# Long-term evaluation of soluble solids content of apples with biological variability by using near-infrared spectroscopy and calibration transfer method

**6 authors**, including:

Shuxiang Fan
National Engineering Research Center for Information Technology in Agriculture
**86** PUBLICATIONS   **2,988** CITATIONS

SEE PROFILE

Jiangbo Li
National Engineering Research Center of Intelligent Equipment for Agriculture
**100** PUBLICATIONS   **3,231** CITATIONS

SEE PROFILE

Yu Xia
Shaanxi University of Science and Technology
**19** PUBLICATIONS   **644** CITATIONS

SEE PROFILE

Xi Tian
China Agricultural University
**53** PUBLICATIONS   **1,401** CITATIONS

SEE PROFILE

# Long-term evaluation of soluble solids content of apples with biological variability by using near-infrared spectroscopy and calibration transfer method

Shuxiang Fan[a,b,c], Jiangbo Li[a,b,c], Yu Xia[a,b], Xi Tian[a,b,c], Zhiming Guo[d], Wenqian Huang[a,b,c,*]

[a] *Beijing Research Center of Intelligent Equipment for Agriculture, Beijing, 100097, China*
[b] *National Research Center of Intelligent Equipment for Agriculture, Beijing, 100097, China*
[c] *Key Laboratory of Agri-Informatics, Ministry of Agriculture, Beijing, 100097, China*
[d] *School of Food and Biological Engineering, Jiangsu University, Zhenjiang, 212013, China*

## ARTICLE INFO

## ABSTRACT

The long-term performance of a near-infrared (NIR) calibration model for soluble solids content (SSC) prediction has been investigated using apples with biological variability collected from 2012 to 2018. The NIR spectrum in the range of 4000–10,000 cm$^{-1}$ was acquired around equator position for each sample. Partial least squares (PLS) was used to develop calibration model based on the samples harvested in 2012 and 2013. The model was then applied to predict the SSC of samples in five separate data sets collected from 2014 to 2018, resulting in a lower performance with higher RMSEP values in the range of 0.704–1.716%. After applying the slope and bias (S/B) correction method, ten samples were selected from each prediction set and used to adjust the model; the prediction results for five independent prediction sets were improved, with RMSEP values ranging from 0.501% to 0.654%. Subsequently, competitive adaptive reweighted sampling (CARS) and successive projections algorithm (SPA) methods were used to select the most effective wavelengths for the determination of SSC. The calibration model built with 15 wavelengths, combined with the S/B correction method, could replace the full spectral range to detect the SSC of apples over a long period of time, with $R_p$ and RMSEP for five prediction sets being 0.919, 0.937, 0.908, 0.896, 0.924 and 0.592, 0.637, 0.513, 0.523, 0.500%, respectively. Overall, the proposed method in this study could make the model valid and robust over a long time and make the biological variability a negligible interference for SSC prediction, thereby providing potential for SSC prediction in practical application.

## 1. Introduction

Near-infrared (NIR) spectroscopy is a non-destructive, low-cost, accurate and reliable method, and can be used to analyze multiple attribute simultaneously without generating waste, and can be used to measure a wide variety of agricultural products and parameters (Nicolaï et al., 2007). When NIR radiation penetrates the product, the spectral responses change as a result of scattering and absorption processes. The spectral change depends on both the internal chemical composition of the product and the light scattering properties. Therefore, the recorded NIR spectroscopy contains both physical and chemical information of the samples (Fan et al., 2016a; Oliveira et al., 2014). Multivariate statistical techniques are then applied to extract useful information about quality attributes from spectroscopic signals.

Most NIR research on apple fruit have focused on assessing internal quality attributes, especially the SSC determination (Bobelyn et al., 2010; Giovanelli et al., 2014; Guo et al., 2016; Lu et al., 2000; Zou et al., 2007). However, a multivariate spectral calibration model is restricted in prediction capabilities when an existing model is applied to spectra that were measured under new environmental and instrumental conditions, or on a separate instrument, limiting the technique use in practical application (Feudale et al., 2002). A variety of biological properties (cultivar, season, shelf-life and geographical origin) of apple fruit influence the optical propagation properties and interaction behaviors with incident light, thus decreasing the inspection accuracy and robustness of calibration models (Zhang et al., 2017). In order to get a more stable and robust NIR calibration model of apple SSC for practical applications, the effect of biological variability on the spectra and robustness of NIR models for SSC prediction has been studied (Bobelyn et al., 2010; Fan et al., 2015b; Peirs et al., 2003). A dedicated

---

calibration model built for each season, cultivar or origin, show good prediction abilities if the predicted samples were from the same season, cultivar, or origin. However, it is time consuming and complicated to establish a calibration model for each biological property. Moreover, the dedicated models decreased their universalities, practicalities and accuracy when they were used to measure the SSC of apples from the other cultivars, seasons, or origins. Hence, dedicated models for each biological property are not suitable to the SSC prediction during commercial applications. On the other hand, a global model developed by a calibration set covering these variations could increase the resistance to small deviations from different samples, and be useful to develop a robust model. Accurate creation of a global model requires the analyst to foresee new sources of variance in the data, however, which is often difficult to determine. In addition, a global tend to be very complex and the predictions from a global model can be less accurate than those from a dedicate model built for a certain biological property if small nonmodeled variations are present in the signal (Feudale et al., 2002). There has been little investigation of the combined effects of factors on calibration model establishment and application, limiting developments of a robust global model. In addition, model invalidation can arise from changes in the instrumental response functions such as aging of light sources, probes, and detectors; instability of the signal over time or nonlinearities in the spectra can pose challenges to long-term application of a calibration model (Fearn, 2001).

Calibration transfer addresses the issue of adjusting spectra and/or the primary model such that the model continually predicts with the same quality using the newly acquired spectra, and at the same time avoids expensive and time-consuming full recalibration (Fearn, 2001; Feudale et al., 2002). The broader context, this includes the concept of using the calibration model over time on multiple instruments, or involve building a model in the original primary condition and then maintaining this model to form accurate predictions of new samples measured in new noncalibrated conditions (Kunz et al., 2010; Workman, 2018). Several strategies have been proposed and widely used to transfer a calibration model between different conditions or spectral instruments, including standardization (DS), piecewise direct standardization (PDS) (Wang et al., 1991), Orthogonal Signal Correction (OSC) (Sjöblom et al., 1998), ShenK–Westerhaus (Fearn, 2001) and Rank-Kennard-Stone-PDS (Liang et al., 2016). Some of these methods have been successfully used in the calibration transfer for fruit internal quality detection. A calibration model for SSC of apple developed on a Fourier transform based spectrophotometer has been successfully transferred to a diode array (DA) spectrophotometer using the PDS method (Alamar et al., 2007). Prediction of SSC after the PDS procedure was improved, with the RMSEP decreased from 3.41% to 0.64%. Pu et al. (2018) applied the PDS algorithm for calibration transfer between a handheld micro NIR spectrometer and a desktop hyperspectral imaging for predicting SSC in banana. They obtained comparable prediction accuracies for the two instruments. In addition to correcting the spectral response between different instruments, these methods can also be used to transfer a calibration model between spectra measured at different temperatures (Lin, 1998). However, the methods mentioned above require a small set of samples (standardization set) to be measured in both conditions or instruments and allow mapping from one spectral domain to another. These methods would become more challenging if the long-term application of a calibration model is considered because the transfer samples, whose properties should be stable and not change over a long time, need to be prepared and measured to correct the spectral responses of different instruments or in different conditions. Slope/Bias (S/B) correction is a simple calibration transfer method by adjusting the final predicted results rather than adjusting the spectral space. Since S/B is a univariate approach based on linear correction of predicted values, it works well when the differences between spectral responses are rather simple (Zhang et al., 2018). Consequently, applying the existing model to detect new samples with different biological properties by using S/B calibration transfer method,

might overcome the drawbacks of dedicated or global SSC prediction models, thereby making the model transferable and avoiding full recalibration.

Modern spectrometers usually possess high resolution, with hundreds or thousands of spectral variables including collinearity, redundancies, and noise, thus increasing the complexity of calibration model built with the full range of spectra, and hindering the computing speed (Wang et al., 2015). Hence, the models built using full spectra are not suitable for on-line or real-time detection in a rapid and non-destructive manner. In addition, the irrelevant information within spectra would affect the accuracy and robustness of the model. Therefore, numerous methods have developed around the theme of selection of effective wavelengths, including Monte Carlo based uninformative variable elimination (MC-UVE) (Cai et al., 2008), genetic algorithm (GA) (Durand et al., 2007), competitive adaptive reweighted sampling (CARS) (Li et al., 2009), successive projection algorithm (SPA) (Araújo et al., 2001), and random frog (Li et al., 2012). Among these methods, CARS is a cost-effective method of removing variables with small means of regression coefficients by using the effective principle 'survival of the fittest' on which Darwin's Evolution Theory is based. Previous studies have tested the effectiveness of this technique in building simplified and high-performance calibration models, such as for the prediction of caffeine content of coffee beans (Zhang et al., 2016), total theaflavins content in black tea (Ouyang et al., 2019), and SSC of pear and apple (Fan et al., 2016a; Travers et al., 2014). Additionally, CARS has been proved to be better than MC-UVE in selecting effective wavelengths for multivariate analysis (Li et al., 2009). However, the number of wavelengths selected by CARS was still too high when it was dealing with a large number of variables (Travers et al., 2014). SPA is a novelty variable selection algorithm by using simple projection operation to select the variables with minimum of collinearity (Araújo et al., 2001). Because the disadvantages of variables selection by SPA were its low signal to noise ratio (S/N) and useless variables for establishment of the model when the full spectrum was considered, the informative variables were firstly obtained by other wavelengths selection methods before SPA was performed (Li et al., 2014; Xu et al., 2012). Therefore, the combination of CARS and SPA, which combined the advantages of CARS and SPA, might obtain more effective wavelengths for SSC prediction compared with they were used individually. Although wavelength selection could simplify the model and improve detection efficiency, there remains a paucity of analysis of discovering chemical compounds matching to the selected wavelengths. More importantly, the principles and applications of different variable selection methods usually lead to different results, causing there is no single, universally optimal technique for selecting key wavelengths for a specific application (Liu et al., 2014). Therefore, validation of the effectiveness of the selected wavelengths became more important but still was limited.

The objectives of this study were: (1) to evaluate the long-term performance of a calibration model for SSC prediction using the S/B correction method based on the data sets with biological variability collected in several successive years, and (2) to select and validate the effective wavelengths for SSC prediction using different independent prediction sets.

## 2. Materials and methods

### 2.1. Samples

A total of 1053 'Fuji' apples (*Malus domestica* Borkh. cv. Fuji) harvested from 2012 to 2018 were used in this study (Table 1). The apples collected in 2016 were from an orchard in Shandong province, the samples selected in 2015 and 2017 were from two different orchards in Beijing, and the remaining samples were purchased in several local markets in Beijing. After being shipped to the Agricultural Bio-sensing Laboratory at the Beijing Research Center of Intelligent Equipment for Agriculture, all of the samples were stored in laboratory for 24 h before

**Table 1**
Distribution of the number of measured samples and their statistic values of SSC (%) in different years.

| Year | No. of samples | Min. | Max. | Mean | Std. |
|------|----------------|------|------|------|------|
| 2012 | 208 | 10.79 | 21.89 | 15.70 | 2.52 |
| 2013 | 130 | 7.71 | 18.35 | 11.58 | 1.73 |
| 2014 | 160 | 10.81 | 17.13 | 13.45 | 1.28 |
| 2015 | 102 | 9.72 | 17.42 | 14.09 | 1.68 |
| 2016 | 143 | 11.6 | 17.2 | 14.06 | 1.10 |
| 2017 | 150 | 8.0 | 13.6 | 11.43 | 1.05 |
| 2018 | 160 | 10.1 | 16.6 | 13.17 | 1.22 |

the experiment to allow the samples to reach room temperature to reduce the effect of apple temperature on the spectra measurement (Fan et al., 2016b).

## 2.2. Spectra collection and SSC measurement

The NIR spectral data of apples were acquired by an AntarisII FT-NIR spectrometer (Thermo Electron., USA) fitted with an integrating sphere working in diffuse reflectance. Apples were placed upon the device with the stem–calyx axis horizontal. The marketed location around equator was illuminated by a tungsten lamp and the diffuse reflectance spectra was collected by a high sensitivity InGaAs detector covering the spectral range of 10,000–4000 cm$^{-1}$ at 1.928 cm$^{-1}$ resolution, yielding a spectrum of 3112 wavelengths. For each apple, a mean spectrum was recorded as log(1/R) (R = reflectance) by averaging spectra of 16 scans. It is necessary to preprocess the spectral data to remove any irrelevant information and to improve a calibration model's performance (Cen and He, 2007). Several spectral preprocessing algorithms, including moving average (smoothing window of 7 points), multiple scatter correction (MSC), standard normal variate (SNV) transformations, first derivative (smoothing window of 91 points and second-order filtering), and second derivative (smoothing window of 101 points and second-order filtering), were applied to preprocess the spectral data by using the Unscrambler v9.7 software (CAMO PROCESS AS, Oslo, Norway).

Immediately after spectrum acquisition, the SSC values of the tested samples were measured by traditional destructive test. A piece of flesh with peel were cut from the same position where NIR spectroscopy measurements had been carried out. Apple juice was squeezed using a manual fruit squeezer and a refractometer (ARIAS 500, Reichert Technologies, New York, USA) used to measure the SSC of samples harvested from 2012 to 2015. The SSCs were determined by the refractometer (PAL-1, ATAGO, Tokyo, Japan) for the apples harvested in 2016, 2017, and 2018.

## 2.3. Partial least squares regression

Partial least squares (PLS) regression is an efficient tool to model the linear relationship between the multivariate predictor variables X (spectral data) and response variables Y (the properties of interest) (Andersson, 2009; Wold et al., 2001). It has been widely used for building calibration models in NIR regression analysis. PLS could extract a smallest possible set of latent variables (LVs) ordered according to their relevance for predicting the Y variable from a highly correlated and collinear original spectral data. One important step in building a PLS model is to determine the number of LVs, which is often optimized by cross validation of the calibration samples. Leave-one-out cross validation (LOOCV) only changes one sample in the training set in each cross validation cycle, having the consequence that the measure of the predictive ability can be overly optimistic (Giovanelli et al., 2014). Therefore, the optimal number of LVs was determined using 10-fold cross-validation until the root mean square error of cross validation (RMSECV) reached a minimum. The PLS analysis was performed in the

Matlab2016a with libPLS toolbox available at http://www.libpls.net/ (Li et al., 2018).

To evaluate the model performance, the following parameters were calculated: the correlation coefficient of calibration ($R_c$) and prediction ($R_p$), root mean square error of calibration (RMSEC) and prediction (RMSEP), and the ratio of the reference data standard deviation and RMSEP (RPD) (Chang et al., 2001; Ferreira et al., 2013). The calculations of $R_c$, $R_p$, RMSEC and RMSEP are defined in the following equations:

$$R_c = \sum_{i=1}^{n_c} (y_{mi} - \bar{y}_m)(y_{pi} - \bar{y}_p) \left/ \sqrt{\sum_{i=1}^{n_c} (y_{mi} - \bar{y}_m)^2} \sqrt{\sum_{i=1}^{n_c} (y_{pi} - \bar{y}_p)^2} \right. \tag{1}$$

$$R_p = \sum_{i=1}^{n_p} (y_{mi} - \bar{y}_m)(y_{pi} - \bar{y}_p) \left/ \sqrt{\sum_{i=1}^{n_p} (y_{mi} - \bar{y}_m)^2} \sqrt{\sum_{i=1}^{n_p} (y_{pi} - \bar{y}_p)^2} \right. \tag{2}$$

$$RMSEC = \sqrt{\frac{1}{n_c} \sum_{i=1}^{n_c} (y_{pi} - y_{mi})^2} \tag{3}$$

$$RMSEP = \sqrt{\frac{1}{n_p} \sum_{i=1}^{n_p} (y_{pi} - y_{mi})^2} \tag{4}$$

$$RPD = \frac{SD}{RMSEP} \tag{5}$$

where $y_{pi}$ is the predicted value of SSC in fruit number $i$, $y_{mi}$ is the measured value of SSC in fruit number $i$, $\bar{y}_m$ and $\bar{y}_p$ are the mean value of measured and predicted SSC, respectively in the calibration set or prediction set, $n_c$ and $n_p$ are the number of fruits in the calibration set and prediction set, respectively. SD represents the standard deviation of measured SSC in prediction set. Generally, a good model should have higher $R_c$, $R_p$, RPD values, lower RMSEC and RMSEP values (Fan et al., 2015a).

## 2.4. Calibration transfer method

S/B correction is a simple calibration transfer method based on the predicted results rather than the spectral space (Feudale et al., 2002). When the developed model was used to predict the SSCs of an independent prediction set with $M$ samples and $P$ spectral variables, a subset of $N_s$ apple samples was selected randomly from the prediction set and then the spectral matrix ($N_s \times P$) were multiplied by the regression coefficients ($\beta$) of the PLS model, obtaining the predicted SSCs ($Yp$) of those $N_s$ samples. Combined with the measured SSC values ($Y_m$) of the $N_s$ samples determined by destructive test, a univariate linear model fitting these $N_s$ discrete 2-D data points ($Y_m$, $Y_p$) was computed by least-squares.

$$Y_m = aY_p + b \tag{6}$$

where $a$ and $b$ were the slope and the bias of this linear model, respectively. The combination of the regression coefficients and slope/bias correction could be used directly to compute the corrected SSC prediction values ($Y_{correct}$) for the remaining samples in the prediction set.

$$Y_{correct} = a(X_r \times \beta) + b \tag{7}$$

where $X_r$ is an ($M$–$N$) × $P$ matrix containing $P$ spectral responses of remaining $M$–$N$ samples. In order to investigate the influence of the number of selected samples in prediction set on S/B correction performance, we have considered the following four cases: $N_s$ is set to 5, 10, 15 and 20. For each value of $N_s$, the S/B correction procedure was repeated ten times. Only the average prediction results from the ten runs were used in statistical analysis. This would ensure more consistent evaluation of the S/B correction method, since its performance is influenced, to a certain extent, by how many samples are selected. The Kruskal-Wallis test was conducted to test the significant differences of

prediction results acquired by S/B correction method with different number of selected samples. This statistical method is a nonparametric comparison test, which allows us to compare several population means without requiring normality and variance homogeneity (Yao et al., 2013), as in other statistical analysis tests, such as ANOVA (Vargha and Delaney, 1998). A p < 0.05 is considered statistically significant for Kruskal-Wallis test.

### 2.5. Effective wavelengths selection

Wavelength selection reduces the original set composed by hundreds or even thousands of spectral variables by eliminating those variables not contributing to improve the model's performance, obtaining a subset of variables that better define the sample class or better correlate with properties of interest. In this way, the wavelength selection procedure could reduce measurement costs, facilitate model interpretation and improve the quantitative and qualitative results (Mehmood et al., 2012).

CARS is an innovative algorithm to select optimal combination of the effective wavelengths coupled with PLS regression (Li et al., 2009). CARS selects $N$ subsets of wavelengths by $N$ Monte Carlo sampling runs in an iterative manner. In each sampling run, some samples are first randomly chosen in a fixed ratio to build a PLS model. As the coefficients of wavelengths in PLS model could reflect the contribution to the prediction of the properties of interest, they are used as an index for evaluating the importance of each wavelength. Thus, a subset of wavelengths with large absolute coefficients was determined and then retained for the next sampling run. After $N$ sampling runs, the subset with the lowest RMSECV was chosen. The number of Monte Carlo sampling runs was set to 100 and the number of variables to be selected is determined by 10-fold cross validation in this study. The procedure of CARS was performed in the Matlab2016a with libPLS toolbox.

SPA is a flexible technique for variable selection in multivariate calibration (Araújo et al. (2001). It is a forward selection method which starts with one wavelength and incorporates a new one at each iteration, yielding wavelengths whose information content is minimally redundant. The optimal number of variables can be determined according to the smallest mean square error of cross validation (MSECV) in the validation set of MLR calibration. The variables selected by SPA can be used as inputs of multiple linear regression (MLR), PLS (Liu et al., 2009), and Least squares support vector machines (LS-SVM) (Cheng and Sun, 2015) to build linear and nonlinear models in spectral analysis. The procedure of SPA was carried out by a graphical user interface (GUI_SPA) in Matlab2016a software available at http://www.ele.ita.br/~kawakami/spa/.

## 3. Results and discussion

### 3.1. spectra features

The original mean FT-NIR spectra of 'Fuji' apples from different years in the region of 4000–10,000 cm$^{-1}$ are shown in Fig. 1A. As FT-NIR spectrometer provided accurate, stable, and reproducible spectra, the average spectra of apples for each year had the similar trend and absorbance values. Two obvious absorption peaks were found around 6915 cm$^{-1}$ and 5204 cm$^{-1}$, which were associated with the absorption of water (Fan et al., 2016a). A slight absorption peak was sited at about 8298 cm$^{-1}$ ascribed to the second overtone of C−H stretching (Magwaza et al., 2011). Similar FT-NIR spectra of apples were also observed in previous literatures (Giovanelli et al., 2014; Zou et al., 2007).

### 3.2. Spectral preprocessing based PLS modeling

Prior to the PLS modeling, the raw spectral data were preprocessed with different methods. The correlation coefficient of cross validation (R$_{CV}$) and RMSECV for apples harvested from 2012 to 2018 were

calculated using PLS analysis based on different preprocessing methods. The results listed in Table 2 are only those preprocessing methods that gave better cross validation results. Although the K–W test indicated that there was no significant difference in RMSECVs obtained by raw and different preprocessed spectra, it was found that the combination of moving average and SNV reduced the RMSEP for all seven data sets compared with raw spectral data, with RMSECV ranging from 0.328 to 0.551%. For solid samples, undesired systematic variations of spectra are primarily caused by light scattering and differences in the effective path length. Given the mode of spectra collection in this study did not cause the variation of path length. Therefore, the light scattering was the main factor adding the systematic variations of spectra. SNV could eliminate the deviations caused by particle size and scattering, plus smoothing is an effective approach for removing high-frequency noise from a spectrum and improving the signal-to-noise ratio (Rinnan et al., 2009). That could explain why the combination of moving average smoothing and SNV obtained better cross validation results compared with other preprocessing methods. Compared with the raw spectral data in Fig. 1A, the spectral difference between different years became negligible after the process of moving average and SNV (Fig. 1B). Since the large enough range of SSC of samples was helpful to develop a good calibration model (Li et al., 2013), the 338 samples collected in 2012 and 2013 were used as the calibration set to develop a PLS model for SSC prediction. The remaining data sets were used as five independent sets to verify the calibration model.

### 3.3. Transfer of calibration model for SSC prediction

The PLS regression model for SSC prediction was constructed depending on the 2012 and 2013 data preprocessed by moving average and SNV. The optimal number of LVs was determined as 16 by 10-fold cross validation. Fig. 2 shows the calibration and cross validation results obtained by the developed model, with fairly good correlation coefficients and low prediction errors. To evaluate the robustness and accuracy of the developed model, the data sets obtained from 2014 and 2018 were used as independent prediction sets to validate it (Table 3). However, higher RMSEP values were obtained for all the prediction sets compared with calibration and cross validation results, especially for 2018, whose RMSEP was as high as 1.716%. These results indicated that the model for SSC prediction was not reliable for practical implementation because of high errors. A possible explanation for this might be the biological variability of those tested apples (including harvest season, year and geographical origin).

However, one interesting finding is that the Rp values for all the prediction sets were around 0.9, representing a very strong relationship between the measured and spectra-predicted SSC values. Those results indicated that the developed model could sort the apples into different SSC classes but could not measure the exact SSCs. Therefore, the model should be corrected to reduce the effects of biological properties on model performance and improve its predicting ability.

For S/B correction method, the number of selected samples that used to calculate the slope and bias should be determined at first. The following four cases were considered: the number is set to 5, 10, 15, and 20. For each case and prediction set, 10 replicates of S/B correction were executed and the averaging prediction results were calculated. The statistical box-plots of the RMSEP values after S/B correction with different number of selected samples are shown in Fig. 3.

The prediction results improved when more samples were selected and used to calculate the slope and bias. However, more samples are selected in the prediction set means more samples have to be analyzed by destructive test to obtain the measured SSCs, thus increasing the complexity of correction process. The statistical tests showed that the RMSEP values obtained by S/B correction with 5 selected samples were statistically higher than those when 20 samples were selected for the prediction sets obtained in 2014, 2016, and 2017. On the other hand, no statistically significant difference was observed between the RMSEP
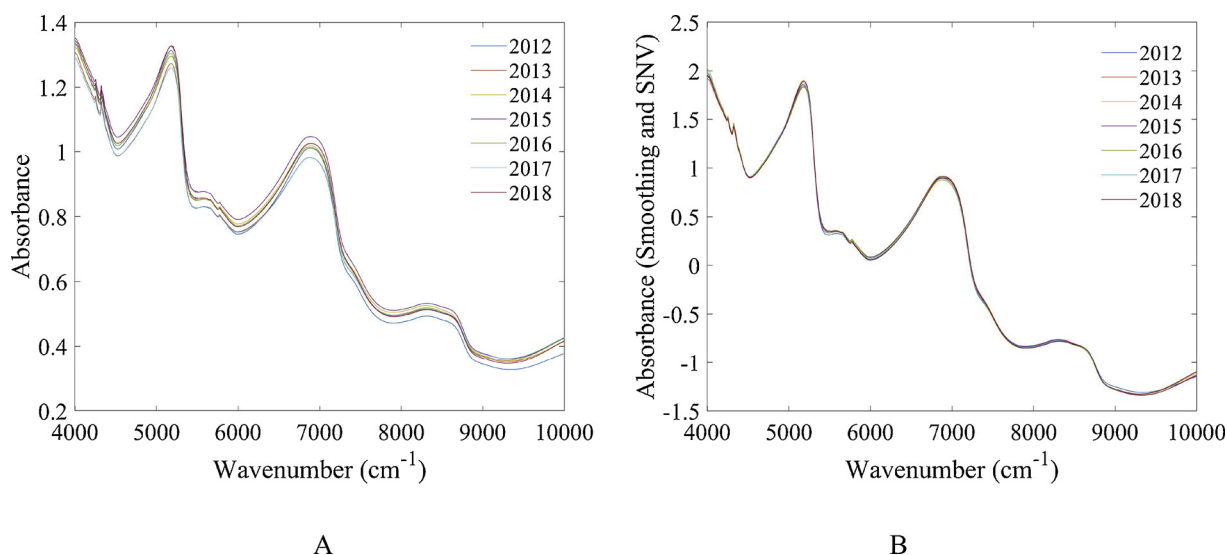
**Fig. 1.** A. Raw absorbance spectra and B. preprocessed spectra after moving average and SNV preprocessing.

**Table 2**
Cross validation results of PLS models built with different preprocessing methods.

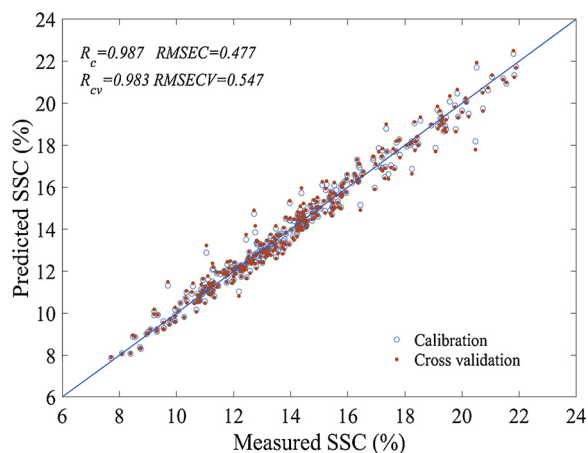| Year | Raw spectra | | Average + SNV | | 1st Derivative | | 2nd Derivative | |
|---|---|---|---|---|---|---|---|---|
| | $R_{cv}$ | RMSECV | $R_{cv}$ | RMSECV | $R_{cv}$ | RMSECV | $R_{cv}$ | RMSECV |
| 2012 | 0.974 | 0.571 | 0.975 | **0.551** | 0.973 | 0.574 | 0.971 | 0.600 |
| 2013 | 0.952 | 0.532 | 0.956 | **0.507** | 0.952 | 0.527 | 0.955 | 0.513 |
| 2014 | 0.935 | 0.455 | 0.940 | **0.434** | 0.937 | 0.448 | 0.935 | 0.454 |
| 2015 | 0.937 | 0.591 | 0.948 | **0.531** | 0.942 | 0.564 | 0.940 | 0.573 |
| 2016 | 0.955 | 0.325 | 0.955 | **0.318** | 0.959 | 0.314 | 0.960 | 0.309 |
| 2017 | 0.885 | 0.488 | 0.889 | **0.479** | 0.900 | 0.457 | 0.902 | 0.451 |
| 2018 | 0.909 | 0.508 | 0.915 | **0.491** | 0.914 | 0.493 | 0.945 | 0.492 |



**Fig. 2.** The prediction results of SSC by PLS model in calibration set and cross validation set.

values based on 10 selected samples and 20 selected samples for all the prediction sets. Hence, in the following analysis, the number of selected samples was set to 10 by balancing the prediction results and complexity of correction process.
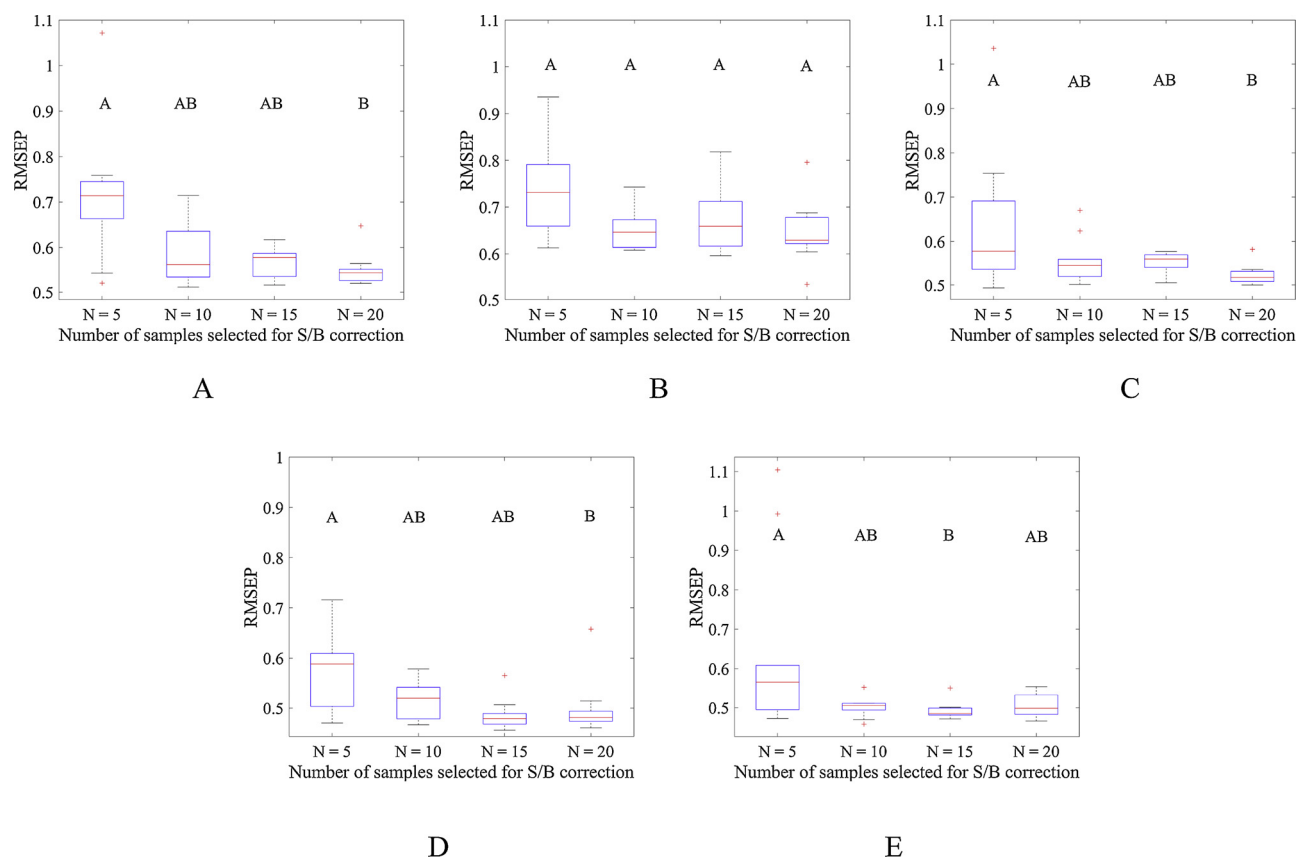
The prediction results after S/B correction, including Rp and RMSEP, were summarized in Table 3 for all the data sets collected from 2014 to 2018. Although the Rp values were comparable to those obtained by the model without S/B correction, the RMSEP values decreased dramatically to 0.501–0.654%, suggesting better prediction
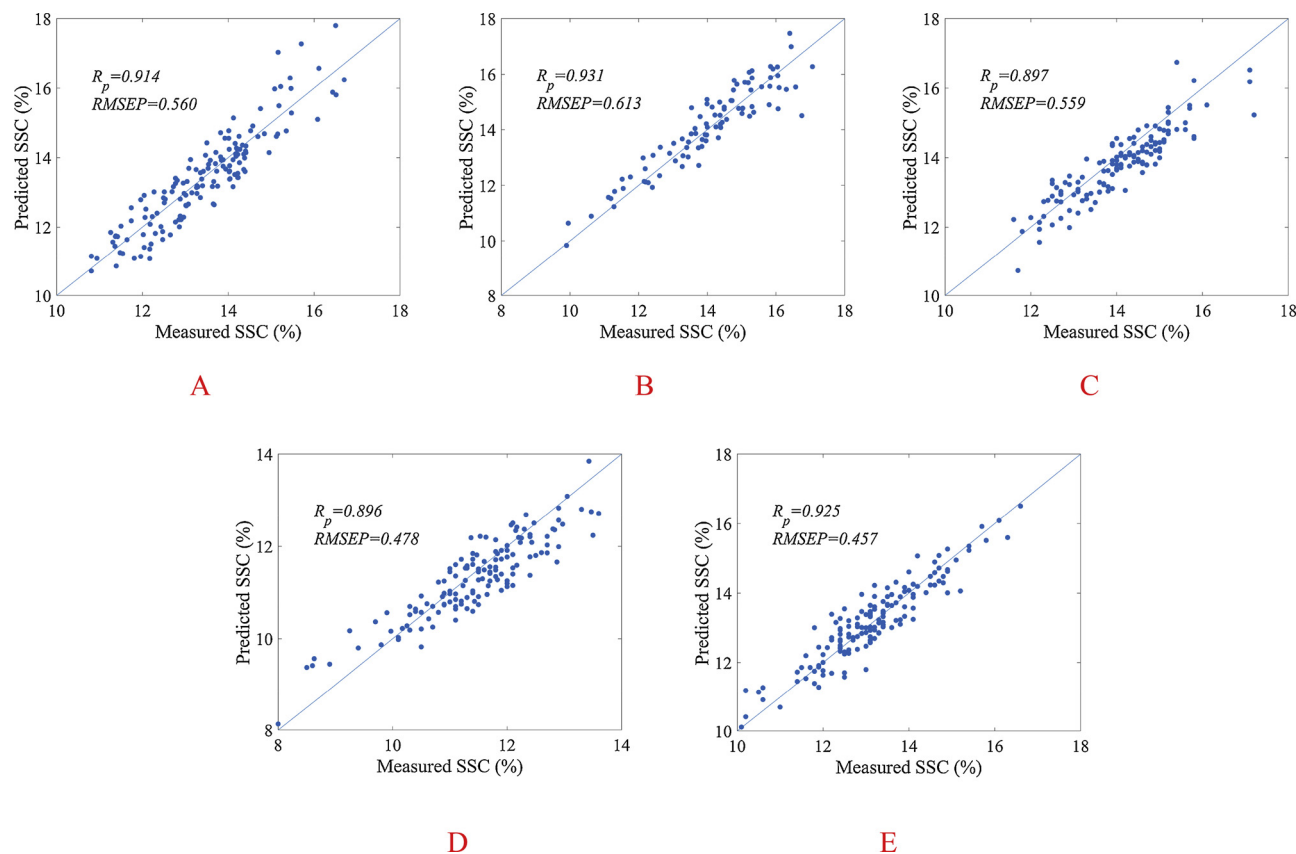
**Table 3**
The prediction results of the calibration model for SSC prediction built with the fusion of 2012 and 2013 data after preprocessing by moving average and SNV. $\overline{R_p}$, $\overline{RMSEP}$ and $\overline{RPD}$: average correlation coefficients and root mean square error of prediction, respectively, over 10 S/B corrections.

| Year | Without S/B correction | | | S/B correction | | |
|---|---|---|---|---|---|---|
| | $R_p$ | RMSEP | RPD | $\overline{R_p}$ | $\overline{RMSEP}$ | $\overline{RPD}$ |
| 2014 | 0.918 | 0.932 | 1.373 | 0.917 | 0.584 | 2.214 |
| 2015 | 0.933 | 0.704 | 2.390 | 0.934 | 0.654 | 2.550 |
| 2016 | 0.895 | 1.052 | 1.048 | 0.895 | 0.554 | 2.004 |
| 2017 | 0.898 | 1.184 | 0.884 | 0.897 | 0.518 | 2.028 |
| 2018 | 0.921 | 1.716 | 0.710 | 0.924 | 0.501 | 2.436 |

results were obtained after S/B correction. The predicted SSC values of five prediction sets predicted by the developed calibration model and S/B correction are plotted against the actual measurements for one of the 10 S/B corrections (Fig. 4). The solid line represents the ideal regression line, as the closer the points are to this line, the better the model is. The S/B correction method is a simpler standardization and this approach is not influenced by problems occurring in the spectral space because it is based on a univariate linear correction of predicted values. These results indicated that the developed calibration model, combined with S/B correction method, has the potential to detect the SSC of apples with a wide variety of biological properties over a long period of time.

**Fig. 3.** Boxplot of RMSEP obtained by the calibration model after S/B correction with 5, 10, 15, and 20 selected samples for the prediction sets obtained in 2014, 2015, 2016, 2017 and 2018, respectively (A–E). Treatments with different letters are statistically significant with each other (K–W test, p < 0.05).



**Fig. 4.** Prediction of SSC obtained by the calibration model after S/B correction for the data set collected in 2014, 2015, 2016, 2017, and 2018, respectively (A–E).
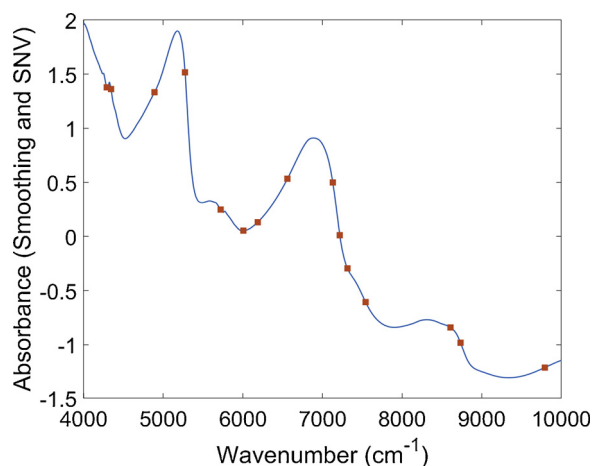
**Fig. 5.** Effective wavelengths selected by CARS-SPA from the spectral data after preprocessing by moving average and SNV.

### 3.4. Effective wavelengths selection and validation

CARS wavelength selection method was used to select the effective wavelengths for the determination of SSC based on the data collected in 2012 and 2013. After CARS processing, the number of wavelengths decreased to 42, but with respect to the online detection of SSC, there were still too many variables for application. In addition, some collinear variables which contain a number of redundant information might still exist in the spectral data. Therefore, to further simplify the model and improve the robustness, SPA was carried out on 42 selected wavelengths for further wavelengths selection. As a result of CARS-SPA calculation, the following 15 wavelengths were selected from 3112 full-spectrum wavelengths: 4290.84, 4344.84, 4890.59, 5274.36, 5719.83, 6007.17, 6186.52, 6560.64, 7129.54, 7218.25, 7314.68, 7542.23, 8606.75, 8734.03, and 9792.76 cm$^{-1}$. The distribution of wavelengths selected by CARS-SPA for apple SSC prediction on the full spectrum is shown in Fig. 5.

Because of the complicated nature of NIR spectra involving large molecules, most NIR band assignments are made on an empirical basis. In complex mixtures, such as in fruit, the presence of multiple bands and the effect of peak broadening result in NIR spectra with broad envelopes and few sharp peaks (Fu and Ying, 2016). Hence, the more reasonable way to evaluate the efficiency of those key wavelengths is to verify the model built with effective wavelengths through different independent prediction sets. At first, the PLS models based on the selected 15 wavelengths were developed after moving average and SNV preprocessing for each of the seven datasets separately (Table 4). The performance of the simplified PLS models was comparable to the models built with full-range spectra in terms of RMSECV and R$_{cv}$ but it used only about 0.5% of the variables used by the full-range spectra (15 vs. 3112), thereby simplifying the prediction model and satisfying the requirements for online detection. The results implied that the key wavelengths selected were effective for SSC prediction.

The second method used to evaluate the effectiveness of these selected wavelengths was to verify the model built with selected wavelengths by using other data sets. On the basis of the selected 15

**Table 4**
Cross validation results of PLS models built with effective wavelengths.

| Cross validation results | Year | | | | | | |
|---|---|---|---|---|---|---|---|
| | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 |
| R$_{cv}$ | 0.981 | 0.961 | 0.941 | 0.951 | 0.959 | 0.902 | 0.921 |
| RMSECV | 0.487 | 0.475 | 0.434 | 0.517 | 0.313 | 0.452 | 0.473 |

wavelengths, the spectral data in calibration set, which was composed of 338 samples collected in 2012 and 2013, was reduced a matrix with a dimension of 338 × 15. A simplified PLS regression model (CARS-SPA-PLS model) was then built with the reduced calibration set, with the R$_{CV}$ and RMSECV of 0.987 and 0.479%, respectively. Compared with the full-spectra PLS model (Fig. 2), the RMSECV was reduced from 0.547 to 0.479%, suggesting again the effectiveness of the selected wavelengths. The simplified CARS-SPA-PLS model is shown as follows:

$$\begin{aligned} Y_{SSC} = &-96.8X_{4290.84cm^{-1}} + 53.19X_{4344.84cm^{-1}} + 26.01X_{4890.59cm^{-1}} \\ &- 34.33X_{5274.36cm^{-1}} - 19.00X_{5719.83cm^{-1}} \\ &- 247.75X_{6007.17cm^{-1}} + 367.34X_{6186.52cm^{-1}} - 204.59X_{6560.64cm^{-1}} \\ &+ 133.46X_{7129.54cm^{-1}} - 271.95X_{7218.25cm^{-1}} \\ &+ 228.76X_{7314.68cm^{-1}} - 198.22X_{7542.23cm^{-1}} + 112.72X_{8606.75cm^{-1}} \\ &- 88.96X_{8734.03cm^{-1}} - 72.14X_{9792.76cm^{-1}} - 23.96 \end{aligned} \tag{8}$$

where $X_{icm^{-1}}$ is the spectral value after preprocessing of moving average and SNV at the wavelength of $i$ cm$^{-1}$ and $Y_{SSC}$ is the predicted SSC value. The performance of the model was then evaluated by five prediction sets collected from 2014 to 2018. The prediction results of the simplified model are summarized in Table 5. The prediction accuracy of the model was still unsatisfied, with the range of RMSEP values being 0.802–1.368%. However, the model showed good performance after S/B correction as it yielded accurate prediction results for all the data sets, with the average RMSEP values of 0.500–0.637%. To compare the performance of full spectra-PLS model and CARS-SPA-PLS model, the samples used for S/B corrections in a specific prediction set were the same when the prediction set was used to evaluate the two models. Compared with the prediction results obtained by the PLS model based on full wavelengths (Table 3), there was no significant difference between the RMSEP values based on full wavelength and effective wavelength. These results suggest again that the selected wavelengths were effective for predicting SSC of apple. Therefore, the 15 effective wavelengths could represent most of the features and characteristics of the original spectra, and could be applied instead of the whole wavelength region to predict the SSC of apple.

## 4. Discussion

A typical RMSEP of SSC prediction of apples obtained by NIR seems to be around 0.5%, but in the few applications where external validation sets with different biological properties or temperature were used to calculate the RMSEP, it is considerably higher (1–1.5%) (Nicolaï et al., 2007). The RMSEP values calculated in this study after S/B correction were 0.500–0.637% for five independent prediction sets, which were comparable to the typical RMSEP. Therefore, the model combined with S/B correction method was not sensitive to the biological variability of the new samples to be predicted and had a high quality of prediction over a long period of time. In postharvest quality sorting and grading, it is sufficient to sort apples into different classes based on their SSCs while unnecessary to determine the SSC exactly.

**Table 5**
The validation results of the CARS-SPA-PLS model for SSC prediction. $\overline{R_p}$, $\overline{RMSEP}$ and $\overline{RPD}$: average correlation coefficients and root mean square error of prediction, respectively, over 10 S/B corrections.

| Year | Without S/B correction | | | S/B correction | | |
|---|---|---|---|---|---|---|
| | $R_p$ | RMSEP | RPD | $\overline{R_P}$ | $\overline{RMSEP}$ | $\overline{RPD}$ |
| 2014 | 0.919 | 0.915 | 1.398 | 0.919 | 0.592 | 2.173 |
| 2015 | 0.936 | 0.802 | 2.095 | 0.937 | 0.637 | 2.618 |
| 2016 | 0.909 | 0.983 | 1.122 | 0.908 | 0.513 | 2.161 |
| 2017 | 0.897 | 1.028 | 1.017 | 0.896 | 0.523 | 2.013 |
| 2018 | 0.924 | 1.368 | 0.890 | 0.924 | 0.500 | 2.439 |

According to the obtained RPD values ($> 2.0$), quantitative predictions for SSC can be obtained (Nicolaï et al., 2007). Hence, it would be feasible to sort apples into different classes according to their SSCs and the proposed method would be helpful during commercial applications.

The effective wavelengths selection in the spectral range of 4000–10,000 $cm^{-1}$ (1000–2500 nm) for SSC prediction of apple has been reported in several studies. Zou et al. (2007) selected 44 key wavelengths from FT-NIR spectra of 'Fuji' apples after preprocessing of multiplicative scatter correction (MSC) and mean centering. Those wavelengths mainly appeared in the following wavelength ranges: 5016–5028 $cm^{-1}$, 5250–5319 $cm^{-1}$, 5403–5529 $cm^{-1}$, and 5760–5901 $cm^{-1}$. In our previous studies, 66 wavelengths were determined from FT-NIR spectra of Fuji apples after 2nd derivative preprocessing (Fan et al., 2016a). Those wavelengths are mainly concentrated in the regions of 4539–4545, 5957–5968, 7357–7364, 7386–7391, 8572–8575, and 8629–8633 $cm^{-1}$. The number of selected wavelengths in this study is much smaller than that in previous studies. The difference of the selected wavelengths might because of the different wavelength selection and preprocessing methods used. The 15 wavelengths selected by CARS-SPA relevant to SSC in NIR spectral region include first and second overtones of C−H stretching at 5719.83, 6007.17, 6186.52 8606.75, and 8734.03 $cm^{-1}$, first overtones of O−H stretching at 7129.54 $cm^{-1}$, the O−H combination band at 4890.59 $cm^{-1}$, C−H combination band at 4290.84 $cm^{-1}$ (Magwaza et al., 2011). These features make great contributions to SSC prediction of apple. In NIR spectroscopy it is not possible to assign a molecular vibration to one specific wavelength, since a spectral value is composed of a combination of different molecular vibrations (Lammertyn et al., 1998). Preprocessing of the spectra makes this problem even more complex. Therefore, it is difficult to interpret the NIR spectra and the selected wavelengths. Light interaction with the turbid fruit tissue involves absorption and scattering, which are characterized by two optical properties: absorption ($\mu_a$) and reduced scattering coefficients ($\mu_s'$) (Qin and Lu, 2008). Knowledge of the optical properties, combined with numerical methods, such as Monte Carlo simulation, could enable us to gain insight about the interaction of light with the fruit tissue, and would be valuable in the development of an optical technique for evaluating properties and characteristics that are indicative of specific quality attributes (Qin and Lu, 2009). However, the structural complexity, diversity, and heterogeneity of fruit, plus the error in data acquisition, instrument calibration, and inverse algorithm implementation, present critical challenges to accurate measurement of optical properties (Hu et al., 2015). Hence, it is still difficult to gain insight about the actual interaction of light with the fruit tissue. Therefore, up to now, the most reasonable way of validating the effectiveness of selected wavelengths and calibration models was using different separate data sets. At the same time, more fundamental research is also required to provide a physicochemical basis of calibration models.

Besides the biological variability and changes in the instrumental response function over time, fruit temperature fluctuations, temperature and humidity variations of experimental conditions could also have a strong influence on the robustness and predictivity of a calibration model. Therefore, in real world application, those factors need to be strictly controlled to avoid their negative impacts on the SSC measurements. When this is not the case, the developed calibration model, selected variables and proposed method in this study need to be further explored to detect SSC of apples with different temperature or experimental conditions in our future studies.

## 5. Conclusion

A PLS model for SSC prediction was developed, and its robustness and accuracy were investigated by different separate data sets from five successive years with different biological variability. The model resulted in a lower performance with higher RMSEP values when it was used directly to predict the SSC for all the separate data sets. But the model showed good performance after S/B correction method, which 10 samples were selected and proven to be enough to adjust the model, with the range of RMSEP values of 0.500–0.637%. In addition, 15 effective wavelengths for SSC prediction has been selected and validated. The calibration model based on 15 effective wavelengths, combined with the S/B correction method, could replace the destructive method to detect the SSC of apples quickly and correctly over a long period of time, having a great potential for apple SSC prediction in the packing line.

## Acknowledgments

## References

Alamar, M.C., Bobelyn, E., Lammertyn, J., Nicolaï, B.M., Moltó, E., 2007. Calibration transfer between NIR diode array and FT-NIR spectrophotometers for measuring the soluble solids contents of apple. Postharvest Biol. Technol. 45, 38–45. https://doi.org/10.1016/j.postharvbio.2007.01.008.

Andersson, M., 2009. A comparison of nine PLS1 algorithms. J. Chemometrics 23, 518–529. https://doi.org/10.1002/cem.1248.

Araújo, M.C.U., Saldanha, T.C.B., Galvão, R.K.H., Yoneyama, T., Chame, H.C., Visani, V., 2001. The successive projections algorithm for variable selection in spectroscopic multicomponent analysis. Chemometr. Intell. Lab. Syst. 57, 65–73. https://doi.org/10.1016/S0169-7439(01)00119-8.

Bobelyn, E., Lammertyn, J., Nicolai, B.M., Saeys, W., Serban, A.S., Nicu, M., 2010. Postharvest quality of apple predicted by NIR-spectroscopy: study of the effect of biological variability on spectra and model performance. Postharvest Biol. Technol. 55, 133–143. https://doi.org/10.1016/j.postharvbio.2009.09.006.

Cai, W., Li, Y., Shao, X., 2008. A variable selection method based on uninformative variable elimination for multivariate calibration of near-infrared spectra. Chemometr. Intell. Lab. Syst. 90, 188–194. https://doi.org/10.1016/j.chemolab.2007.10.001.

Cen, H., He, Y., 2007. Theory and application of near infrared reflectance spectroscopy in determination of food quality. Trends Food Sci. Technol. 18, 72–83. https://doi.org/10.1016/j.tifs.2006.09.003.

Chang, C.-W., Laird, D.A., Mausbach, M.J., Hurburgh, C.R., 2001. Near-infrared reflectance spectroscopy–principal components regression analyses of soil properties. Soil Sci. Soc. Am. J. 65, 480–490. https://doi.org/10.2136/sssaj2001.652480x.

Cheng, J.H., Sun, D.W., 2015. Rapid and non-invasive detection of fish microbial spoilage by visible and near infrared hyperspectral imaging and multivariate analysis. LWT-Food Sci. Technol. 62, 1060–1068. https://doi.org/10.1016/j.lwt.2015.01.021.

Durand, A., Devos, O., Ruckebusch, C., Huvenne, J., 2007. Genetic algorithm optimisation combined with partial least squares regression and mutual information variable selection procedures in near-infrared quantitative analysis of cotton–viscose textiles. Anal. Chim. Acta 595, 72–79. https://doi.org/10.1016/j.aca.2007.03.024.

Fan, S., Huang, W., Guo, Z., Zhang, B., Zhao, C., 2015a. Prediction of soluble solids content and firmness of pears using hyperspectral reflectance imaging. Food Anal. Methods 8, 1936–1946. https://doi.org/10.1007/s12161-014-0079-1.

Fan, S., Huang, W., Guo, Z., Zhang, B., Zhao, C., Qian, M., 2015b. Assessment of influence of origin variability on robustness of near infrared models for soluble solid content of apples. Chin. J. Anal. Chem. 43, 239–244. https://doi.org/10.11895/j.issn.0253-3820.140707.

Fan, S., Zhang, B., Li, J., Huang, W., Wang, C., 2016a. Effect of spectrum measurement position variation on the robustness of NIR spectroscopy models for soluble solids content of apple. Biosyst. Eng. 143, 9–19. https://doi.org/10.1016/j.biosystemseng.2015.12.012.

Fan, S., Zhang, B., Li, J., Liu, C., Huang, W., Tian, X., 2016b. Prediction of soluble solids content of apple using the combination of spectra and textural features of hyperspectral reflectance imaging data. Postharvest Biol. Technol. 121, 51–61. https://doi.org/10.1016/j.postharvbio.2016.07.007.

Fearn, T., 2001. Standardisation and calibration transfer for near infrared instruments: a review. J. Near Infrared Spectrosc. 9, 229–244. https://doi.org/10.1255/jnirs.309.

Ferreira, D.S., Pallone, J.A.L., Poppi, R.J., 2013. Fourier transform near-infrared spectroscopy (FT-NIRS) application to estimate Brazilian soybean [*Glycine max* (L.) Merril] composition. Food Res. Int. 51, 53–58. https://doi.org/10.1016/j.foodres.2012.09.015.

Feudale, R.N., Woody, N.A., Tan, H., Myles, A.J., Brown, S.D., Ferré, J., 2002. Transfer of multivariate calibration models: a review. Chemometr. Intell. Lab. Syst. 64, 181–192. https://doi.org/10.1016/S0169-7439(02)00085-0.

Fu, X., Ying, Y., 2016. Food safety evaluation based on near infrared spectroscopy and imaging: a review. Crit. Rev. Food Sci. Nutr. 56, 1913–1924. https://doi.org/10.1080/10408398.2013.807418.

Giovanelli, G., Sinelli, N., Beghi, R., Guidetti, R., Casiraghi, E., 2014. NIR spectroscopy for the optimization of postharvest apple management. Postharvest Biol. Technol. 87, 13–20. https://doi.org/10.1016/j.postharvbio.2013.07.041.

Guo, Z., Huang, W., Peng, Y., Chen, Q., Ouyang, Q., Zhao, J., 2016. Color compensation and comparison of shortwave near infrared and long wave near infrared spectroscopy for determination of soluble solids content of 'Fuji' apple. Postharvest Biol. Technol. 115, 81–90. https://doi.org/10.1016/j.postharvbio.2015.12.027.

Hu, D., Fu, X., Wang, A., Ying, Y., 2015. Measurement methods for optical absorption and scattering properties of fruits and vegetables. Trans. ASABE 58, 1387–1401. https://doi.org/10.13031/trans.58.11103.

Kunz, M.R., Kalivas, J.H., Erik, A., 2010. Model updating for spectral calibration maintenance and transfer using 1-norm variants of Tikhonov regularization. Anal. Chem. 82, 3642–3649. https://doi.org/10.1021/ac902881m.

Lammertyn, J., Nicolaï, B., Ooms, K., De Smedt, V., De Baerdemaeker, J., 1998. Nondestructive measurement of acidity, soluble solids, and firmness of Jonagold apples using NIR-spectroscopy. Trans. ASAE 41, 1089–1094. https://doi.org/10.1016/j.jfoodeng.2005.06.036.

Li, H., Liang, Y., Xu, Q., Cao, D., 2009. Key wavelengths screening using competitive adaptive reweighted sampling method for multivariate calibration. Anal. Chim. Acta 648, 77–84. https://doi.org/10.1016/j.aca.2009.06.046.

Li, H.D., Xu, Q.S., Liang, Y.Z., 2012. Random frog: an efficient reversible jump Markov Chain Monte Carlo-like approach for variable selection with applications to gene selection and disease classification. Anal. Chim. Acta 740, 20–26. https://doi.org/10.1016/j.aca.2012.06.031.

Li, J., Huang, W., Zhao, C., Zhang, B., 2013. A comparative study for the quantitative determination of soluble solids content, pH and firmness of pears by vis/NIR spectroscopy. J. Food Eng. 116, 324–332. https://doi.org/10.1016/j.jfoodeng.2012.11.007.

Li, J., Huang, W., Chen, L., Fan, S., Zhang, B., Guo, Z., Zhao, C., 2014. Variable selection in visible and near-infrared spectral analysis for noninvasive determination of soluble solids content of 'Ya' pear. Food Anal. Methods 7, 1891–1902. https://doi.org/10.1007/s12161-014-9832-8.

Li, H., Xu, Q., Liang, Y., 2018. libPLS: an integrated library for partial least squares regression and linear discriminant analysis. Chemometr. Intell. Lab. Syst. 176, 34–43. https://doi.org/10.1016/j.chemolab.2018.03.003.

Liang, C., Yuan, H.-f., Zhao, Z., Song, C.-f., Wang, J.-j., 2016. A new multivariate calibration model transfer method of near-infrared spectral analysis. Chemometr. Intell. Lab. Syst. 153, 51–57. https://doi.org/10.1016/j.chemolab.2016.01.017.

Lin, J., 1998. Near-IR calibration transfer between different temperatures. Appl. Spectrosc. 52, 1591–1596. https://doi.org/10.1366/0003702981943095.

Liu, F., Jiang, Y., He, Y., 2009. Variable selection in visible/near infrared spectra for linear and nonlinear calibrations: a case study to determine soluble solids content of beer. Anal. Chim. Acta 635, 45–52. https://doi.org/10.1016/j.aca.2009.01.017.

Liu, D., Sun, D.-W., Zeng, X.-A., 2014. Recent advances in wavelength selection techniques for hyperspectral image processing in the food industry. Food Bioprocess Technol. 7, 307–323. https://doi.org/10.1007/s11947-013-1193-6.

Lu, R., Guyer, D.E., Beaudry, R.M., 2000. Determination of firmness and sugar content of apples using near-infrared diffuse reflectance. J. Texture Stud. 31, 615–630. https://doi.org/10.1111/j.1745-4603.2000.tb01024.x.

Magwaza, L.S., Opara, U.L., Nieuwoudt, H., Cronje, P.J.R., Saeys, W., Nicolaï, B., 2011. Nir spectroscopy applications for internal and external quality analysis of citrus fruit—a review. Food Bioprocess Technol. 5, 425–444. https://doi.org/10.1007/s11947-011-0697-1.

Mehmood, T., Liland, K.H., Snipen, L., Sæbø, S., 2012. A review of variable selection methods in Partial Least Squares Regression. Chemometr. Intell. Lab. Syst. 118, 62–69. https://doi.org/10.1016/j.chemolab.2012.07.010.

Nicolaï, B.M., Beullens, K., Bobelyn, E., Peirs, A., Saeys, W., Theron, K.I., Lammertyn, J., 2007. Nondestructive measurement of fruit and vegetable quality by means of NIR spectroscopy: a review. Postharvest Biol. Technol. 46, 99–118. https://doi.org/10.1016/j.postharvbio.2007.06.024.

Oliveira, G.A.D., Bureau, S., Renard, C.M.G.C., Pereira-Netto, A.B., Castilhos, F.D., 2014. Comparison of NIRS approach for prediction of internal quality traits in three fruit species. Food Chem. 143, 223–230. https://doi.org/10.1016/j.foodchem.2013.07.122.

Ouyang, Q., Yang, Y., Wu, J., Liu, Z., Chen, X., Dong, C., Chen, Q., Zhang, Z., Guo, Z., 2019. Rapid sensing of total theaflavins content in black tea using a portable electronic tongue system coupled to efficient variables selection algorithms. J. Food Compos. Anal. 75, 43–48. https://doi.org/10.1016/j.jfca.2018.09.014.

Peirs, A., Tirry, J., Verlinden, B., Darius, P., Nicola, B.M., 2003. Effect of biological variability on the robustness of NIR models for soluble solids content of apples. Postharvest Biol. Technol. 28, 269–280. https://doi.org/10.1016/S0925-5214(02)00196-5.

Pu, Y.-Y., Sun, D.-W., Riccioli, C., Buccheri, M., Grassi, M., Cattaneo, T.M.P., Gowen, A., 2018. Calibration transfer from micro nir spectrometer to hyperspectral imaging: a case study on predicting soluble solids content of bananito fruit (*Musa acuminata*). Food Anal. Methods 11, 1021–1033. https://doi.org/10.1007/s12161-017-1055-3.

Qin, J., Lu, R., 2008. Measurement of the optical properties of fruits and vegetables using spatially resolved hyperspectral diffuse reflectance imaging technique. Postharvest Biol. Technol. 49, 355–365. https://doi.org/10.1016/j.postharvbio.2008.03.010.

Qin, J., Lu, R., 2009. Monte Carlo simulation for quantification of light transport features in apples. Comput. Electron. Agric. 68, 44–51. https://doi.org/10.1016/j.compag.2009.04.002.

Rinnan, Å., Berg, Fvd., Engelsen, S.B., 2009. Review of the most common pre-processing techniques for near-infrared spectra. Trends Anal. Chem. 28, 1201–1222. https://doi.org/10.1016/j.trac.2009.07.007.

Sjöblom, J., Svensson, O., Josefson, M., Kullberg, H., Wold, S., 1998. An evaluation of orthogonal signal correction applied to calibration transfer of near infrared spectra. Chemometr. Intell. Lab. Syst. 44, 229–244. https://doi.org/10.1016/S0169-7439(98)00112-9.

Travers, S., Bertelsen, M.G., Petersen, K.K., Kucheryavskiy, S.V., 2014. Predicting pear (cv. Clara Frijs) dry matter and soluble solids content with near infrared spectroscopy. LWT-Food Sci. Technol. 59, 1107–1113. https://doi.org/10.1016/j.lwt.2014.04.048.

Vargha, A., Delaney, H.D., 1998. The kruskal-wallis test and stochastic homogeneity. J. Educ. Behav. Stat. 23, 170–192. https://doi.org/10.3102/10769986023002170.

Wang, Y., Veltkamp, D.J., Kowalski, B.R., 1991. Multivariate instrument standardization. Anal. Chem. 63, 2750–2756. https://doi.org/10.1021/ac00023a016.

Wang, H., Peng, J., Xie, C., Bao, Y., He, Y., 2015. Fruit quality evaluation using spectroscopy technology: a review. Sensors 15, 11889–11927. https://doi.org/10.3390/s150511889.

Wold, S., Sjöström, M., Eriksson, L., 2001. PLS-regression: a basic tool of chemometrics. Chemometr. Intell. Lab. Syst. 58, 109–130. https://doi.org/10.1016/S0169-7439(01)00155-1.

Workman, J.J., 2018. A review of calibration transfer practices and instrument differences in spectroscopy. Appl. Spectrosc. 72, 340–365. https://doi.org/10.1177/0003702817736064.

Xu, H., Qi, B., Sun, T., Fu, X., Ying, Y., 2012. Variable selection in visible and near-infrared spectra: application to on-line determination of sugar content in pears. J. Food Eng. 109, 142–147. https://doi.org/10.1016/j.jfoodeng.2011.09.022.

Yao, Y., Chen, H., Xie, L., Rao, X., 2013. Assessing the temperature influence on the soluble solids content of watermelon juice as measured by visible and near-infrared spectroscopy and chemometrics. J. Food Eng. 119, 22–27. https://doi.org/10.1016/j.jfoodeng.2013.04.033.

Zhang, C., Jiang, H., Liu, F., He, Y., 2016. Application of near-infrared hyperspectral imaging with variable selection methods to determine and visualize caffeine content of coffee beans. Food Bioprocess Technol. 10, 213–221. https://doi.org/10.1007/s11947-016-1809-8.

Zhang, B., Dai, D., Huang, J., Zhou, J., Gui, Q., Dai, F., 2017. Influence of physical and biological variability and solution methods in fruit and vegetable quality nondestructive inspection by using imaging and near-infrared spectroscopy techniques: a review. Crit. Rev. Food Sci. Nutr. 1–20. https://doi.org/10.1080/10408398.2017.1300789.

Zhang, F., Zhang, R., Ge, J., Chen, W., Yang, W., Du, Y., 2018. Calibration transfer based on the weight matrix (CTWM) of PLS for near infrared (NIR) spectral analysis. Anal. Methods 10, 2169–2179. https://doi.org/10.1039/c8ay00248g.

Zou, X., Zhao, J., Huang, X., Li, Y., 2007. Use of FT-NIR spectrometry in non-invasive measurements of soluble solid contents (SSC) of 'Fuji' apple based on different PLS models. Chemometr. Intell. Lab. Syst. 87, 43–51. https://doi.org/10.1016/j.chemolab.2006.09.003.