

RSC Advances



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. This *Accepted Manuscript* will be replaced by the edited, formatted and paginated article as soon as this is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

A green method for the quantification of polysaccharides in

Dendrobium officinale

Yong-Huan Yun^{a,1}, Yang-Chao Wei^{a,1}, Xing-Bing Zhao^b, Wei-Jia Wu^b, Yi-Zeng Liang^a,

Hong-Mei Lu^{a,*}

^a College of Chemistry and Chemical Engineering, Central South University, Changsha

410083, PR China

^b Hunan Longshishan Dendrobium Candidum Wall.ex Lindl base Co., Ltd, Changsha 410205,

PR China

Abstract: Polysaccharides is one of active component of *Dendrobium officinale* (*D. officinale*) and its content is used as one of main quality assessment criteria. The polysaccharides quantification existing methods involve sample destruction, tedious sample processing, high cost, and non-environment friendly pretreatment. The aim of this study is to develop a simple, rapid, green and nondestructive analytical method based on near infrared (NIR) spectroscopy and chemometrics methods. A set of 84 *D. officinale* samples from different origins was analyzed by NIR spectroscopy. The potential outlying samples were initially removed from the collected NIR data in two steps by Monte Carlo sampling (MCS) method. The spectral data preprocessing were studied in the construction of partial least squares (PLS) model. To eliminate uninformative variables and improve the performance of model, the pretreated full spectrum was calculated by different wavelength selection methods, including competitive adaptive reweighted sampling (CARS), Monte Carlo-uninformative variable elimination (MC-UVE) and interval random frog (iRF). The selected wavelengths model have met the three following points as: 1) improve the prediction performance; 2) reduce the number of variables; (3) provide a better understanding and interpretation, which proves that it is necessary to conduct wavelength selection in the NIR analytical systems. When comparing the three wavelength selection methods, the results show that CARS has the best performance with the lowest root mean square error of prediction (RMSEP) on the independent test set and

*E-mail: hongmeilu@csu.edu.cn; Tel: +86 731 88830831

¹The first two authors contributed equally to this work

28 least number of latent variables (nLVs). This study demonstrates that the NIR spectral
29 technique with wavelength selection algorithm CARS could be used successfully for
30 quantification of polysaccharides content in *D. officinale*.

31

32 **Keywords:** *Dendrobium officinale*, Near-infrared spectroscopy, Polysaccharides, Green
33 analytical method, Partial least squares (PLS), Competitive adaptive reweighted sampling
34 (CARS)

35 1. Introduction

36 *Dendrobium officinale* (*D. officinale*) is one of the most precious and famous traditional
37 Chinese medicinal material in China. It is claimed to have the function of maintaining gastric
38 tonicity, nourishing Yin and enhancing production of body fluid.^{1,2} It also has been used as a
39 therapeutic agent for curing cataract, throat inflammation, fever and chronic superficial
40 gastritis.³ Many studies suggested that these properties were related to its polysaccharides,
41 one main active component of *D. officinale*.⁴⁻⁷

42 The content of polysaccharides is used as one of quality assessment criteria (no less than
43 0.2500 g glucose per g dry weight) in Chinese pharmacopoeia.⁸ It varies with geographical
44 origin and harvest time. By far, quantification of the polysaccharides in *D. officinale* is mainly
45 performed by the colorimetric method, such as phenol-sulphuric acid method or
46 anthrone-sulphuric acid method. However, those methods involve sample destruction, tedious
47 sample processing, high cost, and non-environment friendly pretreatment, because they
48 require severe conditions of high temperature and strong acid. Therefore, a simple, rapid,
49 green and nondestructive analytical technique is in great demand to determine
50 polysaccharides content in *D. officinale*.

51 Nowadays, as a rapid, green, cost-effective and nondestructive analytical technique, near

52 infrared (NIR) spectroscopy has been widely applied to qualitative and quantitative analysis
53 in agriculture pharmaceuticals, polymer production and food quality evaluation.⁹⁻¹⁸ Recently,
54 NIR spectroscopy has been employed to study traditional Chinese herbs.¹⁹ Some studies on
55 the quantitative analysis of total polysaccharides by NIR were also reported.²⁰⁻²² NIR spectra
56 assess chemical structures through the analysis of the molecular bonds (e.g. C-H, N-H and
57 O-H, which are the primary structural components of organic molecules) in the NIR region,
58 and their characteristic spectra are comprised of different overtone and combination vibrations
59 that are attributable to the molecule's make-up.²³ As a powerful technique, NIR spectroscopy
60 has gained wide acceptance in many fields by virtue of its advantages over other analytical
61 techniques, such as high efficiency, economy, easy operation, and the most salient of its
62 ability to record spectra for solid and liquid samples without any ample preparation. However,
63 NIR spectroscopy usually encounters a collinearity problem because the strongly overlapped
64 and broad absorption bands.²⁴ To address this problem, partial least squares (PLS)²⁵ has been
65 proposed to make a calibration model with NIR data. Typically, the establishment of a
66 calibration model usually covers all the measured wavelengths. It is obvious that such a full
67 spectrum model may contain useless or irrelevant information, which may worsen the
68 predictive ability of the developed model. Liang et al. have demonstrated the importance and
69 necessity of wavelength selection in NIR analytical system.^{26, 27} Many papers have also
70 proved that it is very important and essential to conduct wavelength selection to gain better
71 prediction performance.²⁸⁻³¹ The aim and significance of wavelength selection can be
72 summarized in three points: (1) improving the prediction performance of the calibration
73 model, (2) providing faster and more cost-effective predictors by reducing the curse of

74 dimensionality, (3) providing a better understanding and interpretation of the underlying
75 process that generated the data.^{32, 33}

76 In this work, the first work is to establish the PLS calibration model between the NIR
77 full spectrum data of *D. officinale* and its polysaccharides. Then compare the prediction
78 results of wavelength selection methods and full spectrum. Three recent and often-used
79 wavelength selection methods, including competitive adaptive reweighted sampling (CARS)³⁴,
80 Monte Carlo-uninformative variable elimination (MC-UVE)³⁵ and interval random frog
81 (iRF)³⁶, were employed to compare. Finally determine the best wavelength selection based on
82 the prediction performance and model complexity to develop a calibration model for
83 prediction of polysaccharides in *D. officinale*.

84 2. Materials and methods

85 2.1. Samples collection and reagents

86 A total of 84 *D. officinale* samples were collected from different locations of China in the
87 period of April 2012-April 2014, which are shown in Table 1. It provided a representative set
88 of *D. officinale* consumed in China, which comprised enough variations to make the
89 quantitative model to be robust. Analytical grade *D*-glucose was purchased from
90 Sigma-Aldrich (Sigma, St. Louis, MO, USA). Water was purified by a Milli-Q academic
91 water purification system (Milford, MA, USA). Sulphuric acid of guaranteed reagent grade
92 was purchased from Sinopharm Chemical Reagent Co., Ltd. (Shanghai, China). Other
93 reagents including phenol and ethanol were of analytical grade.

94

95

Insert Table 1

96

97 2.2. Samples preparation and quantitative analysis

98 All the samples were dried at 55 °C in a forced-draught oven from Shanghai Pharmacy
99 Machine Co. (Shanghai, China). After brushing off soil dust from the surface, the samples
100 were ground to fine pieces with a blender and screened through a 60-mesh sieve (particle size
101 ≤ 0.2 mm). These sieved powders were used for further analysis.

102 *D. officinale* polysaccharides content was firstly measured with the phenol-sulphuric
103 acid method provided by Chinese pharmacopoeia (State Pharmacopoeia Committee 2010).
104 Glucose calibration curve was firstly prepared. The glucose (0.255 g) dried to constant weight
105 at 105 °C was placed in 250 ml volumetric flask, added water to obtain 100 µg/ml solution.
106 Accurately draw 0.0, 0.2, 0.4, 0.6, 0.8, 1.0 ml of glucose solution in 10 ml test tube with lid
107 respectively, added water to 1 ml. Then added 1 ml of 5% phenol solution, mixed, quickly
108 added sulphuric acid 5.0 ml, shook, bathed in 90 °C water for 20 min, put in an ice bath for 5
109 min. A BTT miniature array spectrophotometer (B&W Tek, Newark, DE, USA) equipped
110 with glass or quartz cells of 1 cm path length was used for measurement of absorbance spectra.
111 A Lenovo personal computer was used to control the spectrometer and collect data via a
112 BWSpec4 Software. Absorbance unit was recorded at wavelength 488.02 nm. The calibration
113 curve was made according to absorbance unit and glucose concentration.

114 Polysaccharides measurement was as follow. An accurately weighted, powdered *D.*
115 *officinale* sample (0.3 g) was loaded into a standard apparatus set, refluxed for 2 h with 200
116 ml water. Subsequently, the sample was cooled to room temperature and transferred to a 250
117 ml volumetric flask, added water to the scale, shook and filtered. Then 2 ml of filtrate was

118 precipitated by ethanol (10 ml) at 4 °C, followed by centrifugation for 30 min at 4000 r/min.
119 The precipitate was washed twice with 8 ml of 80% ethanol. The precipitate obtained after
120 filtering was dissolved in water and collected in a 25 ml volumetric flask. The following
121 operation was based on the calibration curve of glucose aforementioned. Results were
122 expressed as grams of glucose equivalents per gram of dry weight (g glucose per g DW)
123 through the calibration curve with glucose. Each sample was determined in triplicate, and the
124 mean of three measurements was used for further analysis.

125 2.3.NIR spectroscopy measurement

126 With the Integrating Sphere module of the Antaris II Fourier transform near infrared (FT-NIR)
127 analyzer (Thermo Scientific, Madison, USA), the NIR diffuse reflection spectra were
128 collected from 10,000 to 4000 cm^{-1} (1557 wavelength points). The reference spectrum is the
129 gold foil. Each sample was scanned for 32 times with a resolution of 8 cm^{-1} using a
130 background of the air and the average of spectrum of 32 scans was taken as one result. The
131 environment temperature was controlled at 25 ± 1 °C with an air conditioner.

132 The standard sample cup was used to collect spectra of *D. officinale* samples. It was the
133 standard accessory as sample's holder, specifically designed by Thermo Electron Co.. About
134 0.5 g of the sample in powder form was filled into the sample cup in the standard procedure.
135 In order to avoid errors from uneven samples, the sample cup was rotated 120° to record
136 another spectrum after each record. Each sample was collected three times. The mean of three
137 spectra which were collected from the same sample was used for the following analysis.

138 A set of 84 *D. officinale* samples from different origins in China was analyzed by NIR
139 spectroscopy. The generated spectra of 84 samples are shown in Fig. 1(a).

140 Insert Fig. 1

141 2.4. Outlier detection and spectral data preprocessing

142 Constructing a high-quality model depends on the execution of several steps. One important
143 step is outlier detection. The step of outlier detection should be prior to establish the
144 calibration model. Outliers are abnormal ones in some sense. They may present
145 non-representative samples that could introduce great errors to a model. In this work, a novel
146 strategy which was termed as the Monte Carlo sampling (MCS) method was used for the
147 outlier detection. According to the method, there may be three types of outliers.³⁷ The first
148 one is the outliers in the dependent variable y direction. It breaks away from the normal
149 distribution of y and will cause a large error sum of squares. The second one is the outliers in
150 the predictor or independent variable X direction. This sort of outliers is far away from the
151 main body of the samples. The third type of outliers, so called outliers towards the model, can
152 be found only after building the regression model. They represent a different relationship
153 between X and y . In the MCS method, the number of latent variables (nLVs) was firstly
154 determined using cross-validation in PLS. With the help of the MCS method, the whole data
155 set was randomly divided into two parts, the calibration set and independent test set,
156 respectively. After that the calibration set was used to establish the model using the optimal
157 nLVs. The independent test set was used for prediction. The prediction error would be
158 obtained for each test sample. This cycle was executed in 1000 times. Finally, the prediction
159 error distribution for each sample was obtained. The histograms of these distributions were
160 plotted and their statistic features were used to detect the outliers.

161 In addition to useful information, spectral signal contains systematic noise, such as

162 baseline variation, sample background, light scattering and so on.³⁸ In order to build a robust
163 and reliable model, some preprocess must be taken to weaken and eliminate interference in
164 spectra. In this study, eight different signal pre-treatment methods were evaluated and
165 compared, including multiplicative scattering correction (MSC), standard normal
166 transformation (SNV), first and second derivatives computed by Savitzky-Golay (S-G)
167 method, and the combinations of MSC (or SNV) with the derivatives. MSC is an important
168 procedure for the correction of scatter light caused by different particle sizes. It is also used to
169 correct the additive and multiplicative effects in the spectra. SNV is a mathematical
170 transformation method of the log (1/R) spectra used to remove slope variation and to correct
171 for scatter effects.^{39, 40} Compared to SNV, first and second derivative are used to reduce peak
172 overlap and remove constant and linear baseline drift, respectively. Thus, they are often used
173 to eliminate baseline drifts and enhance small spectral differences between samples.⁴¹

174 2.5. Multivariate calibration methods

175 2.5.1. Partial least squares (PLS) regression

176 PLS is a commonly used multivariate calibration method. It investigates the fundamental
177 relations between the response vector (the properties of interest), y , and the spectral data
178 matrix, X . In this method, data is compressed into orthogonal factors, which have similar
179 properties to PCs in principal component analysis (PCA).^{42, 43} Here, the purpose of PLS is to
180 establish a regression model to make the prediction of chemical constituent concentrations. It
181 extends and improves the potential application of spectroscopy technique in food industry by
182 extracting features from spectra.⁴⁴

183 Three different wavelength selection methods combined with PLS, including competitive

184 adaptive reweighted sampling (CARS), Monte Carlo-uninformative variable elimination
185 (MC-UVE) and interval random frog (iRF) were employed to compare and determine the
186 effective wavelengths.

187 CARS³⁴ is a novelty variable selection algorithm, which is similar to the “survival of the
188 fittest” principle in Darwin’s Evolution Theory. The wavelengths with large absolute
189 coefficients that selected by CARS were defined as the key wavelengths. In each sampling
190 run, CARS contains four successive steps: (1) use MC sampling method to select modeling
191 samples randomly; (2) employ exponentially decreasing function (EDF) to remove the
192 wavelengths which are of relatively small absolute regression coefficients by force; (3) adopt
193 adaptive reweighted sampling (ARS) to realize a competitive selection of wavelengths; (4)
194 employ cross-validation to evaluate the subset and finally to choose the subset with the lowest
195 root mean squared error of cross validation (RMSECV). For CARS, the number of sampling
196 run was set to 100.

197 MC-UVE³⁵ is a useful variable selection algorithm, which combined Monte Carlo (MC)
198 strategy with uninformative variable elimination (UVE) method. The MC-UVE method builds
199 a large number of PLS sub-models with randomly selected calibration samples at first, and
200 each variable is evaluated with a stability of the corresponding regression coefficient.
201 Variables with poor stability are known as uninformative variable and eliminated. The number
202 of MC sampling run was set to 1000 in this study.

203 iRF³⁶ is a wavelength interval selection method that considers the continuity of spectra.
204 It is based on random frog⁴⁵ that employs reversible jump Markov Chain Monte Carlo
205 (RJMCMC)-like search algorithm in the model space through both fixed-dimensional and

206 trans-dimensional between different models. The objective function is to find the subset
207 which has the maximum regression coefficient. Spectra are first divided into sub-intervals on
208 the whole spectra by a moving window of a fixed width and thus it can obtain all the possible
209 continuous spectral intervals. Each interval is regarded as the variable and then is input into
210 the RJMCMC algorithm. A pseudo-MC MC chain is used to compute selection probability of
211 each interval, and then rank all the intervals based on the selection probability. Afterwards,
212 choose the best intervals with the lowest RMSECV. In this work, with 1557 full spectral
213 points, the width of the interval is set to 20 resulting in 1538 intervals in total and each
214 interval has 20 variables.

215 2.6. Data division and model performance evaluation

216 After sample outlier detection and the best pretreatment selection, the next step was to divide
217 the whole data set into calibration and independent test set, which are used to build and
218 validate the model, respectively. To assure that the division of calibration set and independent
219 test set was well proportioned, a procedure based on the Duplex algorithm was used to split
220 the data set.^{46,47}

221 In this work, selection was performed using a splitting ratio of 2:1 (50 samples were
222 taken into calibration set, and the remaining 25 samples served for the independent test set).
223 The statistical values of the polysaccharides content in calibration and independent test sets
224 are listed in Table 2. After the division, the content values in the calibration and independent
225 test sets covered a wide range, which is helpful to develop a robust model.

226 The calibration set is used for building a PLS model and wavelength selection, and the
227 independent test set is used for external validation. The optimal nLVs on the calibration set

228 was determined by 10-fold cross validation as the maximum nLVs was set to 15. The built
229 model was then used to predict the calibration set and test set, generating with a root mean
230 squared error of fitting on the calibration set (RMSEC) value and a root mean squared error of
231 prediction on the independent test set (RMSEP) value. Thus, RMSEC, R_{cal}^2 , RMSEP and
232 R_{pre}^2 (R^2 on the test set), were employed to assess the performance of the generated model.
233 RMSECV and R_{cv}^2 were used to determine spectral data preprocessing method.

234 2.7. Software

235 NIR spectra were collected using an Antaris II FT-NIR spectrometer. The instrument was
236 equipped with the spectral acquisition software, called “Results”. After NIR spectra were
237 collected, spectra were imported directly into MATLAB (Version 2013A, the MathWorks,
238 Inc) on a general-purpose computer with Intel® Core® i5 3.2GHz CPU and 3GB RAM, with
239 operating system Microsoft Windows XP. The spectral data preprocessing and multivariate
240 calibration were implemented by the written codes in MATLAB, which can be downloaded
241 freely in the website: <http://www.libpls.net/>.

242 3. Results and discussion

243 3.1. Polysaccharides content measurement

244 The polysaccharides content in all 84 samples were determined by the reference method (see
245 Section 2.2). The glucose calibration equation was $Y=0.0094X+0.0016$, $R^2=0.9998$, which
246 showed a good linear relationship between 0.0 and 0.1 mg/ml glucose content and absorbance
247 unit. After the outlier removal (see Section 3.2.1), there are 75 samples for PLS modeling.
248 The polysaccharides contents in the 75 *D. officinale* samples were calculated according to
249 glucose calibration equation and absorbance unit, and were shown in Table 2. It was

250 0.4006±0.1329 g glucose per g DW. The polysaccharides content of some samples were less
251 than 0.2500 g glucose per g DW, the threshold value restricted by pharmacopoeia. Therefore,
252 it was necessary to monitor the quality of *D. officinale*.

253

254 Insert Table 2

255 3.2. Model building

256 3.2.1. Deletion of outlying samples

257 The results of outlier detection by the MCS method are shown in Fig. 2. From Fig. 2(a), the
258 three samples (12, 28, 29) in top left area are outliers in X direction which have a large
259 standard deviation of prediction errors, and the lower right one gives two outliers (57,70) in y
260 direction, which have a large mean value of prediction errors. As mentioned above, the
261 division of samples is based on MCS method, so the first result may be not really show all the
262 outliers. In order to further detect the potential outliers, the MCS method was run for the
263 remaining samples once again after the last outlier detection. Similar to Fig. 2(a), Fig. 2(b)
264 shows the result for the data set including two different types of outliers. From this plot, it can
265 be seen that the entire datum is clearly divided into three parts, and different type of outliers
266 compactly clustered together, respectively. The result shows that two samples (69, 71) in the
267 lower right area are outliers in y direction, and the top right two samples (27, 47) are outliers
268 both in the X and y directions. From Fig. 2(b), the four samples which are not shown
269 significantly in the first step are far away from the main body of the data with higher mean
270 values or higher deviations of prediction errors. The MCS method was first used in two steps

271 to reveal the potential outliers in this study. After the removal of outliers, the remaining 75
272 samples were used for the following analysis.

273 Insert Fig. 2

274

275 3.2.2. Selection of spectral data preprocessing methods

276 PLS full spectrum model were developed with different data preprocessing methods. A
277 10-fold cross-validation was used to select the nLVs and the most suitable spectral data
278 preprocessing using the whole samples (75 samples). The spectral preprocessing was
279 optimized based on the lowest RMSECV, highest R_{cv}^2 and few nLVs. According to the Table
280 3, the best one was found to be built with data pretreated by SNV combined with the SG 1st
281 derivative (11 points, 3rd order polynomial) and as it has the lowest RMSECV, 0.0543 highest
282 R_{cv}^2 , 0.8309 and only 6 nLVs, which is consistent with the work from.²¹ When there are
283 overlapping peaks in the original NIR spectra, the SNV 1st derivative for data pretreatment is
284 usually useful to enhance the resolution, correct for scatter effects and for the baseline
285 correction. The reason might be that the SG 1st derivative calculation removed both additive
286 and multiplicative effects in the spectra. The preprocessed spectra are shown in Fig. 1(b). It
287 can be seen that most absorbance values were zero approximately, and the overlapping peaks
288 and baseline effect were removed. The spectral differences of the samples were observed in
289 several different regions of around 4000-4300 cm^{-1} and 5750 cm^{-1} .

290

291 Insert Table 3

292

293 3.2.3. Full spectrum and wavelength selection models

294 There are 1557 variables in the NIR full spectral data. The full spectrum calibration model on
295 the calibration set was developed and then used to make a prediction for validation on the
296 independent test set. In addition, iRF, CARS and MC-UVE were employed to select
297 wavelengths. All methods were conducted 100 times to get the best one because Monte Carlo
298 sampling they used would lead to different results in each time.

299 When compared to the full spectrum model, the selected wavelengths model should meet
300 the three following points as: 1) improve the prediction performance; 2) reduce the number of
301 wavelengths; (3) provide a better understanding and interpretation. The calibration and
302 validation results of full spectrum and wavelength selection methods are shown in Table 4.

303 For the prediction of the full spectrum model, RMSEP and R_{pre}^2 are 0.0542 and 0.7978,
304 respectively. The nLVs is 10. It can be observed that all the wavelength selection methods
305 perform better than full spectrum PLS model based on the RMSEP, R_{pre}^2 and nLVs, which
306 satisfies the first point that improve the prediction performance and. Moreover, the number of
307 selected wavelengths by the CARS, MC-UVE and iRF, are 39, 339 and 364, which are also
308 much less than full spectrum with 1557 wavelengths. Thus, it demonstrates that the model can
309 obtain good prediction performance when eliminating the variables that are uninformative and
310 have irrelevant information.

311 CARS and MC-UVE are the discrete wavelength selection methods, while iRF is
312 wavelength interval selection method. All of them are based on the PLS regression coefficient.
313 Here we do not aim to prove that whether discrete wavelength selection or wavelength
314 interval selection method is better. The performances of all the wavelength selection methods

315 are data dependent. In this work, for the determination of the polysaccharides content in *D.*
316 *officinale*, when in comparison of three wavelength selection methods, the overall results
317 indicated that CARS obtains the best prediction performance with the lowest RMSEP and
318 R^2_{pre} . The least nLVs also indicate that CARS can establish the most parsimonious PLS model.
319 The reason may be that there are too many irrelevant variables in the full spectral data. CARS
320 is an effective procedure to eliminate uninformative variables and improved the predictive
321 precision of the model. Based on exponentially decreasing function, CARS firstly eliminated
322 large number of wavelengths in the first stage and then in a refined way to select wavelength.
323 Although CARS runs fast, it is not stable. Thus, CARS should be conducted many times to
324 obtain the best one.

325 As Polysaccharides belong to carbohydrates, it contains aliphatic cyclic groups with
326 attached OH groups and either linkages. In order to understand and interpret the selected
327 wavelengths by all the wavelength selection methods for polysaccharides, they are displayed
328 in Fig. 3. The wavelengths selected by MC-UVE are very scattered, resulting in that
329 MC-UVE performs a little better than full spectrum model. CARS and iRF have a lot of
330 common selected regions. As CARS performs the best in this work, the interpretation of
331 selected wavelengths focuses on CARS. We can see that the selected wavelengths by CARS
332 are mostly concentrated on the region of 4000-4200 cm^{-1} , 4300-4450 cm^{-1} , 4700-5250 cm^{-1} ,
333 5750-7300 cm^{-1} , 7900-8950 cm^{-1} , 9000-10000 cm^{-1} . The absorption at 4000-4200 cm^{-1} is
334 related to C-H stretching and C-C and C-O-C stretching combination.⁴⁸ 4300-4450 cm^{-1} is
335 corresponding to C-H stretching and CH₂ deformation combination, while 4700-5100 cm^{-1} is
336 corresponding to O-H bending, O-H stretching, C-O stretching combination and HOH

337 bending combination. ⁴⁸ 5750-7300 cm⁻¹ is related to the first overtone of C-H stretching. ⁴⁸
338 7900-8950 cm⁻¹ could be attributed to the first overtone of O-H in polysaccharides. ⁴⁹
339 9000-1000 cm⁻¹ is corresponding to the second overtone of O-H. ⁵⁰

340 From the above points, it can be proved that wavelength selection is necessary and
341 essential in multivariate calibration for the NIR analytical system.

342

343 Insert Fig. 3

344 Insert Table 4

345

346 Fig. 4 shows the correlation between the values determined by phenol-sulphuric acid
347 method and the values predicted by the NIR full spectrum model (Fig. 4a) and CARS (Fig.
348 4b). The blue and red circles correspond to the calibration and independent test set,
349 respectively. The diagonal line represents the ideal results. The closer the points are to the
350 diagonal line, the better the model is. It can be found that the samples are distributed more
351 closely to the diagonal line in Fig. 4b, which shows a good spectral analysis performance of
352 CARS. The results demonstrate the feasibility to use NIR spectroscopy combined with CARS
353 for determination of the polysaccharides content of *D. officinale*.

354

355 Insert Fig. 4

356

357 4. Conclusions

358 In this study, a rapid, cost-effective and non-destructive technique, namely NIR, coupled with

359 multivariate calibration method, PLS, for the determination of the polysaccharides content in
360 *D. officinale* was demonstrated. The integrated step including outlier detection, data
361 preprocessing and establishment of calibration model were introduced. Comparing with the
362 full spectrum model, three recent and often-used wavelength selection methods, including
363 MC-UVE, CARS and iRF, were employed to demonstrate the good prediction performance,
364 reduction of the number of variables and a better understanding and interpretation of selected
365 wavelengths. Thus, wavelength selection is necessary in the multivariate calibration model in
366 NIR analytical system. When comparing the three wavelength selection methods, CARS
367 performs the best with the lowest RMSEP, highest R_{pre}^2 and fewest number of latent
368 variables.

369 Therefore, NIR could provide a fast and green alternative to classical reference methods,
370 as it dramatically reduces analysis time without any chemical reagents. The established
371 method will significantly improve the efficiency of quality control. Furthermore, the future
372 work is to develop similar NIR spectroscopy calibration models coupled with CARS
373 algorithm for predicting additional components in *D. officinale*, such as alkaloid,
374 sesquiterpenoid and aromatic compound. It should be noted that more work should be paid
375 attention to robustness of calibration models by collecting more samples and introducing
376 more wavelength selection methods.

377

378 **Acknowledgement**

379 The authors gratefully thank National Natural Science Foundation of China for support
380 of the projects (No. 21175157, 21375151 and 21275164), China Hunan Provincial science

381 and technology department for support of the project (No. 2012FJ4139), Central South
382 University for special support of the basic scientific research project (No. 2010QZZD007),
383 also supported by the Fundamental Research Funds for the Central Universities of Central
384 South University (grants no. 2014zzts014).

385

386

References

387

1. X. Chen, F. Wang, Y. Wang, X. Li, A. Wang, C. Wang and S. Guo, *Sci. China Life Sci.*, 2012, **55**, 1092-1099.

388

2. T. Ng, J. Liu, J. Wong, X. Ye, S. Wing Sze, Y. Tong and K. Zhang, *Appl. Microbiol. Biotechnol.*, 2012, **93**, 1795-1803.

389

3. X. Q. Zha, J. P. Luo and P. Wei, *S. Af. J. Bot.*, 2009, **75**, 276-282.

390

4. X. Xing, S. W. Cui, S. Nie, G. O. Phillips, H. Douglas Goff and Q. Wang, *Bioactive Carbohydrates and Dietary Fibre*, 2013, **1**, 131-147.

391

5. L. Xia, X. Liu, H. Guo, H. Zhang, J. Zhu and F. Ren, *J. Func. Foods.*, 2012, **4**, 294-301.

392

6. L.-H. Pan, X.-F. Li, M.-N. Wang, X.-Q. Zha, X.-F. Yang, Z.-J. Liu, Y.-B. Luo and J.-P. Luo, *Int. J. Biol. Macromol.*, 2014, **64**, 420-427.

393

7. L.-Z. Meng, G.-P. Lv, D.-J. Hu, K.-L. Cheong, J. Xie, J. Zhao and S.-P. Li, *Molecules*, 2013, **18**, 5779-5791.

394

8. S. P. Committee, *Beijing- People's Medical Publishing House*, 2010.

395

9. Y. Ozaki, W. F. McClure and A. A. Christy, *Near-infrared spectroscopy in food science and technology*, John Wiley & Sons, 2006.

396

10. P. Williams and K. Norris, *Near-infrared technology in the agricultural and food industries*, American Association of Cereal Chemists, Inc., 1987.

397

11. O. Escuredo, M. Carmen Seijo, J. Salvador and M. Inmaculada González-Martín, *Food Chem.*, 2013, **141**, 3409-3414.

398

12. M. J. De la Haba, A. Garrido-Varo, J. E. Guerrero-Ginel and D. C. Pérez-Marín, *J. Agr. Food. Chem.*, 2006, **54**, 7703-7709.

399

13. J. Sarembaud, G. Platero and M. Feinberg, *Food Anal. Methods*, 2008, **1**, 227-235.

400

14. E. Fernández-Ahumada, A. Garrido-Varo, J. Guerrero-Ginel, A. Wubbels, C. Van der Sluis and J. Van der Meer, *J. Near. Infrared. Spec.*, 2006, **14**, 27-35.

401

15. L. León, A. Garrido-Varo and G. Downey, *J. Agr. Food. Chem.*, 2004, **52**, 4957-4962.

402

16. M. Manley, G. du Toit and P. Geladi, *Anal. Chim. Acta.*, 2011, **686**, 64-75.

403

17. D. Wu, H. Shi, S. Wang, Y. He, Y. Bao and K. Liu, *Anal. Chim. Acta.*, 2012, **726**, 57-66.

404

18. W. Li, Y. Wang and H. Qu, *Vib. Spectrosc.*, 2012, **62**, 159-164.

405

19. Q. Luo, Y. Yun, W. Fan, J. Huang, L. Zhang, B. Deng and H. Lu, *RSC Advances*, 2015, **5**, 5046-5052.

406

20. H. Yan, B.-x. Han, Q.-y. Wu, M.-z. Jiang and Z.-z. Gui, *Spectrochim. Acta. A.*, 2011, **79**, 179-184.

407

21. Y. Chen, M. Xie, H. Zhang, Y. Wang, S. Nie and C. Li, *Food Chem.*, 2012, **135**, 268-275.

408

22. Y. Wei, W. Fan, X. Zhao, W. Wu and H. Lu, *Anal. Lett.*, 2014, **48**, 817-829.

409

23. J. R. Lucio-Gutiérrez, J. Coello and S. Maspocho, *Food Research International*, 2011, **44**, 557-565.

410

24. M. Blanco, J. Cruz and M. Bautista, *Anal. Bioanal. Chem.*, 2008, **392**, 1367-1372.

411

25. S. Wold, M. Sjöström and L. Eriksson, *Chemometr. Intell. Lab. syst.*, 2001, **58**, 109-130.

412

26. H.-D. Li, Y.-Z. Liang, X.-X. Long, Y.-H. Yun and Q.-S. Xu, *Chemometr. Intell. Lab. syst.*, 2013, **122**, 23-30.

413

27. Y.-H. Yun, Y.-Z. Liang, G.-X. Xie, H.-D. Li, D.-S. Cao and Q.-S. Xu, *Analyst*, 2013, **138**, 6412-6421.

414

28. J. H. Kalivas, N. Roberts and J. M. Sutter, *Anal. Chem.*, 1989, **61**, 2024-2030.

415

29. D. Jouan-Rimbaud, B. Walczak, D. L. Massart, I. R. Last and K. A. Prebble, *Anal. Chim. Acta*, 1995, **304**, 285-295.

416

30. L. Xu and I. Schechter, *Anal. Chem.*, 1996, **68**, 2392-2400.

417

31. C. H. Spiegelman, M. J. McShane, M. J. Goetz, M. Motamedi, Q. L. Yue and G. L. Coté, *Anal. Chem.*, 1998, **70**, 35-44.

418

32. A. Lorber and B. R. Kowalski, *J. Chemometr.*, 1988, **2**, 67-79.

419

33. I. Guyon and A. Elisseeff, *J. Mach. Learn. Res.*, 2003, **3**, 1157-1182.

420

- 421 34. H. Li, Y. Liang, Q. Xu and D. Cao, *Anal. Chim. Acta.*, 2009, **648**, 77-84.
- 422 35. W. Cai, Y. Li and X. Shao, *Chemometr. Intell. Lab. syst.*, 2008, **90**, 188-194.
- 423 36. Y. H. Yun, H. D. Li, L. R. E. Wood, W. Fan, J. J. Wang, D. S. Cao, Q. S. Xu and Y. Z. Liang, *Spectrochim. Acta. A.*, 2013, **111**, 31~36.
- 424 37. D. S. Cao, Y. Z. Liang, Q. S. Xu, H. D. Li and X. Chen, *J. Comput. Chem.*, 2010, **31**, 592-602.
- 425 38. T. Naes, T. Isaksson, T. Fearn and T. Davies, *A user friendly guide to multivariate calibration and classification*, NIR publications, 2002.
- 426 39. Y. He, X. Li and X. Deng, *J. Food Eng.*, 2007, **79**, 1238-1242.
- 427 40. Q. Chen, J. Zhao and H. Lin, *Spectrochim. Acta. A.*, 2009, **72**, 845-850.
- 428 41. Q. Chen, J. Zhao, C. Fang and D. Wang, *Spectrochim. Acta. A.*, 2007, **66**, 568-574.
- 429 42. W. Dong, Y. Ni and S. Kokot, *J. Agr. Food. Chem.*, 2013, **61**, 540-546.
- 430 43. D. Wu, Y. He, P. Nie, F. Cao and Y. Bao, *Anal. Chim. Acta.*, 2010, **659**, 229-237.
- 431 44. L. Xie, X. Ye, D. Liu and Y. Ying, *Food Research International*, 2011, **44**, 2198-2204.
- 432 45. H.-D. Li, Q.-S. Xu and Y.-Z. Liang, *Anal. Chim. Acta*, 2012, **740**, 20-26.
- 433 46. R. D. Snee, *Technometrics*, 1977, **19**, 415-428.
- 434 47. M. Bevilacqua, R. Bucci, A. D. Magri, A. L. Magri and F. Marini, *Anal. Chim. Acta.*, 2012, **717**, 39-51.
- 435 48. J. Workman Jr and L. Weyer, *Practical guide to interpretive near-infrared spectroscopy*, CRC press, 2007.
- 436 49. W. Z. Lu, H. F. Yuan, G. T. Xu and D. M. Qiang, eds., *The Technology of Modern Near Infrared Spectral Analysis*, China Petrochemical Press,
437 Beijing, 2000.
- 438 50. X. B. Zou, J. W. Zhao, M. J. W. Povey, M. Holmes and H. P. Mao, *Anal. Chim. Acta.*, 2010, **667**, 14-32.
- 439
- 440

441 **Tables**442 Table 1. *D. officinale* samples information.

Sample no.	Origin	Collected time
1-6	Yunnan	Feb.,2013-Mar.,2013
7-12	Zhejiang	Apr.,2012-Oct.,2012
13-14	Hunan	Sep.,2012-Jul.,2013
15-16	Zhejiang	Jul., 2013-Aug.,2013
17-20	Henan	Jul.,2013-Aug.,2013
21-32	Hunan	Dec.,2013
33-49	Hunan	Feb.,2014
50-53	Yunnan	Feb.,2014
54-61	Yunan	Mar.,2013
62-67	Zhejiang	Apr.,2012-Jul.,2012
68-84	Hunnan	Apr.,2014

443

444

445 Table 2. The *D. officinale* polysaccharides content measured with the phenol-sulphuric acid
446 method and the number of *D. officinale* samples used in dataset.

447

Data set	Number	Max (g glucose per g DW ^b)	Min (g glucose per g DW)	Mean±S.D ^a (g glucose per g DW)
Total	75	0.7063	0.1863	0.4006±0.1329
Calibration set	50	0.7063	0.1863	0.4111±0.1302
Test set	25	0.6952	0.1925	0.3796±0.1385

448 ^aS.D is standard deviation.449 ^bDW is dry weight.

450

451 Table 3. The 10-fold cross-validation results by PLS with different data preprocessing
452 methods.

453

454

Pretreatment	nLVs	RMSECV	R_{cv}^2
Original	14	0.0558	0.8211
Smooth+MSC	11	0.0539	0.8330
Smooth+SNV	6	0.0585	0.8036
SG 1st	12	0.0540	0.8330
SG 2nd	4	0.0651	0.7571
MSC+SG 1st	6	0.0543	0.8308
MSC+SG 2nd	6	0.0619	0.7800
SNV+SG 1st	6	0.0543	0.8309
SNV+SG 2nd	6	0.0619	0.7802

458

459

460

461 Table 4. Results of *D. officinale* polysaccharides content by PLS models based on different
 462 wavelength selection methods.

463

	Full spectrum	CARS	MC-UVE	iRF
464 N.W ^a	1557	39	339	364
465 nLVs	10	8	10	9
RMSECV	0.0549	0.0156	0.0260	0.0423
466 R^2_{cv}	0.8397	0.9872	0.9640	0.9048
RMSEC	0.0101	0.0096	0.0010	0.0025
467 R^2_{cal}	0.9946	0.9952	0.9999	0.9997
RMSEP	0.0542	0.0468	0.0533	0.0486
468 R^2_{pre}	0.7978	0.8495	0.8044	0.8373

469 ^aN.W is the number of wavelengths.

470

471 **Figure Captions**

472 Fig. 1. (a)The raw NIR spectra of 84 *D. officinale* samples; (b) Preprocessed spectra by
473 SNV+SG 1st derivative of 75 *D. officinale* samples.

474 Fig. 2. The results of variance of residuals versus mean of residuals on the
475 polysaccharides content of *D. officinale* samples. (a) The first step of MCS; (b) The second
476 step of MCS.

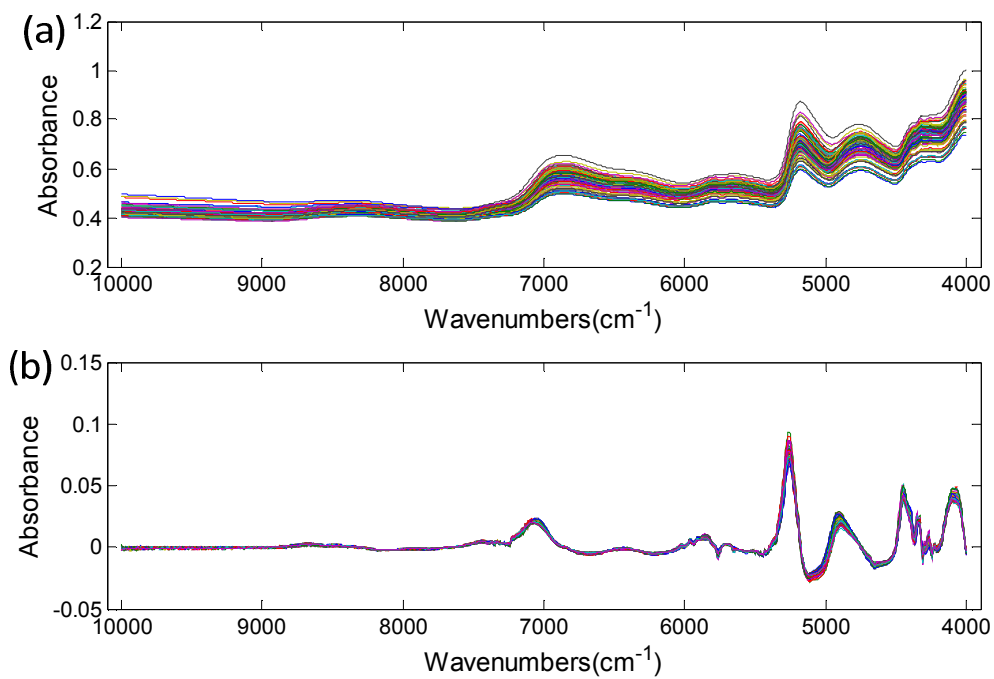
477 Fig. 3. The distribution of the selected variables obtained by different wavelength
478 selection methods.

479 Fig. 4. The correlation between predicted value and measured value of polysaccharides
480 content based on (a) The full spectra PLS model; (b) 39 selected wavelengths by CARS.

481

482

483



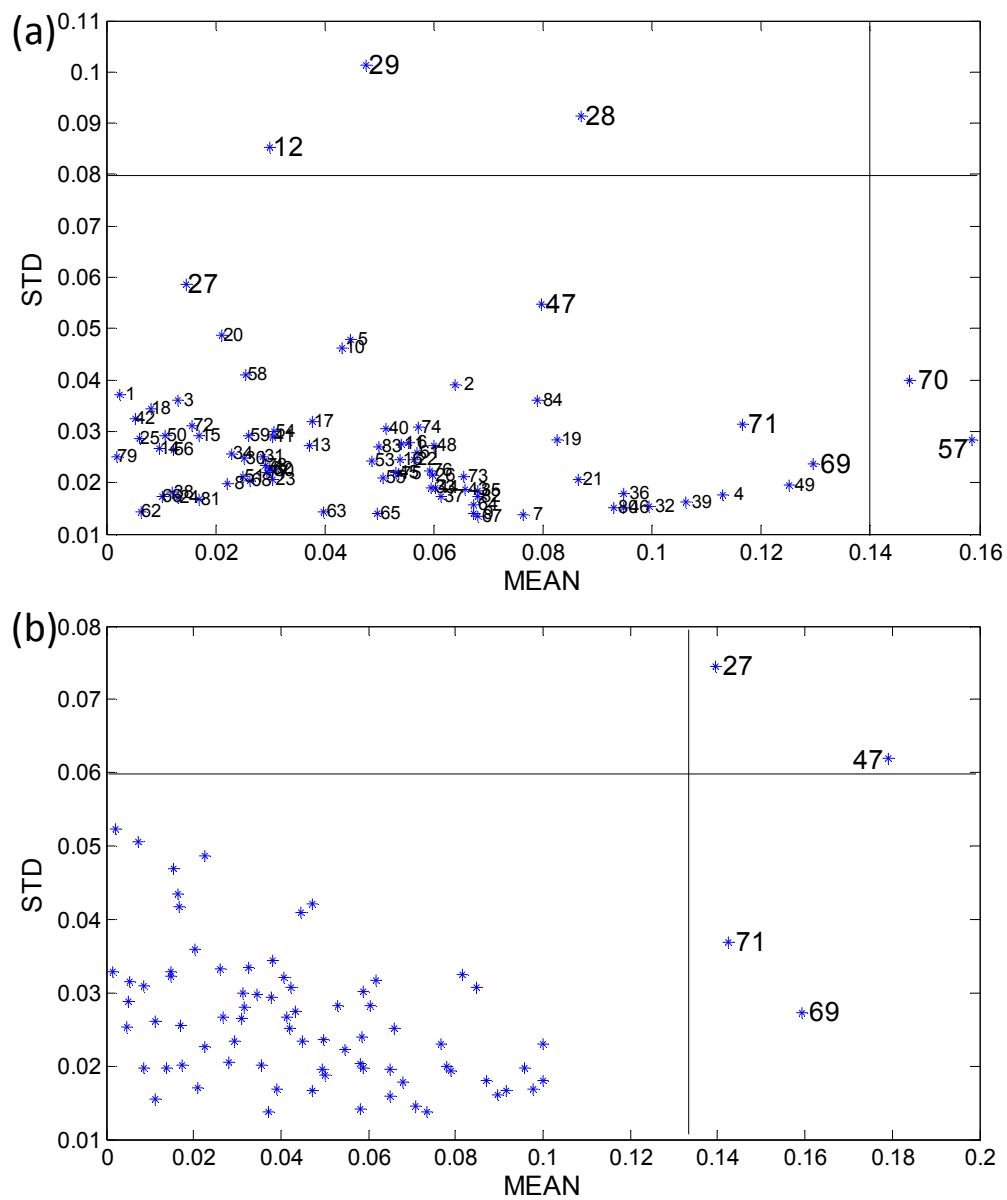
484

485 Fig. 1. (a)The raw NIR spectra of 84 *D. officinale* samples; (b) Preprocessed spectra by

486

SNV+SG 1st derivative of 75 *D. officinale* samples.

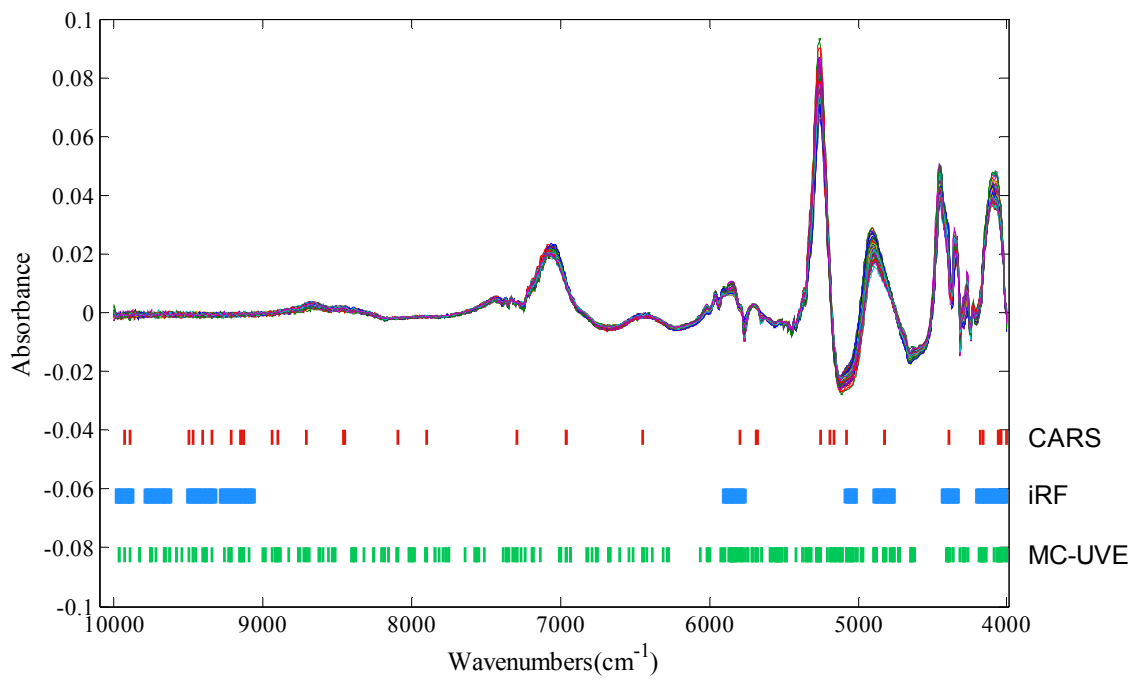
487



488

489 Fig. 2. The results of variance of residuals versus mean of residuals on the polysaccharides

490 content of *D. officinale* samples. (a) The first step of MCS; (b) The second step of MCS.



491

492

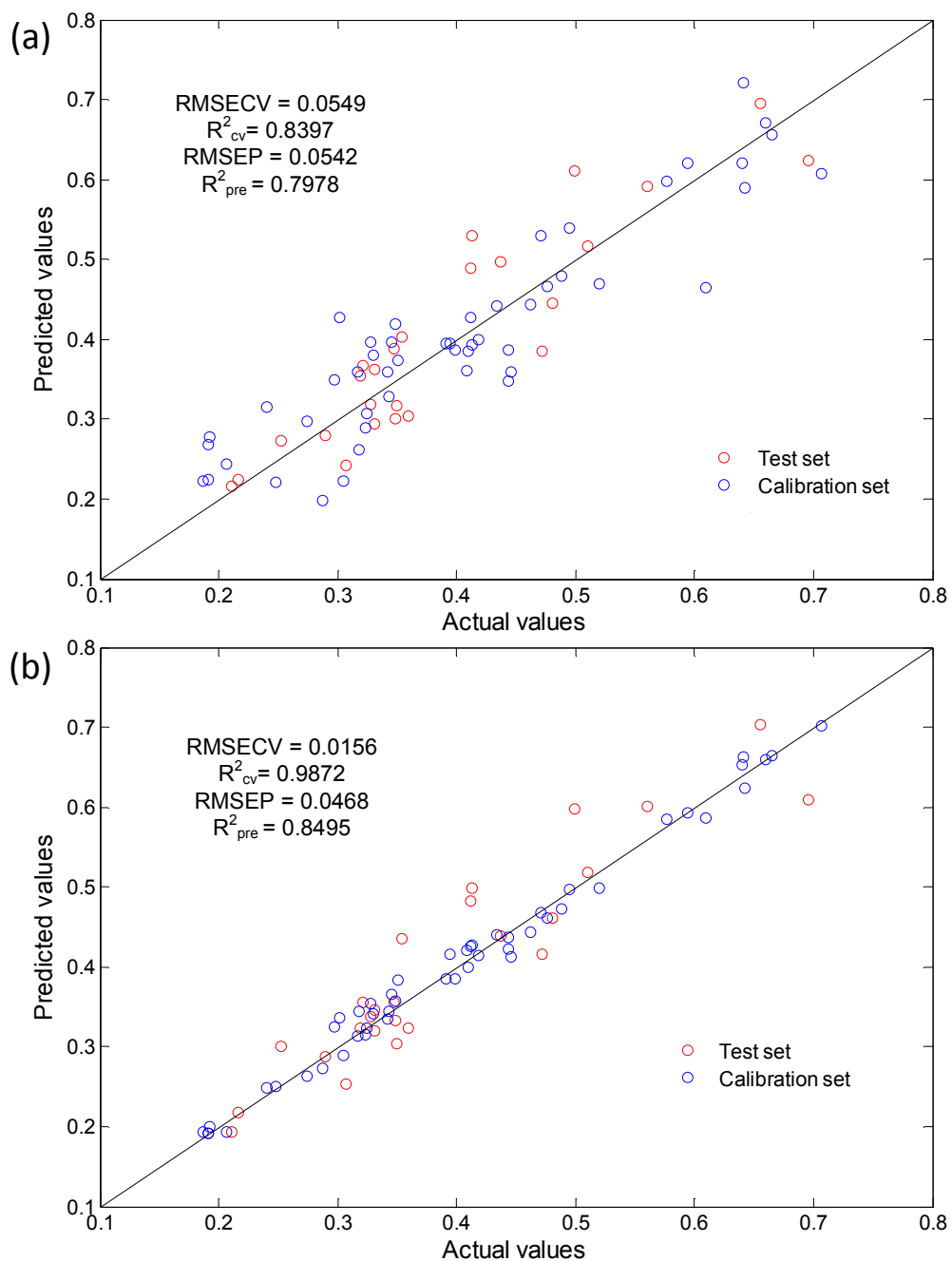
493

Fig. 3. The distribution of the selected variables obtained by different wavelength

494

selection methods.

495



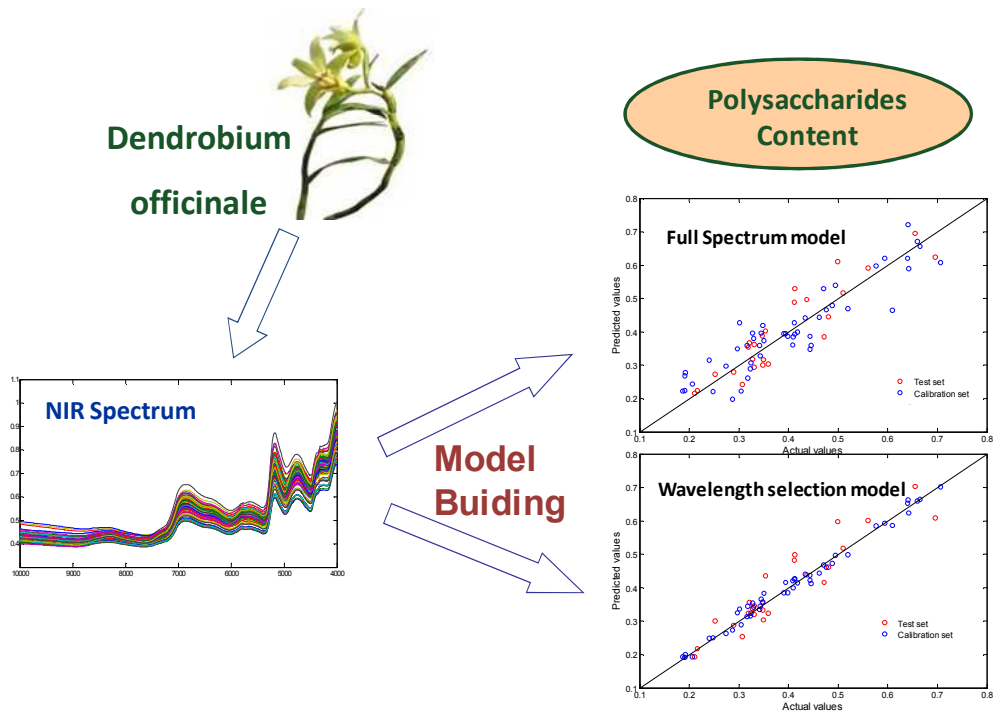
496

497

498 Fig. 4. The correlation between predicted value and measured value of polysaccharides

499 content based on (a) The full spectra PLS model; (b) 39 selected wavelengths by CARS.

500



Graphic Abstract

NIR spectroscopy method for the quantification of polysaccharides in

Dendrobium officinale by PLS calibration model.