

Analytical Methods

Accepted Manuscript



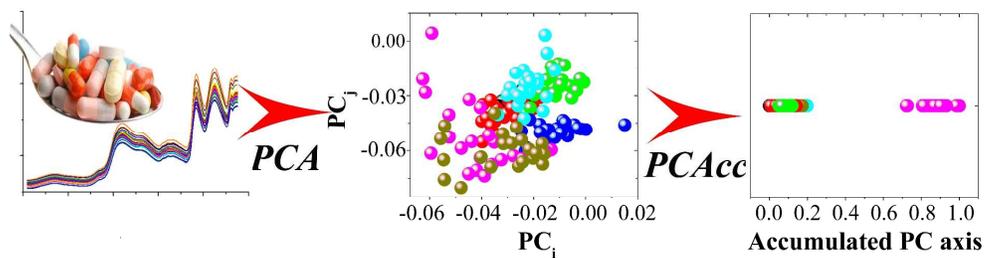
This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

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.

Graphical abstract



Near-infrared (NIR) spectroscopy combined with principal component accumulation (PCAcc) method was used to identify 12 classes of different Chinese patent medicines.

1
2
3
4 1 Discrimination of Chinese patent medicines using near-infrared
5
6 2 spectroscopy and principal component accumulation method
7
8
9 3

10
11 4 Ruifeng Shan^a, Zhiyi Mao^a, Lihui Yin^b, Wensheng Cai^a and Xueguang Shao^{a*}
12

13
14 ^aCollaborative Innovation Center of Chemical Science and Engineering (Tianjin),
15
16 6 State Key Laboratory of Medicinal Chemical Biology, Research Center for Analytical
17
18 7 Sciences, College of Chemistry, Nankai University, Tianjin, 300071, P. R. China

19
20
21 8 ^bNational Institutes for Food and Drug Control, Beijing, 100050, P. R. China
22
23
24 9

25
26
27
28
29 11 Corresponding address:
30

31 12 College of Chemistry,
32

33
34 13 Nankai University,
35

36
37 14 Tianjin, 300071, P. R. China
38

39 15 Tel: +86-22-23503430
40

41 16 Fax: +86-22-23502458
42

43
44 17 E-mail: xshao@nankai.edu.cn
45
46
47 18

48
49
50
51
52
53
54
55
56
57
58
59
60

* Corresponding author.

19 **Abstract**

20 Discrimination of pharmaceutical products has been an important task in
21 pharmaceutical industry and pharmaceutical safety. In this study, principal component
22 accumulation (PCAcc) method was investigated for discrimination of Chinese patent
23 medicines. In PCAcc method, an accumulation strategy is utilized to combine the
24 classification information contained in multiple PC subspaces by using a rotation, a
25 projection and a summation operation. To improve the performance of classification,
26 continuous wavelet transform (CWT) is applied as the pretreatment method to
27 eliminate the background. The results show that, among the 12 classes of Chinese
28 patent medicines, 8 classes are correctly classified, and a total of ten samples are
29 misclassified for the other four classes. Compared with the results obtained by
30 principal component analysis (PCA), radial basis function artificial neural network
31 (RBF-ANN) and partial least squares discriminant analysis (PLSDA), PCAcc
32 produces the best classification.

33

34 Introduction

35 Near-infrared (NIR) spectroscopy is a fast and nondestructive analytical
36 technique and has been widely used in food industry, agriculture, petroleum industry,
37 and etc.^{1,2} In the last decades, the technique has attracted considerable attention in
38 pharmaceutical industry for quantitative analysis, qualitative analysis and on-line
39 control of pharmaceutical products.³⁻⁵ Due to its flexibility in measurement, NIR
40 spectroscopy is suitable for analysis of samples in different pharmaceutical forms.
41 The technique has been extensively studied for quantifying active principal
42 ingredients (API),^{6,7} excipients,⁸⁻¹⁰ and water content¹¹⁻¹³ in pharmaceutical products.
43 Moreover, because NIR spectroscopy has advantages of rapid and nondestructive
44 analysis, it has been used to monitor the production process of pharmaceutical
45 products, e.g., assessing tableting process,¹⁴⁻¹⁶ monitoring blend uniformity of solid
46 dosage forms¹⁷ or API concentration in powder mixing process.¹⁸ On the other hand,
47 NIR spectroscopy has a characteristic that could capture both chemical and physical
48 information of the samples. The parameters of pharmaceutical products such as
49 hardness, particle size, compaction force and dissolution rate can be determined by
50 the technique. NIR spectroscopy was also used to provide the information of
51 polymorphic form,¹⁹⁻²¹ which affects dissolution property of the pharmaceutical
52 products.

53 Discrimination of geographic origins or manufacturer and identification of
54 counterfeit drugs have been an important task in pharmaceutical industry. However,
55 in some cases, e.g., the same pharmaceutical product from different manufacturers,

1
2
3
4 56 there is no significant difference. Therefore, efficient methods are needed to classify
5
6 57 the similar samples by exploring the tiny difference between the products. Pattern
7
8 58 recognition techniques combined with NIR spectroscopy have attracted considerable
9
10 59 attention for variety discrimination. Principal component analysis (PCA) is one of
11
12 60 the most popular and straightforward pattern recognition methods, and has been used
13
14 61 to taxonomic discrimination, quality assessment and discrimination of geographic
15
16 62 origins of pharmaceutical products, and etc.²² For example, PCA was employed to
17
18 63 discriminate three types of Indigowoad Root samples from different origins²³ and to
19
20 64 identify counterfeit drugs.²⁴ In the method, the relation between the samples can be
21
22 65 directly observed by the plot of principal components (PCs). In our recent work,
23
24 66 classification of azithromycin tablets from four manufacturers was studied by PCA
25
26 67 and the effect of morphology was examined by preparing the samples in different
27
28 68 forms.²⁵ The results show that both the samples from different manufacturers and the
29
30 69 samples in different forms can be satisfactorily classified with the help of
31
32 70 chemometric methods. Moreover, least-squares support vector machine (LS-SVM)
33
34 71 was adopted for discrimination of Rhizoma Corydalis and mint tea samples from
35
36 72 different sources. The results demonstrated that the method can find the non-linear
37
38 73 relation between the spectra and predicted properties.^{26,27} Furthermore, K-means
39
40 74 method has been used to discriminate tablets from different manufacturers.²⁸
41
42
43
44
45
46
47
48
49
50

51 The aim of this work is to establish an approach for rapid identification of
52
53
54 76 Chinese patent medicines. NIR spectroscopy was used as a tool for fast and
55
56 77 destructive analytical technique to obtain the information of the samples, and
57
58
59
60

1
2
3
4 78 chemometric methods, including continuous wavelet transform (CWT) and principal
5
6 79 component accumulation (PCAcc), were employed to explore the small difference
7
8
9 80 between the spectra. 12 classes of different Chinese patent medicines or the same
10
11 81 medicine from different manufacturers were studied to demonstrate the performance
12
13
14 82 of the method in discrimination of Chinese patent medicines.
15
16
17
18

19 84 **Experiment and data description**

20
21 85 NIR spectral dataset of Chinese patent medicine is supplied by National
22
23
24 86 Institutes for Food and Drug Control. The dataset includes NIR spectra of five
25
26
27 87 Chinese patent medicines produced by different manufacturers. Table 1 summarized
28
29 88 the information of the samples. The samples of one medicine from one manufacturer
30
31 89 were taken as a class. The capital letters A, B, C, D, and E were used to denote the
32
33
34 90 five medicines and a number following the letter was used to code the manufacturer.
35

36 91 **Table 1**

37
38
39 92 The spectra are divided into calibration and prediction set by Kennard-Stone (KS)
40
41 93 method.²⁹ In order to use the same number of spectra in calibration set for the 12
42
43
44 94 classes of medicines, 22 spectra of each class (a total of $12 \times 22 = 264$ spectra) were
45
46
47 95 used and the remaining spectra as listed in Table 1 were taken as prediction set.
48

49 96 All the spectra were recorded on an MPA FT-NIR spectrometer (Bruker,
50
51 97 Germany) in the wavenumber range $3999.7\text{-}11995.3\text{ cm}^{-1}$ with the digitization interval
52
53
54 98 3.857 cm^{-1} . In the calculations, the variables from 4246.6 to 8913.7 cm^{-1} (1211 data
55
56
57 99 points) were used. Fig. 1(a) displays the measured spectra of the samples. It can be
58
59
60

1
2
3
4 100 seen that most of the spectra are similar and highly overlapped. The reason is that
5
6 101 Chinese patent medicines are mixture of several herbs in composite formulae. Thus,
7
8
9 102 there is no significant difference between the medicines. On the other hand, the
10
11 103 chemical constituents in component herbs may vary with harvest season, geographic
12
13
14 104 origin, drying processes and other factors. This may cause the difference between the
15
16 105 samples of the same medicine from different manufacturers.

106 **Figure 1**

107 Generally, signal preprocessing methods such as multiplicative scattering
108 correction (MSC),³⁰ standard normal variate (SNV)³¹ and derivative are used for
109 correcting the scattering effect and background removal. In our previous works, CWT
110 has been proved to be an efficient tool for removing the variant background and
111 noise.^{32,33} Therefore, CWT is applied as the pretreatment method to eliminate the
112 background in this work. In the calculation of CWT, Haar wavelet with a scale
113 parameter 20 was used. Fig. 1(b) shows the preprocessed spectra. It can be seen that
114 the variant background is removed compared with the spectra in Fig. 1(a). However,
115 the spectra are still overlapped. Therefore, it is impossible to distinguish the 12 classes
116 of the medicines directly by the spectra, although there are differences between these
117 spectra.

118 119 **Calculations**

120 PCA has been the basic method for classification or discrimination analysis. In
121 PCA, the information contained in first two or three PCs are generally used for

1
2
3
4 122 inspection of classification. However, high-order PCs may contain the classification
5
6 123 information. In order to use the information sufficiently, PCAcc was proposed in our
7
8
9 124 previous works.^{34,35} The essential of PCAcc method is an accumulation strategy to
10
11 125 accumulate the classification information contained in multiple PC subspaces.
12
13
14 126 Therefore, the difference between the spectra of the samples in any PC subspaces is
15
16 127 used for the classification. A rotation, projection and summation operations are
17
18
19 128 included in the calculations. For building a PCAcc classification model, PCA is
20
21 129 applied on the calibration set. In order to explore the information contained in
22
23
24 130 different PC subspaces, a large number of PCs can be used. By using the information
25
26 131 in PC subspaces, a decision tree can be obtained, in which each node has two
27
28
29 132 branches. One branch contains the samples of one class and the other one contains the
30
31 133 samples of remaining classes. For each decision node, the class with the largest
32
33
34 134 difference from the others is separated out. The process is repeated until only one
35
36 135 class remains. The classifier in each node is built with the accumulation, which
37
38
39 136 includes the following operations: (1) finding the axis maximizing the distance
40
41 137 between one class and the other classes using Fisher criterion³⁶ in each PC subspace;
42
43
44 138 (2) rotating all the PC subspaces to the same direction; and (3) accumulating the
45
46 139 information of the effective subspaces. The effective subspace is defined as the one
47
48
49 140 producing a “minimal increase”³⁴ to the classification. In the end, an accumulation
50
51 141 sequence (of the PC subspaces) and a threshold to produce the best classification can
52
53
54 142 be obtained as the classifier of the node in the decision tree.

55
56 143 For predicting an unknown sample, a decision can be obtained by testing the
57
58
59
60

1
2
3
4 144 sample along the decision tree. In each node, the spectrum of the sample is projected
5
6 145 into the PC subspaces, rotated, and then an accumulation is performed according to
7
8
9 146 the sequence. The classification or discrimination will end when the sample is
10
11 147 classified to the end node, i.e., the node containing only one class, using the classifier
12
13
14 148 and the threshold.

15
16 149 Detail procedures of PCAcc method can be found in our previous works.^{34,35} In
17
18
19 150 this paper, resolution is still employed as a quantitative measure of the difference
20
21 151 between classes, which is defined by³⁴

22
23
24 152
$$R_s = \frac{|m_A - m_B|}{s_A + s_B} \quad (1)$$

25
26

27 153 where R_s is the resolution between class A and the other classes (denoted as B), m_A ,
28
29 154 m_B , s_A and s_B are the mean values and standard deviations of the two classes,
30
31
32 155 respectively.
33
34

35 156

37 157 **Results and discussion**

40 158 **Discrimination using PCA**

41
42 159 PCA is the most commonly used unsupervised pattern recognition method. In
43
44
45 160 this work, PCA is employed to investigate the classification of the medicines. The
46
47
48 161 result shows that the first four PCs explain more than 98% of the variance. Therefore,
49
50 162 most information of the spectra is included in the first four PC subspaces. Fig. 2(a)
51
52 163 and (b) shows the distribution of the calibration samples in PC1-PC2 and PC3-PC4
53
54
55 164 subspace, respectively. It can be seen that, in PC1-PC2 subspace, 5 classes of the
56
57
58 165 medicines can be separated, including E1, E2, C1, A3 and A4, and in Fig. 2(b), A3 is
59
60

1
2
3
4 166 almost separable from the others in PC3-PC4 subspace. For the samples of the other
5
6 167 classes, however, it is too difficult to separate them by the four PCs. This result
7
8 168 evidently indicates that the difference between the samples may be contained in the
9
10 169 high order PCs although explaining very small variance in the spectral data.
11
12 170 Furthermore, even for the separable classes, it is obvious that the in-class variance is
13
14 171 much larger than the between-class variance. The result clearly demonstrates the
15
16 172 difficulty in classification of Chinese patent medicines by using PCA.
17
18
19
20

21 **Figure 2**

22 **Building PCAcc model**

23
24
25
26 175 In order to use the comprehensive information contained in the PCs to improve
27
28 176 the classification, PCAcc model was studied for discrimination of the medicines using
29
30 177 the spectra in the calibration set. In the method, all the possible information for the
31
32 178 classification contained in multiple PCs is used. In order to use more information, 12
33
34 179 PCs were used in the calculations. Therefore, a total of 66 PC subspaces are included
35
36 180 in the accumulation.
37
38
39
40

41 181 To demonstrate the effect of the accumulation, the variation of R_s in the
42
43 182 accumulation process for the last node to separate the class D1 and D2 is shown in Fig.
44
45 183 3(a). Clearly, 21 PC subspaces that increase the separation of the two classes are
46
47 184 accepted and the R_s value increases from 0.97 (the best PC subspace) to 1.60 (the
48
49 185 accumulated value). The result means that 21 of the 66 PC subspaces contain the
50
51 186 effective information for the separation and acceptable classification of the two
52
53 187 classes is obtained. Statistically, when the value of R_s is above 1.5, it can be known as
54
55
56
57
58
59
60

1
2
3
4 188 a complete separation.
5

6 189 **Figure 3**
7

8
9 190 Fig. 3(b) shows the R_s values of the accepted PC subspaces in the accumulation.

10
11 191 Clearly, even for the accepted subspaces, the R_s values vary significantly. The largest

12
13 192 value can be as high as 0.97 and the smallest one can be as low as 0.12. All the values

14
15 193 are lower than 1.0 and more than half of the values are lower than 0.5. This indicates

16
17 194 that the discriminating information in an individual PC subspace is limited. Therefore,

18
19 195 the accumulation is necessary. Moreover, a sharp increase can be seen in Fig. 3(a)

20
21 196 when the PC subspace No. 9, 10, 13 and 14 was accumulated. This indicates that the

22
23 197 PC subspaces with a higher R_s value may have a significant contribution to the

24
25 198 accumulation. However, there are also cases that the accumulation of the PC

26
27 199 subspaces with higher R_s value does not produce significant contribution to the

28
29 200 separation, e.g., the PC subspaces No. 16, 18, and 20 in the figure. This maybe

30
31 201 accounted for by the fact that the information contained in these PCs is similar with

32
33 202 that in the previously accepted PCs.
34

35
36 203 **Figure 4**
37

38
39 204 Figs. 4(a) - (k) shows the discrimination sequence of 12 classes of medicines.

40
41 205 The balls in the bottom line display the situation of the calibration samples. The long

42
43 206 vertical line in the center denotes the threshold of the classification, and two short

44
45 207 vertical lines denote the mean values of the two classes. The position of the long line,

46
47 208 i.e., the threshold, is determined by the two short lines, locating at the middle of the

48
49 209 two lines. Moreover, the accumulated R_s values are labeled in the figure.
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 210 Fig. 4 shows the sequence of the separation in the training process. Clearly, C1 is
5
6 211 the first class to be separated. This is because the R_s value between C1 and the other
7
8
9 212 11 classes is the largest one. After removing the samples of C1, B2 is selected as the
10
11 213 second class to be separated with the reason that the R_s value between B2 and the
12
13
14 214 other 10 classes is the largest one. The remaining classes are separated in an order of
15
16 215 B1, A2, A4, A3, E2, E1, A5, A1, D1 and D2. The sequence forms a decision tree with
17
18 216 11 nodes. From the R_s values labeled in the figure, it can be known that all the values
19
20
21 217 are bigger than 1.5, indicating a good classification. With detail examination of the
22
23
24 218 figures it can be found that all the samples are correctly discriminated except for one
25
26 219 sample in class D1 and three samples in class D2, as shown in Fig. 4 (g) and (k),
27
28
29 220 respectively.

30
31 221

32 33 34 222 **Discrimination of the prediction samples**

35
36 223 In order to validate the efficiency of the PCAcc model for the 12 class medicines,
37
38 224 discrimination was performed using the spectra of the prediction set. Along the
39
40
41 225 decision tree, a spectrum is repeatedly identified with the classifier in the node until
42
43
44 226 the sample is classified into a class. The operation for the identification in each node
45
46 227 includes the projection into the accepted PC subspaces in the node, the rotation, the
47
48
49 228 accumulation, and then identification with the threshold. The results for the samples
50
51 229 in the prediction set are shown by the balls in the upper line in Figs. 4(a) - (k). From
52
53
54 230 the figures, it can be seen that all the samples are correctly classified except for one
55
56 231 sample in class E2, one sample in class A5 and eight samples in class D2, as shown in

1
2
3
4 232 Figs. 4 (g), (i) and (k), respectively. Further investigation shows that all the 10
5
6 233 samples are misclassified from class D1. Therefore, the large diversity of the samples
7
8
9 234 in class D1 is the reason for the misclassification. From Table 1 it can be seen that the
10
11 235 samples in class D1 and D2 are same medicine from different manufacturers. This is
12
13
14 236 the reason for explanation of the eight samples misclassified from D1 into D2.
15

16
17 **Table 2**

18
19 238 To further investigate the performance of the method, the values of the true
20
21 239 positive (TP) and false positive (FP) obtained by PCAcc, PCA, radial basis function
22
23 240 artificial neural network (RBF-ANN) and partial least squares discriminant analysis
24
25
26 241 (PLSDA) were summarized in Table 2. The two parameters are generally used to
27
28
29 242 evaluate the performance of a classifier, which are defined as the ratio of the number
30
31 243 of correctly classified and misclassified samples, respectively, to the total number of
32
33
34 244 the samples in the class. From the table, it can be seen that PCAcc method produces
35
36 245 the best result for the prediction set. Among the 12 classes of Chinese patent
37
38
39 246 medicines, 8 classes are correctly classified, the true positive accuracies for the other
40
41 247 four classes (E2, A5, D2 and D1) are 100%, 100%, 100%, 67.7%, and the false
42
43
44 248 positive accuracies are 10%, 9.1%, 100%, 0.0%, respectively. Clearly, the results for
45
46
47 249 all the 12 classes are acceptable except for the true positive accuracies of class D1 and
48
49 250 the false positive accuracies of class D2. However, only five classes can be classified
50
51 251 by PCA, and it is difficult to obtain acceptable results by RBF-ANN and PLSDA,
52
53
54 252 because the true positive accuracies for some classes are lower than 50% and the false
55
56
57 253 positive accuracies are even higher than 60%.
58
59
60

1
2
3
4 254

5
6 **Conclusions**
7

8
9 256 Discrimination of the 12 classes of Chinese patent medicines was studied using
10
11 257 NIR spectroscopy and PCAcc method. CWT was adopted to eliminate the variant
12
13 258 background in the NIR spectra. Because PCAcc method uses the accumulation of the
14
15 259 information contained in multiple PC subspaces, an acceptable classification was
16
17 260 achieved for different medicines or the same medicine from different manufacturers.
18
19 261 Due to the advantage of the PCAcc in exploring as much as the classification
20
21 262 information in the NIR spectra of the samples, PCAcc produced the best classification
22
23 263 compared with the results of PCA, RBF-ANN and PLSDA.
24
25
26
27

28
29 264

30
31 **Acknowledgements**
32

33
34 266 This study was supported by National Natural Science Foundation of China (No.
35
36 267 21175074).
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

268 **References**

- 269 1. Y. Roggo, P. Chalus, L. Maurer, C. Lema-Martinez, A. Edmond and N. Jent, *J.*
270 *Pharm. Biomed. Anal.*, 2007, **44**, 683–700.
- 271 2. W.J. Dong, Y.N. Ni and S. Kokot, *J. Agric. Food Chem.*, 2013, **61**, 540–546.
- 272 3 L. Tomuta, L. Rus, R. Lovanov and L. L. Rus, *J. Pharm. Biomed. Anal.*, 2013, **84**,
273 285–292.
- 274 4 O. Scheibelhofer, N. Balak, P. R. Wahl, D. M. Koller, B. J. Glasser and J. G. Khinast,
275 *AAPS Pharmscitech*, 2013, **14**, 234-244.
- 276 5 D. S. Bu, B. Y. Wan and G. McGeorge, *Chemom. Intell. Lab. Syst.*, 2013, **120**,
277 84-91.
- 278 6 M. Boiret, L. Meunier and Y-M. Ginot, *J. Pharm. Biomed. Anal.*, 2011, **54**,
279 510–516.
- 280 7 S. S. Rosa, P. A. Barata, J. M. Martins and J.C. Menezes, *Talanta*, 2008, **75**,
281 725–733.
- 282 8 M. Blanco and A. Peguero, *TRAC-Trend. Anal. Chem.*, 2010, **29**, 1127–1136.
- 283 9 A. D. Karande, P. W. S. Heng and C. V. Liew, *Int. J. Pharm.*, 2010, **396**, 63–74.
- 284 10 M. Blanco, R. Cueva-Mestanza and A. Peguero, *Talanta*, 2011, **85**, 2218–2225.
- 285 11 H. Grohganz, D. Gildemyn, E. Skibsted, J. M. Flink and J. Rantanen, *Anal. Chim.*
286 *Acta*, 2010, **676**, 34–40.
- 287 12 C. C. Corredor, D. Bu and D. Both, *Anal. Chim. Acta*, 2011, **696**, 84–93.
- 288 13 X. B. Zhang, Y. C. Feng and C. Q. Hu, *Anal. Chim. Acta*, 2008, **630**, 131–140.
- 289 14 J. Moes, M. M. Ruijken, E. Gout, H. W. Frijlink and M. Ugwoke, *Int. J. Pharm.*,

- 1
2
3
4 290 2008, **357**, 108–118.
5
6 291 15 M. Blanco, M. Alcalá, J.M. Gonzales and E. Torras, *J. Pharm. Sci.* 2006, **95**,
7
8 292 2137–2144.
9
10 293 16 M. Blanco and A. Peguero, *J. Pharm. Biomed. Anal.*, 2010, **52**, 59–65.
11
12 294 17 Y. Sulub, B. Wabuyele, P. Gargiulo, J. Pazdan, J. Cheney, J. Berry, A. Gupta, R.
13
14 295 Shah, H. Wu and M. Khan, *J. Pharm. Biomed. Anal.*, 2009, **49**, 48–54.
15
16 296 18 A.U. Vanarase, M. Alcalá, J. I. Jerez Rozo, F. J. Muzzio, R. J. Romanach, *Chem.*
17
18 297 *Eng. Sci.*, 2010, **65**, 5728–5733.
19
20 298 19 C. M. McGoverin, L. C. H. Ho, J. A. Zeitler, C. J. Strachan, K. C. Gordon and T.
21
22 299 Rades, *Vib. Spectrosc.*, 2006, **41**, 225–231.
23
24 300 20 Y. Hu, A. Erxleben, A. G. Ryder and P. McArdle, *J. Pharm. Biomed. Anal.*, 2010,
25
26 301 **53**, 412–420.
27
28 302 21 A. Heinz, M. Savolainen, T. Rades and C. J. Strachan, *Eur. J. Pharm. Sci.*, 2007, **32**,
29
30 303 182–192.
31
32 304 22 H. A. Gad, S. H. El-Ahmady, M. I. Abou-Shoer and M. M. Al-Azizi, *Phytochem.*
33
34 305 *Anal.* 2013, **24**, 1-24.
35
36 306 23 Y. N. Ni, R. M. Song and S. Kokot, *Spectrochim. Acta A*, 2012, **96**, 252-258.
37
38 307 24 S. H. F. Scafi and C. Pasquini, *Analyst*, 2001, **126**, 2218–2224.
39
40 308 25 P. Li, G. R. Du, W. S. Cai and X. G. Shao, *J. Pharm. Biomed. Anal.*, 2010, **70**,
41
42 309 288–294.
43
44 310 26 Y. H. Lai, Y. N. Ni and S. Kokot, *Vib. Spectrosc.*, 2011, **56**, 154-160.
45
46 311 27 W.J. Dong, Y.N. Ni and S. Kokot, *Appl. Spectros.*, 2014, **68**, 245-254.
47
48
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3
4 312 28 P. Gemperline, *CRC/Taylor & Francis, Boca Raton.*, 2006.
5
6 313 29 R. W. Kennard and L. A. Stone, *Technometrics*, 1969, **11**, 137–148.
7
8
9 314 30 R. J. Barnes, M. S. Dhanoa and S. J. Lister, *Appl. Spectrosc*, 1989, **43**, 772-777.
10
11 315 31 P. Geladi, D. MacDougall and H. Martens, *Appl. Spectrosc*, 1985, **39**, 491-500.
12
13 316 32 X. G. Shao, A. K. M. Leung and F. T. Chau, *Acc. Chem. Res.*, 2003, **36**, 276–283.
14
15
16 317 33 X. G. Shao, G. R. Du, M. Jing and W. S. Cai, *Chemom. Intell. Lab. Syst*, 2012, **114**,
17
18 318 44-49.
19
20
21 319 34 J. J. Liu, W. S. Cai and X. G. Shao, *Sci. China: Chem.*, 2011, **54**, 802–811.
22
23
24 320 35 Y. Wang, X. Ma, Y. D. Wen, J. J. Liu, W. S. Cai and X. G. Shao, *Anal. Methods*,
25
26 321 2012, **4**, 2893-2899.
27
28
29 322 36 R. A. Fisher, *Ann. Hum. Genet.*, 1936, **7**, 179–188.
30
31 323
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

324 **Figure captions**

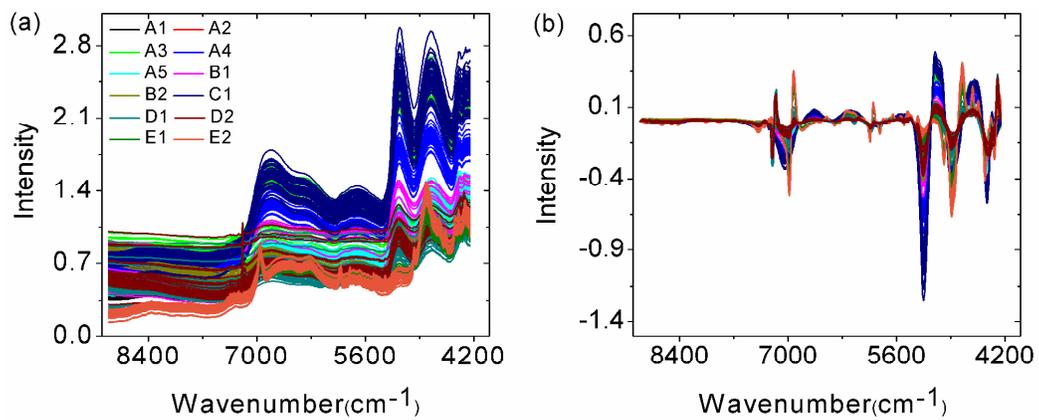
325 Fig. 1 Measured (a) and preprocessed (b) spectra of 12 classes of Chinese patent
326 medicines.

327 Fig. 2 Distribution of the calibration samples in PC1-PC2 (a) and PC3-PC4 (b)
328 subspaces for 12 classes of medicines.

329 Fig. 3 Resolution parameter (R_s) of the accepted subspaces (a) and their accumulated
330 effect (b) in the discrimination of D1 and D2.

331 Fig. 4 Distribution of the samples along the accumulated PC axis for the calibration
332 and prediction set of the medicines.

333



334

335

Figure 1

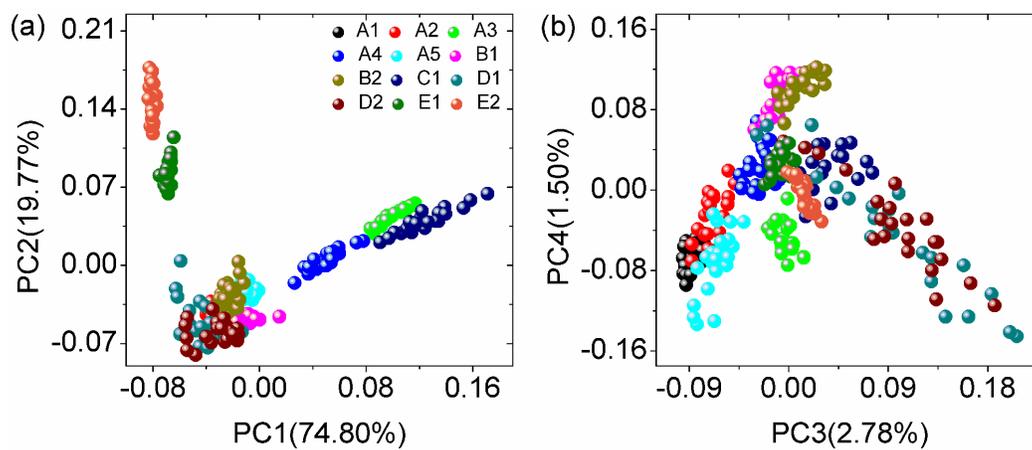


Figure 2

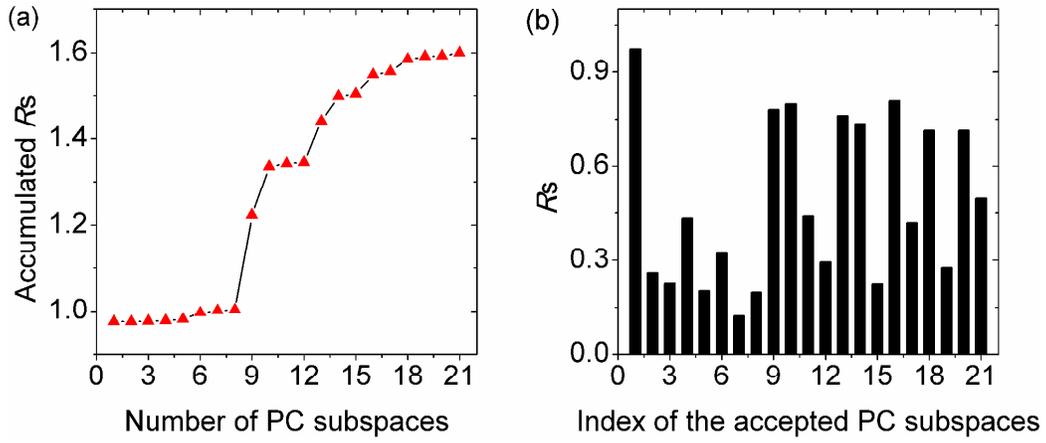


Figure 3

339

340

341

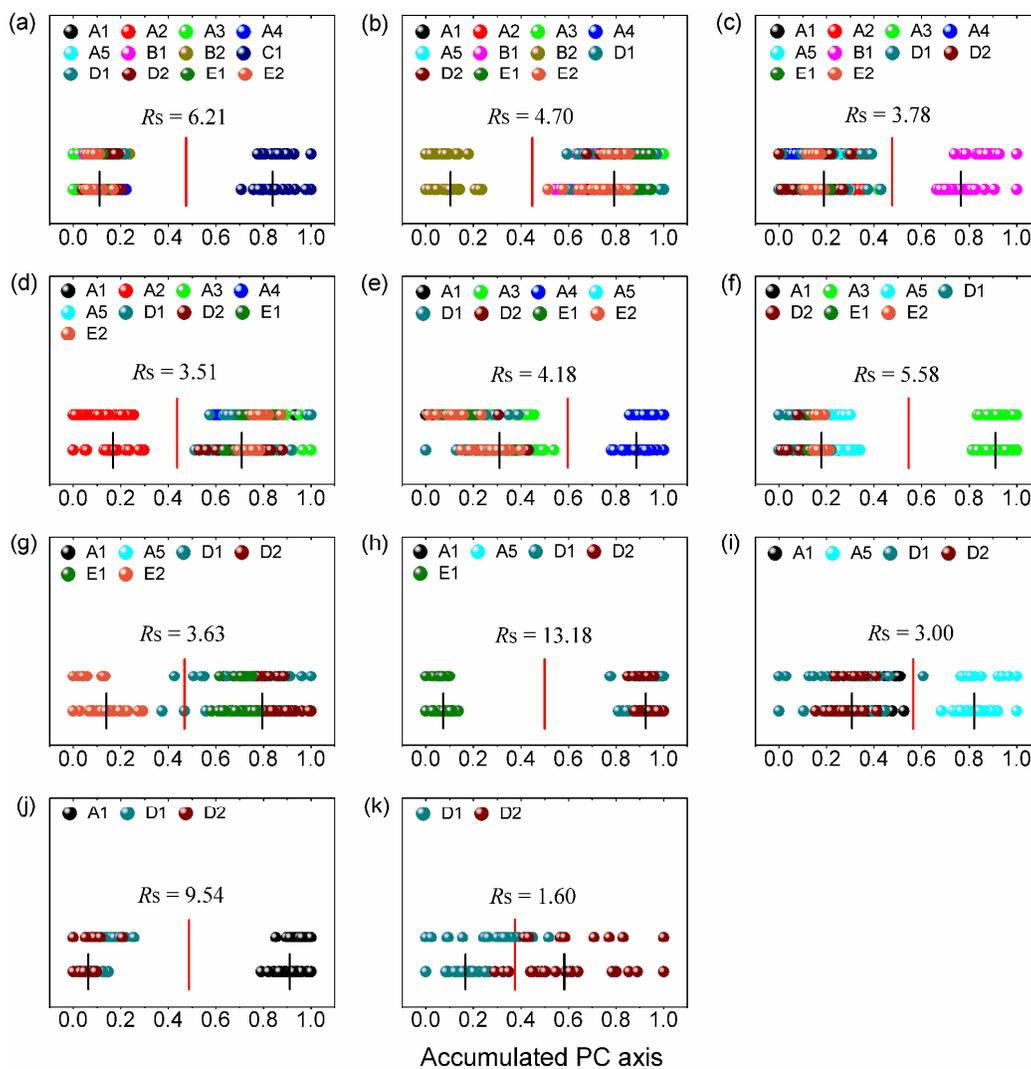


Figure 4

342

343

344

345 Table 1 Information of the samples in the calibration and prediction set

Medicine	Manufacturer	Class Label	Number of samples	Calibration set	Prediction set
	1	A1	40	22	18
	2	A2	45	22	23
A	3	A3	45	22	23
	4	A4	40	22	18
	5	A5	33	22	11
B	1	B1	42	22	20
	2	B2	47	22	25
C	1	C1	52	22	30
D	1	D1	53	22	31
	2	D2	30	22	8
E	1	E1	39	22	17
	2	E2	32	22	10
Total			498	264	234

346

347

348 Table 2 Comparison of the classification accuracy by different methods.

Class	Method ^a	Accuracy (TP, %)		Accuracy (FP, %)	
		Calibration set	Prediction set	Calibration set	Prediction set
C1	PCAcc	100.0	100.0	0.0	0.0
	PCA	100.0	96.7	0.0	3.3
	RBF-ANN	77.3	76.7	22.7	30.0
	PLSDA	86.4	46.7	31.8	33.3
B2	PCAcc	100.0	100.0	0.0	0.0
	RBF-ANN	77.3	92.0	22.7	24.0
	PLSDA	72.7	80.0	27.3	68.0
B1	PCAcc	100.0	100.0	0.0	0.0
	RBF-ANN	90.9	80.0	22.7	0.0
	PLSDA	59.1	80.0	22.7	25.0
A2	PCAcc	100.0	100.0	0.0	0.0
	RBF-ANN	45.5	52.2	45.5	21.7
	PLSDA	54.6	65.2	54.6	43.5
A4	PCAcc	100.0	100.0	0.0	0.0
	PCA	90.9	100.0	0.0	0.0
	RBF-ANN	72.7	83.3	45.5	44.4
	PLSDA	72.7	66.7	50.0	50.0
A3	PCAcc	100.0	100.0	0.0	0.0
	PCA	100.0	100.0	9.1	0.0
	RBF-ANN	90.9	100.0	63.6	43.5
	PLSDA	81.8	56.5	50.0	43.5
E2	PCAcc	100.0	100.0	4.6	10.0
	PCA	100.0	100.0	0.0	0.0
	RBF-ANN	95.5	100.0	0.0	0.0
	PLSDA	86.4	100.0	18.2	30.0
E1	PCAcc	100.0	100.0	0.0	0.0
	PCA	100.0	100.0	0.0	5.9
	RBF-ANN	95.5	100.0	4.6	0.0
	PLSDA	68.2	82.4	31.8	11.8
A5	PCAcc	100.0	100.0	0.0	9.1
	RBF-ANN	45.5	18.2	18.2	45.5
	PLSDA	54.6	54.6	31.8	18.2
A1	PCAcc	100.0	100.0	0.0	0.0
	RBF-ANN	54.6	72.2	13.6	22.2
	PLSDA	54.6	72.2	13.6	27.8
D1	PCAcc	95.5	67.7	13.6	0.0
	RBF-ANN	59.1	54.8	27.3	6.5
	PLSDA	50.0	48.4	31.8	9.7
D2	PCAcc	86.4	100.0	0.0	100.0
	RBF-ANN	81.8	100.0	27.3	75.0
	PLSDA	59.1	75.0	36.4	50.0

^a PC1-PC2 was used in PCA method, and the same latent variable number as in PCAcc was used in PLSDA. 33 hidden neurons were used in RBF-ANN. For the classes that can not be discriminated by PCA, the result was not listed.

349