because MLR supported the probability of the class outputs, the training and testing DI values of the developed MLR were also calculated.

## III. RESULTS AND DISCUSSION

In *Phase 1* of the proposed approach for developing the classification model, 100 random forests were generated to vote for the significance of the 18 interactance responses using (3)-(5). Each RF was the best classification model that was determined from a grid search with parameter configuration listed in Table 3. For our case of small datasets, leave-one-out cross-validation was utilized in the grid search.

Fig. 5 shows the average spectra for grades I, II, and III. Since neither anomalous spectra nor obviously-grade-correlated spectra were found in any grade, the multispectral data acquisition was successfully performed. Table 4 shows

**TABLE 3.** Parameter configuration for grid search.

| Parameter | Values |
|---|---|
| max_depth | 3, 4, 5, None |
| max_features | 'sqrt', 'log2', None |
| n_estimators | 100 |
| oob_score | True |

**TABLE 4.** Vote results for variable importance.

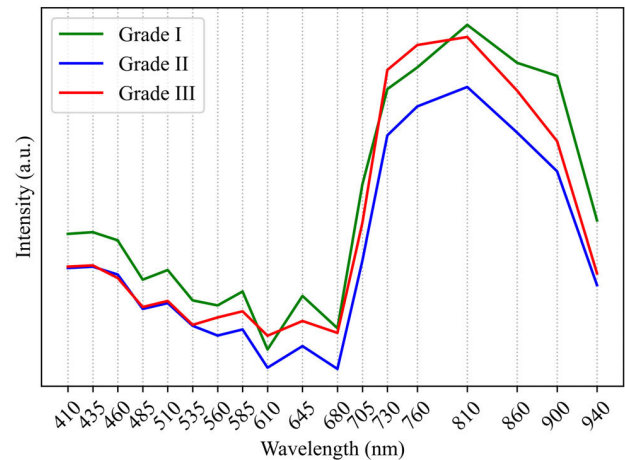| Index $i$ | Wavelength of spectral response $\overline{I}_i$ (nm) | Number of votes | Selected |
|---|---|---|---|
| 1 | 410 | 37 | No |
| 2 | 435 | 97 | Yes |
| 3 | 460 | 92 | Yes |
| 4 | 485 | 2 | No |
| 5 | 510 | 11 | No |
| 6 | 535 | 4 | No |
| 7 | 560 | 13 | No |
| 8 | 585 | 2 | No |
| 9 | 610 | 79 | Yes |
| 10 | 645 | 0 | No |
| 11 | 680 | 32 | No |
| 12 | 705 | 2 | No |
| 13 | 730 | 98 | Yes |
| 14 | 760 | 80 | Yes |
| 15 | 810 | 69 | Yes |
| 16 | 860 | 70 | Yes |
| 17 | 900 | 73 | Yes |
| 18 | 940 | 93 | Yes |



**FIGURE 5.** Average spectra for various mango grades.

the vote results of all the variables. Nine variables were determined as important, which were the interactance responses at wavelengths of 435, 460, 610, 730, 760, 810, 860, 900, and 940 nm. These wavelengths were relatively close to the respective narrow absorption bands of isolated Chl *a* (near 428 nm), and Chl *b* (near 453 nm), and they were in the broad absorption band (between 400 and 500 nm) of isolated yellow carotenoids [28].

For mangoes, their green pigmentation and yellowish color are attributed to the presence of Chls and carotenoids [1], [51], respectively. As the mango ripened, a decrease in Chl content in the mango pulp was found accompanied by

a coordinate biosynthesis and accumulation of carotenoids with a relatively high content [26]. Therefore, the significance of the spectral responses at 435 nm and 460 nm could be confirmed because the mango color is mainly characterized by Chls and carotenoids at different ripening stages, and it is highly correlated with the sugar content according to [27].

Three significant interaction responses at 730, 900, 940 nm were found associated with sugar absorbance bands due to various combinations of major C–H and O–H vibrations [23], [52]–[54]. Two significant responses at 760 and 860 nm were close to the reported wavelengths at 755 and 850 nm that were strongly correlated with organic acids [53] whose content might have been decreased according to the increase in the sugar content during mango ripening [55].

When the important interactance responses had been determined, the significant spectral data $\mathbf{X}$ (67 x 9) was extracted from the normalized spectral data $\mathbf{I}$ (67 x 18). PLS regression was applied with input data $\mathbf{X}$ and Brix data $\mathbf{y}$.
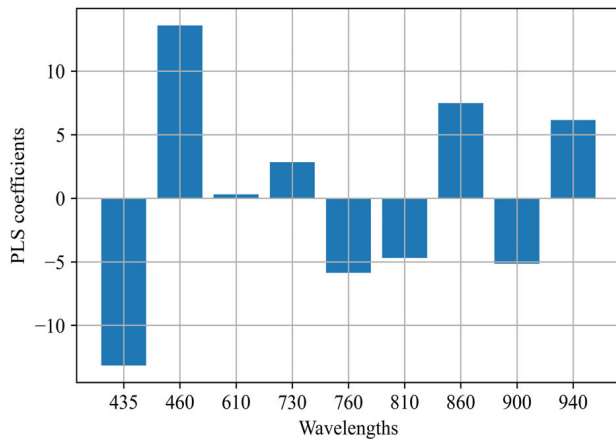
**FIGURE 6.** PLS regression b-coefficients.

Fig. 6 shows the plot of PLS b-coefficients. Comparably large and small coefficients were both observed, which suggested that an optimal subset of latent variables might be obtained. Fig. 7 shows the plot of the coefficient of determination after leave-one-out cross-validation had been performed with different numbers of latent variables. The seven-LV model had the largest coefficient of determination. Thus, the seven optimal LVs were determined. Accordingly, the latent-feature data $\mathbf{T}$ (67 x 7) that best explained for Brix data $\mathbf{y}$ could be obtained.

The latent features $\mathbf{T}$ (i.e., the outcome of *Phase 2* of the proposed approach) and the desired mango grade $\mathbf{g}$ were used to develop the classification model for grading mangoes based on their sweetness index. The grid search was performed for 100 iterations using the same parameter configuration as listed in Table 3. The outcomes of the grid search were 100 candidate models from which 10 models with the highest testing accuracy were shortlisted (Table 5).

It was noted that there were models with identical training and testing accuracy, which showed the necessity of some
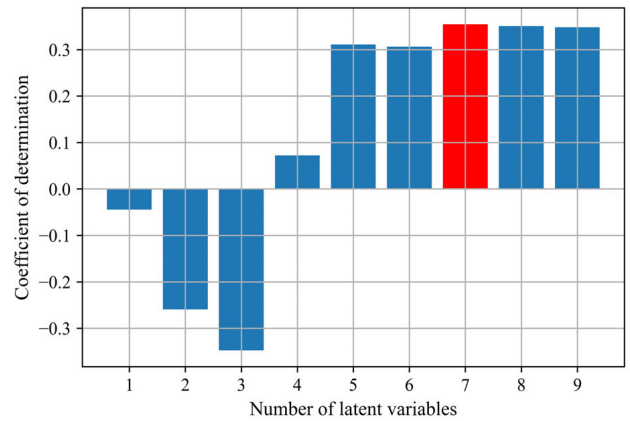
**FIGURE 7.** Coefficient of determination obtained after leave-one-out cross-validation with different numbers of latent variables. Red bar indicated the case of the largest coefficient of determination.

**TABLE 5.** Discriminant indices of candidate models.

| Grid search iteration | Training dataset | | Testing dataset | |
|---|---|---|---|---|
| | Accuracy (%) | DI | Accuracy (%) | DI |
| 68 | 100 | 0.690 | 82.1 | 0.236 |
| 85 | 100 | 0.687 | 82.1 | 0.273 |
| **97** | **100** | **0.703** | **82.1** | **0.287** |
| 0 | 100 | 0.703 | 79.5 | 0.241 |
| 4 | 100 | 0.596 | 79.5 | 0.253 |
| 10 | 100 | 0.571 | 79.5 | 0.194 |
| **32** | 100 | **0.721** | 79.5 | 0.258 |
| 39 | 100 | 0.588 | 79.5 | 0.248 |
| 40 | 98.5 | 0.583 | 79.5 | 0.232 |
| 6 | 97.0 | 0.564 | 79.5 | 0.214 |

benchmark in addition to classification accuracy to determine the best classification model. In this study, we formulated the discriminant index (DI) as an important guideline for optimal model selection. Table 5 showed that the best model was obtained from the 97-th grid search because it had the largest testing DI value of 0.287 with a comparatively higher training DI value of 0.703. Although the model obtained from the 32-th grid search iteration had the highest training DI value, it was not chosen because it had smaller testing accuracy and testing DI value. More importantly, a smaller testing DI value might signify a less reliable model regarding the generalization performance.

The confusion matrix and illustration of classification results obtained from the best model using the testing dataset are depicted in Fig. 8 and Fig. 9, respectively. There were five out of seven misclassified cases in which the samples' Brix values were relatively close to the decision boundaries at the Brix value of 20 and 24 (Fig. 9). There were only three cases of "grade III" samples (i.e., unacceptable samples) that were classified as "grade II" samples (i.e., acceptable samples). However, the Brix values of these samples were very close to the decision boundary between "grade II' and "grade III" regions. Their maximum Brix distance was only 0.7, which was relatively small.

It was noted that one "grade I" sample with a large Brix value was misclassified into "grade II." A "grade II" sample

**TABLE 6.** Discriminant indices of the high-performance classifiers developed with significant interactance responses as inputs.

| Grid search iteration | Training dataset | | Testing dataset | |
|---|---|---|---|---|
| | Accuracy (%) | DI | Accuracy (%) | DI |
| 94 | 100 | 0.684 | 76.9 | 0.218 |
| **97** | **100** | **0.700** | **76.9** | **0.219** |
| 67 | 97.0 | 0.584 | 76.9 | 0.183 |
| 36 | 95.5 | 0.530 | 76.9 | 0.213 |
| 66 | 94.0 | 0.508 | 76.9 | 0.189 |
| 27 | 86.6 | 0.392 | 76.9 | 0.172 |
| 83 | 86.6 | 0.415 | 76.9 | 0.188 |
| 32 | 85.1 | 0.392 | 76.9 | 0.180 |
| 3 | 83.6 | 0.405 | 76.9 | 0.194 |
| 61 | 83.6 | 0.413 | 76.9 | 0.177 |

**TABLE 7.** Performance comparison between the proposed classifier with other classifiers.

| Performance index | Proposed RF classifier | SVM classifier | MLR classifier |
|---|---|---|---|
| Training accuracy | 100% | 86.6% | 83.6% |
| Testing accuracy | 82.1% | 66.7% | 61.5% |
| Training DI | 0.703 | - | 0.414 |
| Testing DI | 0.287 | - | 0.224 |
| Number of close-to-decision-boundary misclassified samples | 5 | 10 | 10 |



**FIGURE 8.** Confusion matrix for testing data.

**TABLE 8.** Misclassified training cases.

| Sample | Brix value | True grade | False predicted grade | |
|---|---|---|---|---|
| | | | SVM classifier | MLR classifier |
| 1 | 24.6 | 1 | - | 2 |
| 2 | 24.8 | 1 | 2 | 2 |
| 3 | 24.4 | 1 | 2 | 2 |
| 4 | 24.0 | 2 | 1 | 1 |
| **5** | **20.4*** | **2** | - | **3** |
| 6 | 21.5 | 2 | 3 | 3 |
| **7** | **20.5** | **2** | **3** | **3** |
| **8** | **20.1** | **2** | **3** | **3** |
| 9 | 21.2 | 2 | 3 | 3 |
| 10 | 22.3 | 2 | 3 | - |
| 11 | 15.2 | 3 | - | 2 |
| 12 | 18.9 | 3 | 2 | 2 |
| Number of misclassified samples | | | 9 | 11 |
| Classification accuracy | | | 86.6% | 83.6% |

\* *Samples whose Brix values were close to the decision boundaries were denoted in bold font style.*



**FIGURE 9.** Illustration of classification results with testing data.
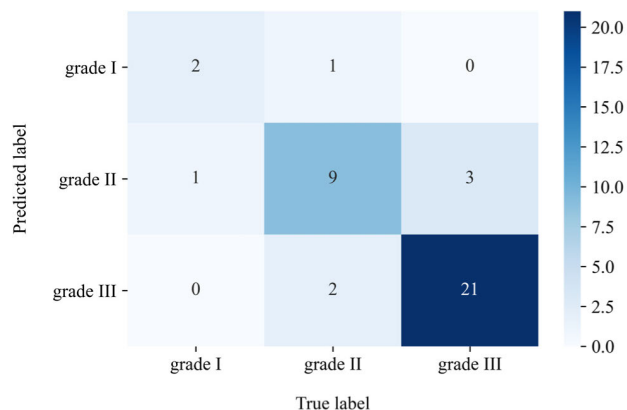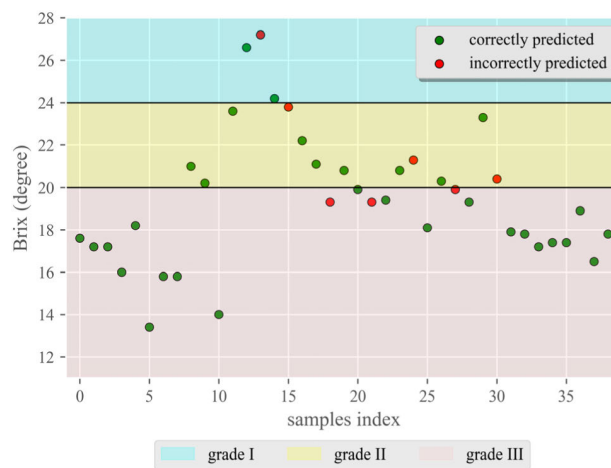
with a Brix value of 21.3 was also misclassified into "grade III." Therefore, an enlarged dataset with more samples in "grade I" region and in the regions close to the decision boundaries (i.e., at 20 and 24 degrees Brix) might help improve classification accuracy.

To confirm the outperformance of latent features over the raw interactance responses, a similar grid search was performed again for 100 iterations using the same parameter configuration (Table 3). Ten models with the highest testing accuracy were obtained and their DI values were also calculated (Table 6). It was obvious that the highest testing accuracy only reached 76.9%, less than the testing accuracy of the best model (i.e., 82.1%). Also, the testing DI values of these models were comparably smaller with the largest testing DI value of only 0.219. These figures further supported the

use of latent features for developing better classifiers and DI for model selection. Therefore, by implementing RF as a variable selector followed by PLS as a feature extractor, significant latent features could be obtained from high-impact spectral responses for noninvasive grading of fruit quality.

Using the same latent features **T**, an SVM classifier and a multinomial logistic regression (MLR) classifier were developed to demonstrate the effectiveness of the proposed RF classifier. The performance comparison of all classifiers was provided in Table 7. The proposed RF classifier outperformed the SVM and MLR classifiers in both training and testing accuracies. Moreover, the training DI value of the RF classifier was much larger than that of the MLR.

The misclassified cases in training and testing were also listed in Table 8 and Table 9, respectively. Most of the misclassified testing samples by the proposed RF classifier were also wrongly graded by the other classifiers. More "grade III" samples were misclassified by SVM and MLR classifiers although they contributed up to 56.7% of the training

**TABLE 9.** Misclassified testing cases.

| Sample | Brix value | True grade | False predicted grade | | |
|---|---|---|---|---|---|
| | | | Proposed RF classifier | SVM classifier | MLR classifier |
| 1 | 27.2 | 1 | 2 | 2 | 2 |
| 2 | 26.6 | 1 | - | 2 | 2 |
| 3 | 24.2* | 1 | - | 2 | 2 |
| 4 | 20.4 | 2 | 3 | 3 | 3 |
| 5 | 21.3 | 2 | 3 | - | 3 |
| 6 | 23.8 | 2 | 1 | - | 1 |
| 7 | 22.2 | 2 | - | - | 3 |
| 8 | 19.9 | 3 | 2 | 2 | 2 |
| 9 | 19.3 | 3 | 2 | 2 | 2 |
| 10 | 19.3 | 3 | 2 | 2 | 2 |
| 11 | 19.9 | 3 | - | 2 | 2 |
| 12 | 19.4 | 3 | - | 2 | - |
| 13 | 19.3 | 3 | - | 2 | 2 |
| 14 | 17.9 | 3 | - | 1 | - |
| 15 | 16.5 | 3 | - | 2 | 2 |
| Number of misclassified samples | | | 7 | 12 | 13 |
| Classification accuracy | | | 82.1% | 66.7% | 61.5% |

\* *Samples whose Brix values were close to decision boundaries were denoted in bold font style.*

dataset. Moreover, five more training and testing samples were misclassified, probably because their Brix values were close to the decision boundaries. Therefore, the RF classifier demonstrated a higher classification power over the SVM and MLR classifiers, and also showed its strong potential for various classification purposes.

## IV. CONCLUSION

Random forest classifier was successfully developed for sweetness grading of Cat Hoa Loc mangoes (Mangifera indica L.) using their interactance responses acquired from low-cost VIS-NIR multispectral sensors. The proposed three-phase RPR approach was very effective for the systematic development of the optimal RF classifier mainly due to the following points:

1) Significant interactance responses could be determined.
2) Good PLS-based latent features could be extracted from only significant interactance spectral responses rather than the whole spectral data.
3) The best model could be identified among the candidates with identical training and testing accuracies using the formulated discriminant index.

Experimental results showed that the discriminant index (DI), formulated based on the class-prediction probability of the RF classifier, was a very effective measure for performance comparison between RF classifiers. Based on both training and testing DI values, the best model was determined from three models with the same training and testing accuracies of 100% and 82.1%, respectively. The classification accuracy of the model was remarkable in terms of the relatively small and imbalanced training and testing datasets, a high degree of sample variability due to different mango harvest regions and time, and the low cost of

the multispectral sensors. Therefore, it is very promising to extend the proposed approach to more complicated quality grading of fruit using cost-effective VIS-NIR multispectral sensors, and using a larger dataset will be very helpful to improve the generalization performance of the classifier.

## REFERENCES

[1] M. Lauricella, S. Emanuele, G. Calvaruso, M. Giuliano, and A. D'Anneo, "Multifaceted health benefits of Mangifera indica L. (Mango): The inestimable value of orchards recently planted in sicilian rural areas," *Nutrients*, vol. 9, no. 5, p. 525, May 2017.

[2] *Mango Global Export and Top Exporting Countries—Tridge*. Accessed: Apr. 10, 2020. [Online]. Available: https://www.tridge.com/intelligences/mango/export

[3] S. K. Bejo and S. Kamaruddin, "Determination of Chokanan mango sweetness (Mangifera indica) using non-destructive image processing technique," *Austral. J. Crop Sci.*, vol. 8, no. 4, pp. 475–480, 2014.

[4] N. Vélez-Rivera, J. Blasco, J. Chanona-Pérez, G. Calderón-Domínguez, M. de Jesús Perea-Flores, I. Arzate-Vázquez, S. Cubero, and R. Farrera-Rebollo, "Computer vision system applied to classification of 'Manila' mangoes during ripening process," *Food Bioprocess Technol.*, vol. 7, no. 4, pp. 1183–1194, Apr. 2014.

[5] R. J. Nissen, N. D. Duc, and N. M. Chau, "'Cat hoa loc mango quality guide,' the AusAID collaboration of agriculture and rural development (CARD) project 050/04 VIE 'Improvement of export and domestic markets for Vietnamese fruit through improved post-harvest and supply chain management,'" Tech. Rep., 2008.

[6] P. P. Subedi and K. B. Walsh, "Assessment of sugar and starch in intact banana and mango fruit by SWNIR spectroscopy," *Postharvest Biol. Technol.*, vol. 62, no. 3, pp. 238–245, Dec. 2011.

[7] A. Zakaria, A. Y. M. Shakaff, M. J. Masnan, F. S. A. Saad, A. H. Adom, M. N. Ahmad, M. N. Jaafar, A. H. Abdullah, and L. M. Kamarudin, "Improved maturity and ripeness classifications of Magnifera indica cv. Harumanis mangoes through sensor fusion of an electronic nose and acoustic sensor," *Sensors*, vol. 12, no. 5, pp. 6023–6048, May 2012.

[8] P. Wanitchang, A. Terdwongworakul, J. Wanitchang, and N. Nakawajana, "Non-destructive maturity classification of mango based on physical, mechanical and optical properties," *J. Food Eng.*, vol. 105, no. 3, pp. 477–484, Aug. 2011.

[9] M. Haggag, S. Abdelhay, A. Mecheter, S. Gowid, F. Musharavati, and S. Ghani, "An intelligent hybrid experimental-based deep learning algorithm for tomato-sorting controllers," *IEEE Access*, vol. 7, pp. 106890–106898, 2019.

[10] Z. Liu, J. Wu, L. Fu, Y. Majeed, Y. Feng, R. Li, and Y. Cui, "Improved kiwifruit detection using pre-trained VGG16 with RGB and NIR information fusion," *IEEE Access*, vol. 8, pp. 2327–2336, 2020.

[11] N. T. Thinh, N. D. Thong, H. T. Cong, and N. T. T. Phong, "Mango classification system based on machine vision and artificial intelligence," in *Proc. 7th Int. Conf. Control, Mechatronics Automat. (ICCMA)*, Nov. 2019, pp. 475–482.

[12] A. Wendel, J. Underwood, and K. Walsh, "Maturity estimation of mangoes using hyperspectral imaging from a ground based mobile platform," *Comput. Electron. Agricult.*, vol. 155, pp. 298–313, Dec. 2018.

[13] S. N. Jha, S. Chopra, and A. R. P. Kingsly, "Determination of sweetness of intact mango using visual spectral analysis," *Biosyst. Eng.*, vol. 91, no. 2, pp. 157–161, Jun. 2005.

[14] T. M. P. Cattaneo and A. Stellari, "Review: NIR spectroscopy as a suitable tool for the investigation of the horticultural field," *Agronomy*, vol. 9, no. 9, p. 503, Sep. 2019.

[15] M. Manley, "Near-infrared spectroscopy and hyperspectral imaging: Nondestructive analysis of biological materials," *Chem. Soc. Rev.*, vol. 43, no. 24, pp. 8200–8214, Dec. 2014.

[16] S. Saranwong, J. Sornsrivichai, and S. Kawano, "Prediction of ripe-stage eating quality of mango fruit from its harvest quality measured nondestructively by near infrared spectroscopy," *Postharvest Biol. Technol.*, vol. 31, no. 2, pp. 137–145, Feb. 2004.

[17] Z. Schmilovitch, A. Mizrach, A. Hoffman, H. Egozi, and Y. Fuchs, "Determination of mango physiological indices by near-infrared spectrometry," *Postharvest Biol. Technol.*, vol. 19, no. 3, pp. 245–252, Jul. 2000.

[18] I. Saranwong, J. Sornsrivichai, and S. Kawano, "Improvement of PLS calibration for brix value and dry matter of mango using information from MLR calibration," *J. Near Infr. Spectrosc.*, vol. 9, no. 4, pp. 287–295, Oct. 2001.

[19] A. A. Munawar, H. Meilina, and Z. Zulfahrizal, "The application of near infrared reflectance spectroscopy as a fast and non-destructive method to determine inner quality parameters of intact mango," in *Proc. IGC*, 2019, pp. 1–8.

[20] Z. Zhang, Z. Cai, G. Han, W. Sheng, and J. Liu, "Estimation of glucose absorption spectrum at its optimum pathlength for every wavelength over a wide range," *Spectrosc. Lett.*, vol. 49, no. 9, pp. 588–595, Oct. 2016.

[21] K. B. Walsh, J. A. Guthrie, and J. W. Burney, "Application of commercially available, low-cost, miniaturised NIR spectrometers to the assessment of the sugar content of intact fruit," *Funct. Plant Biol.*, vol. 27, no. 12, pp. 1175–1186, 2000.

[22] P. P. Subedi, K. B. Walsh, and G. Owens, "Prediction of mango eating quality at harvest using short-wave near infrared spectrometry," *Postharvest Biol. Technol.*, vol. 43, no. 3, pp. 326–334, Mar. 2007.

[23] P. Rungpichayapichet, B. Mahayothee, M. Nagle, P. Khuwijitjaru, and J. Müller, "Robust NIRS models for non-destructive prediction of postharvest fruit ripeness and quality in mango," *Postharvest Biol. Technol.*, vol. 111, pp. 31–40, Jan. 2016.

[24] *F-751 Mango Quality Meter | Tools for Applied Food Science | Felixinstruments.com*. Accessed: Apr. 26, 2020. [Online]. Available: https://felixinstruments.com/food-science-instruments/portable-nir-analyzers/f-751-mango-quality-meter/

[25] A. A. Munawar, R. Hayati, and D. Wahyuni, "The application of near infrared technology as a rapid and non-destructive method to determine vitamin C content of intact mango fruit," *INMATEH-Agric. Eng.*, vol. 58, no. 2, pp. 285–292, 2019.

[26] E. Alós, M. J. Rodrigo, and L. Zacarias, "Ripening and senescence," in *Postharvest Physiology and Biochemistry of Fruits and Vegetables*, E. M. Yahia, Ed. Cambridge, U.K.: Woodhead Publishing, 2019, ch. 7, pp. 131–155.

[27] P. P. S. Gill, S. K. Jawandha, and N. Kaur, "Transitions in mesocarp colour of mango fruits kept under variable temperatures," *J. Food Sci. Technol.*, vol. 54, no. 13, pp. 4251–4256, Dec. 2017.

[28] H. K. Lichtenthaler and C. Buschmann, "Chlorophylls and carotenoids: Measurement and characterization by UV-VIS spectroscopy," *Current Protocols Food Anal. Chem.*, vol. 1, no. 1, pp. F4.3.1–F4.3.8, 2001.

[29] N. T. Tran and M. Fukuzawa, "A portable spectrometric system for quantitative prediction of the soluble solids content of apples with a pre-calibrated multispectral sensor chipset," *Sensors*, vol. 20, no. 20, pp. 1–11, 2020.

[30] AMS. *AS7265x Smart Spectral Sensor*. Accessed: Oct. 6, 2020. [Online]. Available: https://ams.com/as7265x#tab/description

[31] A. S. Franca and L. M. L. Nollet, *Spectroscopic Methods in Food Analysis*. Boca Raton, FL, USA: CRC Press, 2017.

[32] Y.-Y. Pu and D.-W. Sun, "Vis–NIR hyperspectral imaging in visualizing moisture distribution of mango slices during microwave-vacuum drying," *Food Chem.*, vol. 188, pp. 271–278, Dec. 2015.

[33] L. S. Magwaza and U. L. Opara, "Analytical methods for determination of sugars and sweetness of horticultural products—A review," *Scientia Horticulturae*, vol. 184, pp. 179–192, Mar. 2015.

[34] G. Biau and E. Scornet, "A random forest guided tour," *TEST*, vol. 25, no. 2, pp. 197–227, Jun. 2016.

[35] G. Louppe, L. Wehenkel, A. Sutera, and P. Geurts, "Understanding variable importances in Forests of randomized trees," in *Proc. Adv. Neural Inf. Process. Syst.*, 2013, pp. 431–439.

[36] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Vanderplas, and D. Cournapeau, "Scikit-learn: Machine learning in Python," *J. Mach. Learn. Res.*, vol. 12, pp. 2825–2830, Oct. 2011.

[37] H. Abdi, "Partial least squares regression and projection on latent structure regression (PLS regression)," *Wiley Interdiscip. Rev., Comput. Statist.*, vol. 2, no. 1, pp. 97–106, 2010.

[38] K. S. Ng. (2013). *A Simple Explanation of Partial Least Squares*. Accessed: Oct. 6, 2020. [Online]. Available: http://users.rsise.anu.edu.au/~kee/pls.pdf

[39] K. Dunn. (2017). *6.6. Principal Component Regression (PCR)—Process Improvement Using Data*. Accessed: Mar. 19, 2020. [Online]. Available: https://learnche.org/pid/latent-variable-modelling/principal-components-regression

[40] P. D. Alamar, E. T. S. Caramês, R. J. Poppi, and J. A. L. Pallone, "Quality evaluation of frozen guava and yellow passion fruit pulps by NIR spectroscopy and chemometrics," *Food Res. Int.*, vol. 85, pp. 209–214, Jul. 2016.

[41] A. W. Caulk and K. A. Janes, "Robust latent-variable interpretation of *in vivo* regression models by nested resampling," *Sci. Rep.*, vol. 9, no. 1, pp. 1–15, Dec. 2019.

[42] *3.3. Metrics and Scoring: Quantifying the Quality of Predictions—Scikit-Learn 0.23.2 Documentation*. Accessed: Oct. 14, 2020. [Online]. Available: https://scikit-learn.org/stable/modules/model_evaluation.html#r2-score

[43] J. Vaicenavicius, D. Widmann, C. Andersson, F. Lindsten, J. Roll, and T. B. Schön, "Evaluating model calibration in classification," Feb. 2019, *arXiv:1902.06977*. [Online]. Available: https://arxiv.org/abs/1902.06977

[44] *Ncnghiem/Discriminant_Index*. Accessed: May 4, 2020. [Online]. Available: https://github.com/ncnghiem/discriminant_index/

[45] C.-W. Hsu and C.-J. Lin, "A comparison of methods for multiclass support vector machines," *IEEE Trans. Neural Netw.*, vol. 13, no. 2, pp. 415–425, Mar. 2002.

[46] M. Zareef, Q. Chen, M. M. Hassan, M. Arslan, M. M. Hashim, W. Ahmad, F. Y. H. Kutsanedzie, and A. A. Agyekum, "An overview on the applications of typical non-linear algorithms coupled with NIR spectroscopy in food analysis," *Food Eng. Rev.*, vol. 12, no. 2, pp. 173–190, Jun. 2020.

[47] S. Hongmao, "Quantitative structure-activity relationships," in *A Practical Guide to Rational Drug Design*. Amsterdam, The Netherlands: Elsevier, 2016, pp. 163–192.

[48] I. Ahmad, M. Basheri, M. J. Iqbal, and A. Rahim, "Performance comparison of support vector machine, random forest, and extreme learning machine for intrusion detection," *IEEE Access*, vol. 6, pp. 33789–33795, May 2018.

[49] C.-W. Hsu, C.-C. Chang, and C.-J. Lin. (2016). *A Practical Guide to Support Vector Classification*. Accessed: Aug. 19, 2020. [Online]. Available: http://www.csie.ntu.edu.tw/~cjlin

[50] D. Lei, H. Zhang, H. Liu, Z. Li, and Y. Wu, "Maximal uncorrelated multinomial logistic regression," *IEEE Access*, vol. 7, pp. 89924–89935, 2019.

[51] M. E. Maldonado-Celis, E. M. Yahia, R. Bedoya, P. Landázuri, N. Loango, J. Aguillón, B. Restrepo, and J. C. G. Ospina, "Chemical composition of mango (Mangifera indica L.) fruit: Nutritional and phytochemical compounds," *Frontiers Plant Sci.*, vol. 10, p. 1073, Oct. 2019.

[52] M. Golic, K. Walsh, and P. Lawson, "Short-wavelength near-infrared spectra of sucrose, glucose, and fructose with respect to sugar concentration and temperature," *Appl. Spectrosc.*, vol. 57, no. 2, pp. 139–145, Feb. 2003.

[53] A. F. Omar, H. Atan, and M. Z. MatJafri, "NIR spectroscopic properties of aqueous acids solutions," *Molecules*, vol. 17, no. 6, pp. 7440–7450, Jun. 2012.

[54] A. F. Omar, H. Atan, and M. Z. MatJafri, "Peak response identification through near-infrared spectroscopy analysis on aqueous sucrose, glucose, and fructose solution," *Spectrosc. Lett.*, vol. 45, no. 3, pp. 190–201, Apr. 2012.

[55] A. P. Medlicott and A. K. Thompson, "Analysis of sugars and organic acids in ripening mango fruits (Mangifera indica L. Var Keitt) by high performance liquid chromatography," *J. Sci. Food Agricult.*, vol. 36, no. 7, pp. 561–566, Jul. 1985.

**CHANH-NGHIEM NGUYEN** (Member, IEEE) received the M.S. degree in mechatronics from the Asian Institute of Technology, Pathumthani, Thailand, in 2007, and the Ph.D. degree from the Graduate School of Engineering Science, Osaka University, Osaka, Japan, in 2012. Since 2005, he has been a Lecturer with the Department of Automation Technology, College of Engineering Technology, Can Tho University. His research interests include machine vision, microrobotics, embedded control systems, GNSS applications, machine learning, remote sensing, multispectral and hyperspectral imaging, and applications.

**QUOC-THANG PHAN** received certificates for completion of Agriculture Internship Program from the Yanmar Research & Development Center, Maibara, Japan, in 2017, and the KIT Bio Tech & IT Spring School Program 2018 while he was an Exchange Student at the Image Processing Laboratory, Kyoto Institute of Technology, Japan. He received the bachelor's degree in mechatronics from the College of Engineering Technology, Can Tho University, in 2020. His research interests include machine learning, computer vision, and multispectral and hyperspectral sensing.

**NHUT-THANH TRAN** received the B.Eng. degree in mechatronics from Can Tho University, Vietnam, in 2008, and the M.Eng. degree in automation from the Ho Chi Minh City University of Technology, Vietnam, in 2011. He is currently pursuing the Ph.D. degree in engineering design with the Kyoto Institute of Technology, Japan. He has been teaching with the Department of Automation Technology, Can Tho University, since 2012. His research interests include automated systems, precision agriculture, spectroscopy application, and image analysis.

**PHUOC-LOC NGUYEN** received the master's degree in electronic engineering from the Ho Chi Minh City University of Technology and Education, Vietnam, in 2013. He is currently pursuing the Ph.D. degree with the Department of Automation Technology, College of Engineering Technology, Can Tho University. His research interests include machine learning, computer vision, and multispectral and hyperspectral sensing.

**MASAYUKI FUKUZAWA** (Member, IEEE) received the B.Eng. degree from the Department of Electronics, Kyoto Institute of Technology (KIT), in 1992, and the M.Eng. and D.Eng. degrees from the Graduate School of Science and Technology, KIT, in 1994 and 1997, respectively. He worked as a Japan Science Promotion Society (JSPS) Research Fellow from 1996 to 1997. Since 1997, he has been working with KIT, where he is currently an Associate Professor. His specialty is image instrumentation and processing as multidimensional signal. His current research interests include image and video processing for clinical diagnosis, optical instrumentation of semiconductor crystals, and intelligent image sensors.

**CHI-NGON NGUYEN** received the B.S. degree in electrical engineering from Can Tho University, in 1996, the M.S. degree in electrical engineering from the Ho Chi Minh City University of Technology, in 2001, and the Ph.D. degree from the University of Rostock, Germany, in 2007. Since 1996, he has been working with Can Tho University. He is currently a Senior Lecturer with the Department of Automation Technology. He is also working as a position of the Director of the Electrical and Electronic Center, and the Dean of the College of Engineering Technology, Can Tho University. His research interests include intelligent control, medical control, pattern recognition, classifications, speech recognition, and computer vision.

· · ·