

# Analytical Methods

Accepted Manuscript



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. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it 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.

1  
2  
3 **Rapid determination of total polyphenols content and antioxidant activity in**  
4 ***Dendrobium officinale* by near-infrared spectroscopy**

5  
6 Longhui Ma<sup>a</sup>, Zhimin Zhang<sup>a</sup>, Xingbing Zhao<sup>b</sup>, Sufeng Zhang<sup>b</sup>, Hongmei Lu<sup>a\*</sup>

7  
8  
9 <sup>a</sup>College of Chemistry and Chemical Engineering, Central South University, Changsha  
10 410083, PR China

11  
12 <sup>b</sup>Hunan Longshishan Dendrobium Candidum Wall.ex Lindl Base Co., Ltd, Changsha 410205,  
13 PR China  
14  
15

16  
17  
18  
19 **Abstract:** Despite their popularity and extensive use, some herbs have not been  
20 officially recognized in most countries. The main reason is the lack of comprehensive  
21 research data and methods. In this paper, a rapid approach based on near-infrared  
22 spectroscopy (NIR) was developed for the determination of total polyphenols content  
23 (TPC) and antioxidant activity (AA) in *Dendrobium officinale* (*D. officinale*), an  
24 important Chinese herb. Adopting the Folin-Ciocalteu (FC) assay and  
25 2,2-diphenyl-1-picrylhydrazyl radical (DPPH) free radical scavenging activity as the  
26 reference methods, TPC and AA in *D. officinale* samples (n=83) collected from  
27 different locations in China were analyzed. Spectra generated by NIR were pretreated  
28 with different preprocessing methods and analyzed with partial least-square (PLS). To  
29 obtain robust and predictive quantitative model, competitive adaptive reweighted  
30 sampling (CARS) was applied to screen the key variables. The correlation coefficient  
31 of prediction ( $R_{pre}^2$ ) and root mean square error of prediction (RMSEP) by  
32 competitive adaptive reweighted sampling - partial least-square (CARS-PLS) were  
33 0.8412, 0.2905 for TPC and 0.9062, 0.1028 for AA, respectively. The results show  
34 that the combination of NIR spectroscopy with CARS-PLS provides a rapid and  
35 precise alternative to existing chemical analysis for the determination of TPC and AA  
36 in *D. officinale*.  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52

53 **Keywords:** Polyphenols; Antioxidant activity; Near-infrared spectroscopy;  
54 *Dendrobium officinale*  
55  
56

57  
58  
59  
60  

---

\* Corresponding author. E-mail address: hongmeilu@csu.edu.cn

## 1. Introduction

*Dendrobium officinale* (*D. officinale*) is ranked “the first of the nine Chinese fairy herbs” and has been widely used as pharmaceuticals or functional foods for thousands of years in China.<sup>1,2</sup> Despite its popularity and extensive use during the last decade, it has not been officially recognized in most countries. It is far from sufficient to meet the criteria needed to support its use world-wide. The main reason is the lack of comprehensive research data and methods about component, activity and safety etc. Previous studies on *D. officinale* mainly focused on polysaccharides, due to the content of polysaccharides used as one of the quality assessment criteria (no less than 0.2500 g of glucose per g dry weight) in Chinese pharmacopoeia.<sup>3</sup> Barely researches focus on determining the polyphenol, another kind of the main components in *D. officinale*, and its activity.<sup>4</sup> It is well-known that polyphenol compounds are associated with reducing risk of developing chronic diseases, such as aging, chronic gastric, various cancer and cardiovascular.<sup>5, 6</sup> Moreover, many studies have demonstrated that polyphenol compounds also contribute to the antioxidant activity.<sup>7,8</sup> It is important and necessary to qualitatively and quantitatively study polyphenol compounds and antioxidant activity in *D. officinale*. However, existing analytical methods such as colorimetric measurement<sup>9, 10</sup> and HPLC measurement<sup>11</sup> involve tedious sample preparations, intensive labor and expensive solvents,<sup>12</sup> which make them unsuitable for routine analysis.

Recently, our group have successfully applied near-infrared spectroscopy (NIR) for rapid authentication or identification of some herb and food<sup>13, 14</sup> and precise

1  
2  
3  
4 51 prediction of component content.<sup>15, 16</sup> NIR has several desirable advantages over  
5  
6 52 existing analytical methods: e.g. non-invasion, cost-efficiency, less laboring and high  
7  
8  
9 53 efficiency.<sup>17</sup> As the absorption of NIR spectra correspond to molecular overtone and  
10  
11 54 combination vibrations, the absorption peaks are broad and strongly overlapped,  
12  
13  
14 55 which is hard to interpret. In such case, partial least-squares (PLS) regression has  
15  
16 56 been frequently employed to make a calibration model with spectral data .<sup>18</sup> Since  
17  
18  
19 57 some spectral regions contain noise from environment and interference variables,  
20  
21 58 better calibration model can be obtained by selecting effective variables instead of the  
22  
23  
24 59 full-spectrum.<sup>19, 20</sup> With the rapid calculation process, competitive adaptive  
25  
26 60 reweighted sampling (CARS) variable selection method has been proposed by Li et al  
27  
28  
29 61 lately.<sup>21</sup> And lots of literatures have successfully proved that CARS is a powerful and  
30  
31 62 high-performance tool in complex analytical system.<sup>22, 23</sup>

33  
34 In this study, NIR spectroscopy together with CARS variables selection and PLS  
35  
36 64 regression was used to simultaneous determination of total polyphenols content (TPC)  
37  
38  
39 65 and antioxidant activity (AA) in *D. officinale*. To the best of our knowledge, this  
40  
41 66 approach has never been tried before to determinate the *D. officinale* TPC and AA.  
42  
43  
44 67 The specific procedure was outlined as follows: (1) adopting existing chemical  
45  
46 68 analysis methods as the reference methods to analyze TPC and AA in *D. officinale*  
47  
48  
49 69 samples; (2) comparing different spectral data preprocessing methods and screening  
50  
51 70 the key variables using CARS algorithm; (3) establishing the quantitative models for  
52  
53  
54 71 determination of TPC and AA based on NIR spectra.  
55

56  
57 72  
58  
59  
60

## 73 2. Materials and methods

### 74 2.1. Materials and Reagents

75 A total of 83 *D. officinale* samples were collected from different locations in  
76 China during the period from April 2012-April 2014. It provided a representative set  
77 of *D. officinale* consumed in domestic market, which comprised enough variation to  
78 make the quantitative model robust. Folin-Ciocalteu reagent was purchased from  
79 Sigma Aldrich (St. Louis, Mo, USA). 2,2-diphenyl-1-picrylhydrazyl radical (DPPH),  
80 gallic acid (99% purity) and 6-hydroxy-2,5,7,8-tetramethyl-2-carboxylic acid (Trolox)  
81 were obtained from the National Institution for Food and Drug Control (Beijing,  
82 China). Dehydrated alcohol and anhydrous sodium carbonate ( $\text{Na}_2\text{CO}_3$ ) were  
83 purchased from Sinopharm Chemical Reagent Co.Ltd (Shanghai, China). Deionized  
84 water was purified with a Milli-Q system (Millipore, Bedford, MA, USA).

### 85 2.2. Reference analysis

86 The dried samples were ground and passed through a 60 mesh sieve. An accurate  
87 weight (1.0 g) of each powder sample was mixed with 30 mL 80% anhydrous ethanol.  
88 The mixture was then sonicated for 30 min at 35 °C. After the ethanol extract was  
89 filtered, the filtrate was moved into 50 mL volumetric flask. Prepared extracts were  
90 stored at 4 °C. TPC was determined by an improved Folin-Ciocalteu colorimeter<sup>24</sup>  
91 and expressed as mg of gallic acid equivalent per 1g of dry weight sample. Reference  
92 measurement of AA in *D. officinale* was assayed by the improved DPPH radical  
93 scavenging activity.<sup>25</sup> The UV-Vis spectra of *D. officinale* were acquired on a BTT  
94 miniature array spectrophotometer (B&W Tek, Newark, DE, USA) equipped with

1  
2  
3  
4 95 glass or quartz cells of 1 cm path length in the range 200-800 nm.  
5

### 6 96 **2.3. NIR spectra collection**

7  
8  
9 97 NIR spectra of sample powders were acquired by a diffuse reflectance mode  
10  
11 98 using the Antaris II Fourier transform-NIR System (Thermo Scientific Inc, Madison,  
12  
13 99 USA). The number of scans was 32 and the spectral resolution was  $8\text{ cm}^{-1}$ . The range  
14  
15  
16 100 of spectra was from 10000 to  $4000\text{ cm}^{-1}$  and the data were measured in  $3.9\text{ cm}^{-1}$   
17  
18  
19 101 interval, which resulted in 1557 variables. Each sample was scanned three times in a  
20  
21 102 ring cup and the average spectra were collected for subsequent analysis.  
22

### 23 103 **2.4. Chemometric methods**

#### 24 104 **2.4.1. Spectral preprocessing and outlier detection**

25  
26  
27  
28  
29 105 Since physical variations, such as particle size and shape, sample packing and  
30  
31 106 sample surface, could impact on spectra measurement, the raw spectra inevitably  
32  
33 107 consist of systematic noises or background information.<sup>26</sup> Accordingly, a proper  
34  
35  
36 108 spectral preprocessing method is necessary to reduce the unwanted spectral variations.  
37  
38  
39 109 In this study, different kinds of spectral preprocessing methods were compared,  
40  
41 110 including smoothing, standard normal variable (SNV), multivariate scatter correction  
42  
43 111 (MSC), Savitzky-Golay first-derivative (SG1), and the combinations of SNV (or MSC)  
44  
45  
46 112 with the derivative. Smoothing is an averaging algorithm that used to reduce the noise  
47  
48  
49 113 and to enhance the signal-noise ratio (SNR). MSC or SNV method is always  
50  
51 114 performed to remove slope variation and to modify scatter effects. By calculating SG1  
52  
53  
54 115 derivative, the baseline drift are eliminated and small spectral differences are  
55  
56  
57 116 enhanced.  
58  
59  
60

1  
2  
3  
4 117 In general, outliers are incorrect or abnormal ones in some sense compared to the  
5  
6 118 majority of the data. Outliers in the calibration set would lead to severe errors on the  
7  
8  
9 119 model, while outliers in the prediction set would obtain misleading results to evaluate  
10  
11 120 the model performance. Considering the calibration model is so sensitive to outliers,  
12  
13 121 Monte-Carlo (MC) method<sup>27</sup> was applied to eliminate outliers in spectral dataset  
14  
15  
16 122 before quantitative analysis. The core idea of the MC outlier detector is to develop a  
17  
18  
19 123 Monte-Carlo procedure for detecting outlier by studying the distribution of prediction  
20  
21 124 errors of each sample obtained from original data set. The detailed description of  
22  
23  
24 125 scheme and procedure for this method based on predictive errors and Monte-Carlo  
25  
26 126 sampling has been shown in the Ref.<sup>27</sup>

#### 27 28 29 127 **2.4.2. PLS model and evaluation**

30  
31 128 The quantitative models of TPC and AA in *D.officinale* were developed by PLS  
32  
33  
34 129 regression. As a well-known method, it is used to establish relationships between  
35  
36  
37 130 spectra data matrix (**X**) and reference concentration of elements matrix (**y**) with a  
38  
39 131 small number of latent variables (LV).<sup>18, 28</sup> PLS follows a two-step strategy to  
40  
41 132 establish a functional relationship between **X** and **y**. Firstly, the input variable matrix  
42  
43  
44 133 **X** (n, m) and output variable matrix **y** (n, 1) are decomposed as followed:

$$45  
46 134 \mathbf{X} = \mathbf{TP}^T + \mathbf{E} \quad (1)$$

$$47  
48 135 \mathbf{y} = \mathbf{Tq}^T + \mathbf{f} \quad (2)$$

49  
50  
51 136 where **P** (m, k) and **q** (k, 1) are the loadings, **T** (n, k) is scores matrix, **E** (n, m) and **f**  
52  
53  
54 137 (n, 1) are error terms which are not explained by the model and k is the number of LV  
55  
56  
57 138 used in the PLS model.

1  
2  
3  
4 139 Subsequently, a least squares regression is performed on the extracted orthogonal  
5  
6 140 latent variables/score vectors. Input and output scores vectors are related by a multiple  
7  
8  
9 141 lineal regression model:

10  
11 142 
$$\hat{y} = \mathbf{Xb} + \mathbf{f} \quad (3)$$

12  
13  
14 143 where  $\mathbf{b}$  ( $m,1$ ) is an inner coefficient.

15  
16 144 All of the samples were divided into two sets by Duplex algorithm: one for  
17  
18  
19 145 calibrating model and the other for prediction ability. Duplex algorithm is a method  
20  
21 146 for the selection of a representative set of samples, in which calibration and prediction  
22  
23  
24 147 objects are selected alternately, starting with the inclusion of the most distant pair of  
25  
26 148 objects into the calibration set.<sup>29</sup>

27  
28  
29 149 The root mean square error of calibration (RMSEC), correlation coefficient of  
30  
31 150 calibration ( $R_c^2$ ), the root mean square error of prediction (RMSEP) and correlation  
32  
33  
34 151 coefficient of prediction ( $R_{pre}^2$ ) were calculated to evaluate the performance of the  
35  
36 152 final quantitative models. Generally, a robust and accurate model should have low  
37  
38  
39 153 values of RMSEC and RMSEP and high values of  $R_c^2$  and  $R_{pre}^2$ . Besides, these  
40  
41 154 parameters are defined as follows:

42  
43  
44 155 
$$\text{RMSEC or RMSEP} = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}} \quad (4)$$

45  
46  
47  
48 156 
$$R_c^2 \text{ or } R_{pre}^2 = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}} \quad (5)$$

49  
50  
51  
52  
53 157 where  $y_i$  is the measured value and  $\hat{y}_i$  is the prediction value;  $\bar{y}$  is the average  
54  
55 158 measurement.

56  
57  
58 159 **2.4.3. Key variables selection**



1  
2  
3  
4 160 CARS algorithm was applied to locate an optimal combination of the variables  
5  
6 161 for accurately determination of TPC and AA. The major idea of CARS is to employ  
7  
8 162 the mechanic of “survival of the fittest” based on Darwin’s Evolution.<sup>21</sup> Absolute  
9  
10 163 values of regression coefficients of PLS model are used as an index for evaluating the  
11  
12 164 importance of each variable. Then, according to the importance level of each variable,  
13  
14 165 CARS sequentially selected N subsets of variables from N Monte Carlo (MC)  
15  
16 166 sampling run in an iterative and competitive manner. Next, the exponentially  
17  
18 167 decreasing function (EDF) and adaptive reweighted sampling (ARS) are employed to  
19  
20 168 eliminate the variables, which are of relatively small absolute regression coefficients  
21  
22 169 by force. Finally, the subset with the lowest root mean square error of cross validation  
23  
24 170 (RMSECV) is considered as optimal combination of the variables.  
25  
26  
27  
28  
29  
30

## 31 171 **2.5. Software**

32  
33  
34 172 All algorithms were implemented with MATLAB for Windows (Version 2013A,  
35  
36 173 the MathWorks, Inc). The code of CARS was available as open source software in the  
37  
38 174 website: <http://www.libpls.net/>.  
39  
40  
41  
42

## 43 176 **3. Results and discussion**

### 44 177 **3.1. TPC and AA measure results using the reference methods**

45  
46  
47 178 The TPC and AA of 83 samples were determined using the reference methods  
48  
49 179 (see Section 2.2). The content ranges of TPC (mg/g) and AA (mg/g) are shown in  
50  
51 180 Table 1. Obviously, different geographical origins and harvest seasons lead to  
52  
53 181 considerable variability of TPC and AA in *D. officinale*. The content ranges of  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4 182 reference values in calibration set are 0.3148-1.5922 (mg/g) for TPC and  
5  
6 183 1.3606-3.8962 mg/g for AA. The content ranges of TPC and AA in prediction set are  
7  
8  
9 184 0.2332-1.4835 mg/g and 1.1995-4.3435 mg/g, respectively. There are no significant  
10  
11 185 dissimilarities to be observed in means, ranges and standard deviations (SD) between  
12  
13 186 calibration set and prediction set. This indicates that the Duplex algorithm enable the  
14  
15  
16 187 same diversity in both sets.  
17

18  
19 **Insert Table 1**

20  
21 **Insert Fig. 1**

### 22 23 24 190 **3.2. Selection of preprocessing methods**

25  
26 191 **Fig.1** shows the raw NIR spectra of 83 *D. officinale* samples at wavenumbers  
27  
28 192 10000-4000  $\text{cm}^{-1}$ . NIR spectra in this region contain some intensive spectral peaks.  
29  
30  
31 193 These intensive peaks correspond to the vibration of some groups such as the  
32  
33 194 combination of C-H and C-C ( $4000 \text{ cm}^{-1}$ ), the second overtone vibration of the  
34  
35 195 carbonyl group ( $5350 \text{ cm}^{-1}$ ), the first overtone of O-H and N-H ( $6900 \text{ cm}^{-1}$ ), stretching  
36  
37 196 and deformation vibrations of C-H ( $7200 \text{ cm}^{-1}$ ). It is apparently difficult to  
38  
39 197 discriminate samples just by visually examining the raw average spectra. Hence,  
40  
41 198 preprocessing methods are critical to enhance the quality of spectra. Five spectral  
42  
43 199 preprocessing methods were applied and compared on the basis of RMSECV. The  
44  
45 200 results are showed in Table 2. According to the model evaluation standard, different  
46  
47 201 preprocessing methods are chosen for TPC and AA. Smoothing + MSC is the best  
48  
49 202 choice for the developing of TPC model, while SNV + SG1 derivative is the best one  
50  
51 203 for the developing of AA quantitative model. The reason is that weights of  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4 204 wavenumber are different when building PLS models for TPC and AA. SG1 derivative  
5  
6 205 algorithm for TPC calibration model may increase noise, whereas this pretreatment  
7  
8  
9 206 method may reflect more information corresponding to AA.

10  
11 **Insert Table 2**

12  
13 **3.3. Outlier detection for robust models**

14  
15 209 The method based on Monte-Carlo was performed to detect outliers.<sup>27</sup> Three  
16  
17 210 types of outliers should be considered, namely outliers in X, outliers in y, as well as  
18  
19 211 outliers towards the model. Fig. 2a and Fig. 2b show the diagnostic plots of the outlier  
20  
21 212 detection for TPC and AA respectively. As shown in Fig. 2, a threshold of standard  
22  
23 213 deviation was set on the top left area. The samples were outliers in X if their standard  
24  
25 214 deviation is larger than the threshold. In the example, object 28 in Fig. 2a and four  
26  
27 215 objects (10, 11, 27, 28) in Fig. 2b are X outliers. The y outliers appear on the lower  
28  
29 216 right area, namely objects (39, 42) in Fig. 2a and objects (9, 12) in Fig. 2b, which  
30  
31 217 could cause a large error sum of squares. In addition, it can be seen that object 29 on  
32  
33 218 top right area of the two figures is outlier towards the model. Clearly, both objects 28  
34  
35 219 and 29 decrease accuracy of the model and affect the subsequent ability to determine  
36  
37 220 the TPC and the AA. These two outliers may be caused by measurement errors when  
38  
39 221 collected the NIR spectra. Hence, 4 samples (28, 29, 39, 42) and 6 samples (10, 11, 12,  
40  
41 222 27, 28, 29) are removed for TPC and AA, respectively. After samples outlier detection,  
42  
43 223 there are 79 and 76 samples for TPC and AA analysis, respectively. Then the  
44  
45 224 remaining samples are divided using the Duplex algorithm.

46  
47  
48  
49  
50  
51  
52  
53  
54  
55 **Insert Fig. 2**

56  
57  
58 **3.4. Variable selection**

1  
2  
3  
4 227 As stated above, quantitative models are first established based on full spectrum  
5  
6 228 (10000-4000  $\text{cm}^{-1}$ ). The results are satisfying (Table 3). RMSEP are 0.3242, 0.1232  
7  
8  
9 229 and  $R_{\text{pre}}^2$  are 0.8023, 0.8653 for TPC and AA, respectively. However, the input  
10  
11 230 variables are too much, which affect the robust of the quantitative models and  
12  
13  
14 231 increase the calculation time. Thus, CARS has been applied to screen key variables  
15  
16 232 prior to application of PLS and improve the model performance.

17  
18  
19 233 In this study, the number of MC sampling runs was set to 100 as default during  
20  
21 234 the calculation process. In order to guarantee the reliability of the model, the CARS  
22  
23  
24 235 procedure was conducted 100 times and RMSECV values were recorded. The optimal  
25  
26 236 variable subsets were selected for further analysis according to the minimum  
27  
28  
29 237 RMSECV. Finally, the obtained numbers of key variables are decreased from 1557 to  
30  
31 238 26 for TPC analysis and from 1557 to 34 for AA analysis, respectively. Under the  
32  
33  
34 239 selected key variables and spectral preprocessing method, PLS models for TPC and  
35  
36 240 AA are established. The results of both full spectrum-PLS models and competitive  
37  
38  
39 241 adaptive reweighted sampling - partial least-square (CARS-PLS) are listed in Table 3.  
40  
41 242 The  $R_c^2$  values of two methods for TPC are very close. However, CARS-PLS has  
42  
43  
44 243 smaller RMSEP (0.2905) and higher  $R_{\text{pre}}^2$  (0.8412), indicating that CARS-PLS models  
45  
46 244 for TPC has better predictive ability than full spectrum-PLS. In terms of CARS-PLS  
47  
48  
49 245 model for determining the AA, the  $R_c^2$  and  $R_{\text{pre}}^2$  are significantly increased to 0.9854  
50  
51 246 and 0.9062; RMSEC and RMSEP are reduced to 0.0294 and 0.1028. In short, CRAS  
52  
53  
54 247 could eliminate uninformative variables effectively and improve the predictive  
55  
56 248 precision of the model to a certain extent.  
57  
58  
59  
60

1  
2  
3  
4 249 Fig. 3 shows the relationship between NIR predicted values and measure values  
5  
6 250 for TPC and AA obtained by PLS calibration models combining CARS method. The  
7  
8  
9 251 diamond marked points referred to calibration set, and round marked points referred  
10  
11 252 to prediction set. With a close observation, TPC model is a little worse compared with  
12  
13  
14 253 the AA model, because some points in Fig. 3a fall off the bisectrix line. The reason  
15  
16 254 may be that Folin–Ciocalteu assay for TPC is a less accurate or precise method than  
17  
18  
19 255 DPPH radical scavenging assay for AA. Therefore, it is difficult to achieve a more  
20  
21 256 robust model for determination of TPC. In general, the results suggest that PLS  
22  
23  
24 257 models based on CARS for determination TPC and AA are accurate compared with  
25  
26 258 full spectrum-PLS models.

28  
29 **Insert Fig.3**

30  
31 **Insert Table 3**

### 32 33 34 **3.5. Interpretation of key variables**

35  
36 262 The positions of selected key variables are illustrated by marked points in Fig. 4.  
37  
38  
39 263 The polyphenol compounds contain abundant hydrogenous bonds (i.e. C-H and O-H  
40  
41 264 groups, etc) and exhibit antioxidant activity. The information contained in selected  
42  
43  
44 265 spectral regions plays a crucial role in determination of TPC and AA. The selected  
45  
46 266 variables for both TPC and AA prediction are mainly concentrated in the region of  
47  
48  
49 267  $7200\text{-}4000\text{ cm}^{-1}$  and  $9990\text{-}8333\text{ cm}^{-1}$ . The former is assigned to the combination bands  
50  
51 268 of the functional groups  $\text{-C=O}$ ,  $\text{N-H}$ ,  $\text{C-H}$  and  $\text{C-C}$ . The latter is the region of C-H  
52  
53  
54 269 third overtone and combination tone.<sup>17, 23</sup> For example, the selected region of around  
55  
56 270  $4204\text{ cm}^{-1}$  is related to the combination tone of C-H and C-C stretching vibration, and  
57  
58  
59  
60

1  
2  
3  
4 271 the region of 6940-7140  $\text{cm}^{-1}$  is the first overtone of O-H.<sup>30</sup> One can see that the  
5  
6 272 selected variables are associated with needful chemistry feature of the TPC and AA in  
7  
8  
9 273 *D. officinale*. These indicate that CARS method has the ability to screen the key and  
10  
11 274 effective variables.

12  
13  
14 275 **Insert Fig.4**

#### 15 16 276 **4. Conclusion**

17  
18  
19 277 NIR spectroscopy coupled with chemometric methods was successfully utilized  
20  
21 278 for rapid quantification of the total polyphenols content (TPC) and antioxidant  
22  
23 279 activity (AA) in *D. officinale*. The most suitable data preprocessing methods were  
24  
25 280 smoothing + MSC and SNV + SG1 derivative for TPC and AA model, respectively.  
26  
27  
28 281 Combined with the corresponding preprocessing method, CARS was used for  
29  
30  
31 282 screening informative variables and reducing uninformative variables. Twenty-six  
32  
33 283 variables were picked out of 1557 wavenumbers by CARS for the prediction TPC,  
34  
35 284 and 34 variables for the prediction AA. CARS-PLS models achieved the optimal  
36  
37 285 performance with  $R_{\text{pre}}^2$  and RMSEP were 0.8412, 0.2905 for TPC and 0.9062, 0.1028  
38  
39 286 for AA, respectively. The results show that NIR spectroscopy combined CARS with  
40  
41  
42 287 PLS algorithm is able to determine TPC and AA in *D. officinale* in a fast, accurate and  
43  
44 288 reliable way. Based on this study and previous studies, NIR spectroscopy coupled  
45  
46 289 with chemometric methods is a rapid potent approach for simultaneous quantification  
47  
48  
49 290 of chemical constituent content and activities in herbs and foods. These results also  
50  
51  
52 291 provide some fundamental research data and method for *D. officinale*. In the future,  
53  
54  
55 292 more samples from different geographical areas will be studied to obtain a model with  
56  
57  
58  
59  
60

1  
2  
3  
4 293 the greatest applicability.  
5  
6  
7 294

8  
9 295 **Acknowledgments**

10  
11 296 The authors gratefully thank the National Natural Science Foundation of China for  
12  
13 297 support of the projects (No. 21175157, 21375151 and 21305163).  
14  
15  
16 298

17  
18 299 **References**

- 19 300 1. X. Lin, P.-C. Shaw, S. C.-W. Sze, Y. Tong and Y. Zhang, *International immunopharmacology*,  
20 301 2011, **11**, 2025-2032.  
21 302 2. G. Ding, G. Xu, W. Zhang, S. Lu, X. Li, S. Gu and X.-Y. Ding, *European Food Research and*  
22 303 *Technology*, 2008, **227**, 1283-1286.  
23 304 3. Chinese Pharmacopoeia Committee, *Chinese Pharmacopoeia*, People's Medical Publishing  
24 305 House, Beijing, 2010.  
25 306 4. Y. Li, C.-L. Wang, Y.-J. Wang, S.-X. Guo, J.-S. Yang, X.-M. Chen and P.-G. Xiao, *Chemical and*  
26 307 *Pharmaceutical Bulletin*, 2009, **57**, 218-219.  
27 308 5. X. Zhang, J.-K. Xu, J. Wang, N.-L. Wang, H. Kurihara, S. Kitanaka and X.-S. Yao, *Journal of*  
28 309 *Natural Products*, 2007, **70**, 24-28.  
29 310 6. Y. Cai, Q. Luo, M. Sun and H. Corke, *Life sciences*, 2004, **74**, 2157-2184.  
30 311 7. M. P. Kähkönen, A. I. Hopia, H. J. Vuorela, J.-P. Rauha, K. Pihlaja, T. S. Kujala and M. Heinonen,  
31 312 *Journal of agricultural and food chemistry*, 1999, **47**, 3954-3962.  
32 313 8. J. Švarc-Gajić, Z. Stojanović, A. S. Carretero, D. A. Román, I. Borrás and I. Vasiljević, *Journal of*  
33 314 *Food Engineering*, 2013, **119**, 525-532.  
34 315 9. L. M. Magalhães, M. A. Segundo, S. Reis and J. L. Lima, *Analytica chimica acta*, 2008, **613**,  
35 316 1-19.  
36 317 10. S. Georgé, P. Brat, P. Alter and M. J. Amiot, *Journal of Agricultural and Food Chemistry*, 2005,  
37 318 **53**, 1370-1373.  
38 319 11. M.-S. Lee, Y.-S. Hwang, J. Lee and M.-G. Choung, *Food chemistry*, 2014, **158**, 351-357.  
39 320 12. H. Yan, B.-x. Han, Q.-y. Wu, M.-z. Jiang and Z.-z. Gui, *Spectrochimica Acta Part A: Molecular*  
40 321 *and Biomolecular Spectroscopy*, 2011, **79**, 179-184.  
41 322 13. Y. Wei, W. Fan, X. Zhao, W. Wu and H. Lu, *Analytical Letters*, 2014, **48**, 817-829.  
42 323 14. X. Wang, J. Huang, W. Fan and H. Lu, *Analytical Methods*, 2015, **7**, 787-792.  
43 324 15. Y. Yun, Y. Wei, X. Zhao, W. Wu, Y. Liang and H. Lu, *RSC Advances*, 2015.  
44 325 16. Q. Luo, Y. Yun, W. Fan, J. Huang, L. Zhang, B. Deng and H. Lu, *RSC Advances*, 2015, **5**,  
45 326 5046-5052.  
46 327 17. X. B. Zou, J. W. Zhao, M. J. Povey, M. Holmes and H. P. Mao, *Analytica chimica acta*, 2010, **667**,  
47 328 14-32.  
48 329 18. S. Wold, M. Sjöström and L. Eriksson, *Chemometrics and intelligent laboratory systems*, 2001,  
49 330 **58**, 109-130.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 331 19. H.-D. Li, Y.-Z. Liang, X.-X. Long, Y.-H. Yun and Q.-S. Xu, *Chemometrics and Intelligent Laboratory*  
4 332 *Systems*, 2013, **122**, 23-30.  
5 333 20. B.-C. Deng, Y.-H. Yun, P. Ma, C.-C. Lin, D.-B. Ren and Y.-Z. Liang, *The Analyst*, 2015, **140**,  
6 334 1876-1885.  
7 335 21. H. Li, Y. Liang, Q. Xu and D. Cao, *Analytica chimica acta*, 2009, **648**, 77-84.  
8 336 22. D. Wu and D.-W. Sun, *Talanta*, 2013, **116**, 266-276.  
9 337 23. X. Wei, N. Xu, D. Wu and Y. He, *Food and Bioprocess Technology*, 2013, **7**, 184-190.  
10 338 24. K. Tawaha, F. Q. Alali, M. Gharaibeh, M. Mohammad and T. El-Elimat, *Food chemistry*, 2007,  
11 339 **104**, 1372-1378.  
12 340 25. R. Scherer and H. T. Godoy, *Food chemistry*, 2009, **112**, 654-658.  
13 341 26. J. W. Jin, Z. P. Chen, L. M. Li, R. Steponavicius, S. N. Thennadil, J. Yang and R. Q. Yu, *Analytical*  
14 342 *chemistry*, 2012, **84**, 320-326.  
15 343 27. D. S. Cao, Y. Z. Liang, Q. S. Xu, H. D. Li and X. Chen, *Journal of computational chemistry*, 2010,  
16 344 **31**, 592-602.  
17 345 28. P. Geladi and B. R. Kowalski, *Analytica chimica acta*, 1986, **185**, 1-17.  
18 346 29. R. D. Snee, *Technometrics*, 1977, **19**, 415-428.  
19 347 30. J. Workman Jr and L. Weyer, *Practical guide to interpretive near-infrared spectroscopy*, CRC  
20 348 press, Boca Raton, 2007.  
21 349  
22 350  
23 351  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



352 **Tables:**

353 Table 1. TPC and AA measured with the existing chemical methods and the number of  
354 *D. officinale* samples used in dataset.

Constituent	Calibration set				Prediction set			
	S.N. <sup>a</sup>	Mean	S.D. <sup>b</sup>	Range	S.N. <sup>a</sup>	Mean	S.D. <sup>b</sup>	Range
TPC(mg/g)	60	2.2802	0.5634	1.3606-3.8962	19	2.4886	0.7177	1.1995-4.3435
AA(mg/g)	59	0.9535	0.2452	0.3148-1.5922	17	0.8594	0.3456	0.2332-1.4835

355 <sup>a</sup>S.N: sample number; <sup>b</sup>S.D: standard deviation.

356

357 Table 2. Comparison of different spectral preprocessing methods on performance of  
 358 PLS calibration models.

Properties	Methods	PLS results			
		RMSECV	RMSEP	RMSEC	$R_c^2$
TPC	raw data	0.3916	0.4228	0.2440	0.8358
	Smoothing+ SNV	0.3705	0.4158	0.1906	0.8659
	Smoothing+ MSC	0.3287	0.3922	0.1837	0.8766
	SNV+SG1st	0.3639	0.4052	0.2148	0.8509
	MSC+SG1st	0.3617	0.4049	0.1977	0.8522
AA	raw data	0.2558	0.3109	0.1083	0.8641
	Smoothing+ SNV	0.2343	0.2386	0.1009	0.8761
	Smoothing+ MSC	0.2464	0.2841	0.1072	0.8626
	SNV+SG1st	0.2280	0.2291	0.0846	0.9182
	MSC+SG1st	0.2350	0.2394	0.0853	0.9094

359

1  
2  
3  
4 360 Table 3. PLS model results based on full-spectrum and selected wavelengths using  
5  
6 361 CARS.  
7

properties	methods	Variable number	Calibration set		Prediction set	
			$R_c^2$	RMSEC	$R_{pre}^2$	RMSEP
TPC(mg/g)	PLS	1557	0.9121	0.1656	0.8023	0.3242
	CRAS-PLS	26	0.9185	0.1596	0.8412	0.2905
AA(mg/g)	PLS	1557	0.9488	0.0550	0.8653	0.1232
	CRAS-PLS	34	0.9854	0.0294	0.9062	0.1028

8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23 362  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

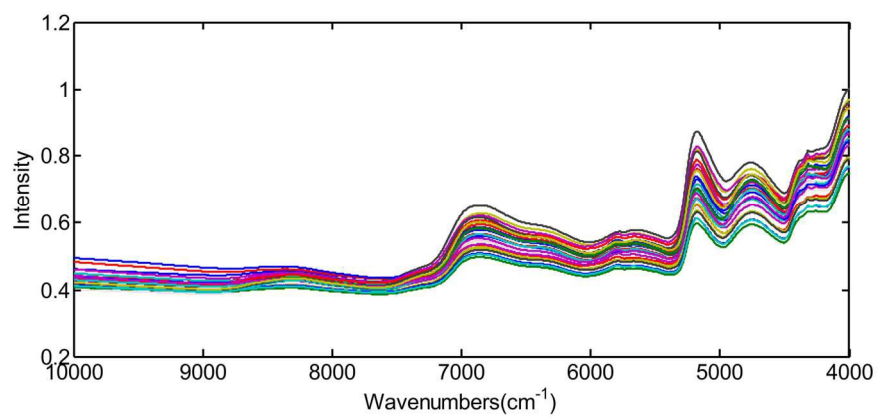
1  
2  
3 363 **Figure captions**

4  
5 364 Fig. 1. NIR spectra of the *D.officinale* samples (n=83).

6  
7 365 Fig. 2. Diagnostic plots for outlier detection based on MCS. (a) TPC; (b) AA.

8  
9  
10 366 Fig. 3. Correlation diagrams between the predicted values and measured values based  
11  
12 367 on CARS-PLS method. (a) TPC; (b) AA.

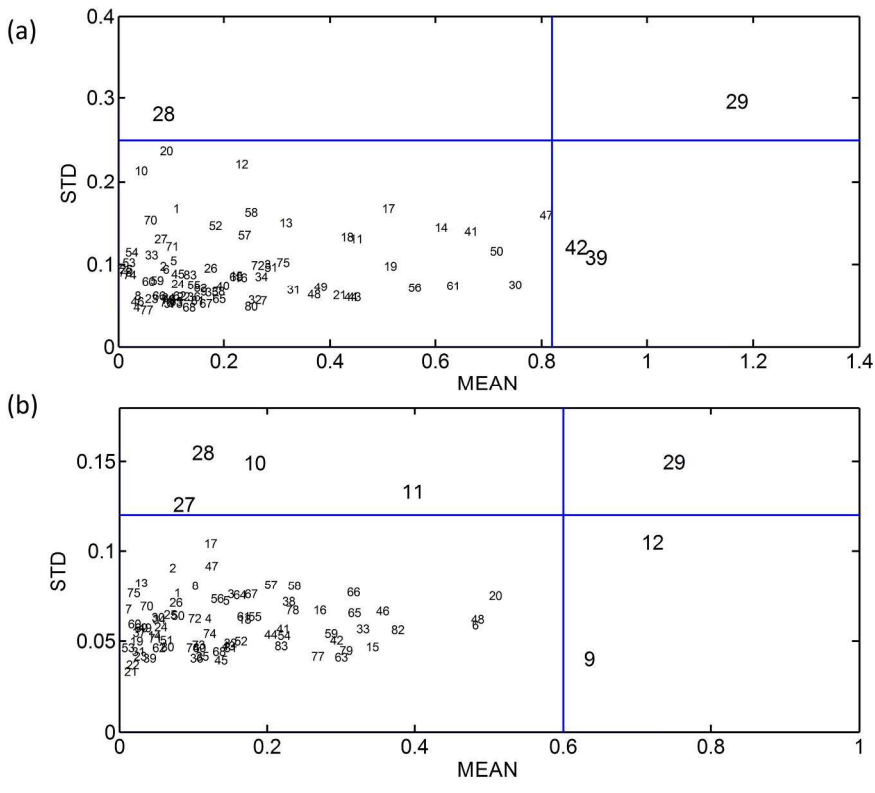
13  
14  
15 368 Fig. 4. Positions of variables selected by CARS-PLS for prediction of TPC (round  
16  
17 369 marked points) and AA (diamond marked points) on the full spectra.



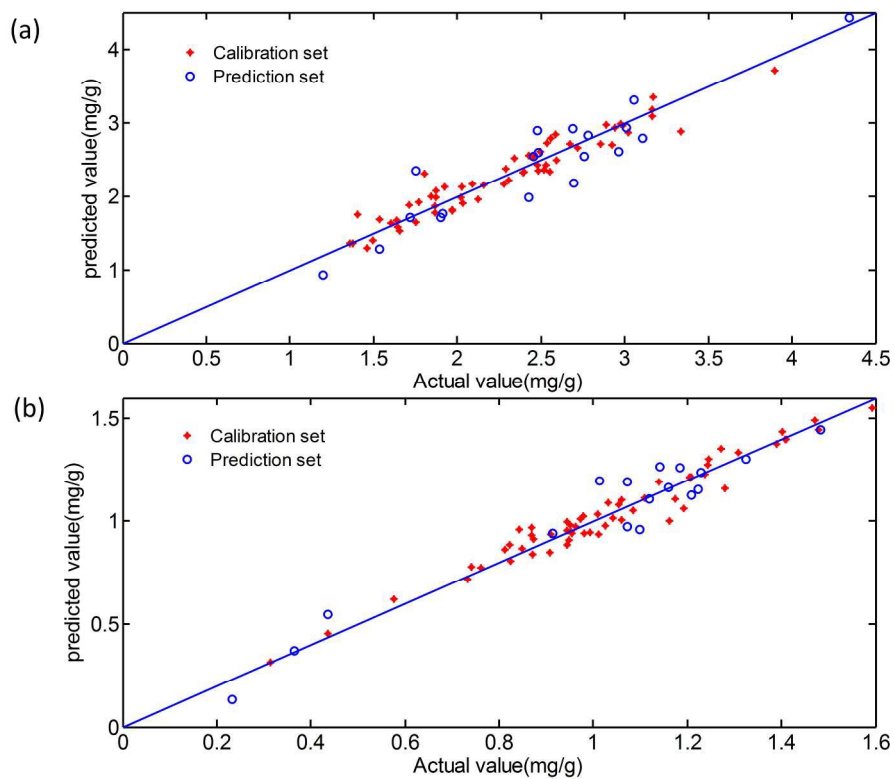
152x64mm (300 x 300 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



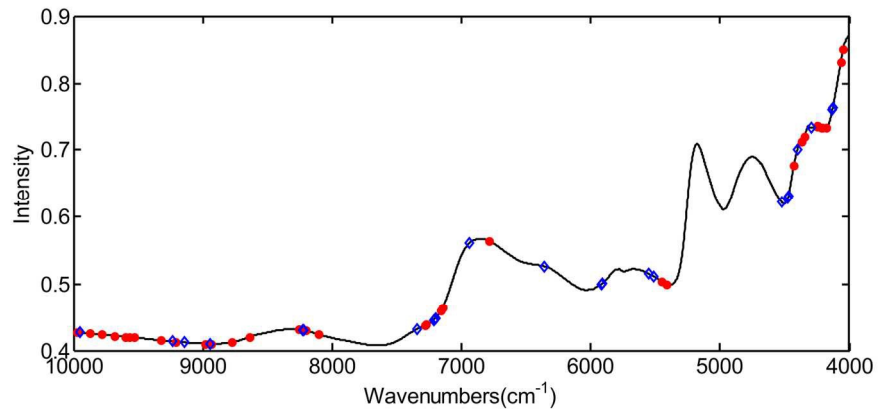
Diagnostic plots for outlier detection based on MCS. (a) TPC; (b) AA.  
231x184mm (300 x 300 DPI)



Correlation diagrams between the predicted values and measured values based on CARS-PLS method.

(a)TPC; (b) AA

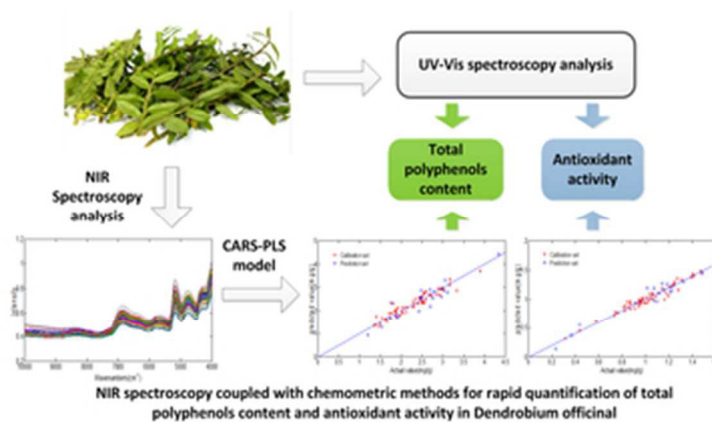
231x184mm (300 x 300 DPI)



Positions of variables selected by CARS-PLS for prediction of TPC (round marked points) and AA (diamond marked points) on the full spectra  
152x64mm (300 x 300 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60





NIR spectroscopy coupled with chemometric methods for rapid quantification of total polyphenols content and antioxidant activity in *Dendrobium officinale*  
36x17mm (300 x 300 DPI)