

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.

Iteratively variable subset optimization for multivariate calibration

Weiting Wang,^a Yonghuan Yun,^{a,**} Baichuan Deng,^{a,b} Wei Fan,^c Yizeng Liang,^{a,*}

^aCollege of Chemistry and Chemical Engineering, Central South University, Changsha

410083, P.R. China

^bDepartment of Chemistry, University of Bergen, Bergen N-5007, Norway

^cJoint Lab for Biological quality and safety, College of Bioscience and Biotechnology, Hunan

Agriculture University, Changsha 410128, P.R. China

Abstract: Based on the theory that a large partial least squares (PLS) regression coefficient on the autoscaled data indicates an important variable, a novel strategy for variable selection called iteratively variable subset optimization (IVSO), is proposed in this study. In addition, we take it into consideration that the optimal number of latent variables generated by cross-validation will make a great difference to the regression coefficients and sometimes the difference can even vary by several orders of magnitude. In this work, the regression coefficients generated in every sub-model are normalized to remove the influence. In each iterative round, the regression coefficients of each variable obtained from the sub-models are summed to evaluate its importance level. A two-step procedure including weighted binary matrix sampling (WBMS) and sequentially addition is employed to eliminate uninformative variables gradually and gently in a competitive way and reduce the risk of losing important variables. Thus,

*Corresponding author. Tel.:86-0731-88830824; fax: +86-0731-88830831.

E-mail address: yizeng_liang@263.net (Yizeng Liang)

** Corresponding author. E-mail address: yunyonghuan@foxmail.com (Yonghuan Yun)

20 IVSO can achieve high stability. Investigated by one simulated dataset and two NIR
21 datasets, IVSO shows much better prediction ability than another two outstanding and
22 commonly used methods, Monte Carlo uninformative variable elimination (MC-UVE)
23 and competitive adaptive reweighted sampling (CARS). The MATLAB code for
24 implementing IVSO is available in the supplemental materials.

25 **Keywords:** Partial least squares, Regression coefficient, Weighted binary matrix
26 sampling, Sequentially addition, Variable selection

27

28 1. Introduction

29 Nowadays, multivariate calibration models have been playing an essential role in
30 multi-component spectral data, such as ultraviolet (UV), near infrared (NIR) and
31 Raman spectroscopy. However, the spectral data obtained from these modern
32 spectroscopic instruments usually contain hundreds or thousands of variables with high
33 colinearity. Latent variable extraction techniques, such as principal component
34 regression (PCR) and partial least squares (PLS) ¹, provide a way to address the high
35 colinearity problem. But the full spectrum used in these methods may bring negative
36 influence on the performance of the calibration model due to the existing of
37 uninformative variables. Many papers have demonstrated that it is critical to conduct
38 variable selection in models instead of using full spectrum.²⁻⁶ The advantages of
39 variable selection have been concluded in the following three aspects: (1) improve the
40 prediction accuracy of the model because of the elimination of uninformative variables
41 that must lead to less precision as proved theoretically; (2) selecting wavelengths
42 probably responsible for the property of interest makes the model more interpretative;
43 (3) enhance the computational efficiency for modeling with a small amount of

44 variables.⁷

45 At present, many methods on variable selection have been employed in
46 multi-component spectral data. In general, these methods can be classified into two
47 categories as static and dynamic approach. The static approaches use one criterion for
48 the whole data space, while the dynamic approaches take into account the result of
49 each iteration. The static approaches includes t-statistics and Akaike information
50 criteria (AIC), uninformative variable elimination (UVE),⁸ Monte Carlo based
51 uninformative variable elimination (MC-UVE),^{9,10} variable importance in projection
52 (VIP),¹¹ selectivity ratio (SR),¹² and moving window partial least squares (MWPLS).¹³
53 The dynamic approaches consists of optimized algorithm-based such as Genetic
54 algorithm (GA),¹⁴⁻¹⁶ particle swarm optimization (PSO),¹⁷ firefly,¹⁸ ant colony
55 optimization (ACO),^{19, 20} gravitational search algorithm (GSA),²¹ and simulated
56 annealing (SA).²² The variable selection methods, such as Random forest,²³ successive
57 projection algorithm (SPA),²⁴ iteratively retaining informative variables (IRIV),²⁵
58 variable combination population analysis (VCPA),²⁶ competitive adaptive reweighted
59 sampling (CARS),²⁷ interval random frog (*i*RF),²⁸ are also the dynamic approaches.
60 The theories of UVE, MC-UVE, CARS, and *i*RF comes from that the larger the
61 absolute regression coefficient on the autoscaled data is, the more important the
62 variable is.^{8,29} In addition to regression coefficient, Kvalheim et al. discussed the usage of SR that
63 can assist in improved algorithm for variable selection in latent variable regression model.³⁰

64 Among all the methods upon regression coefficient, MC-UVE and CARS are
65 adopted extensively in multivariate calibration models for their better prediction. In
66 both MC-UVE and CARS, Monte Carlo sampling technique is applied to the sample
67 space to establish a large number of sub-models, which assures that the number of
68 samples selected randomly for modeling is strictly the same, for example, 80% of all

69 samples is used to build the model. For MC-UVE, after N Monte Carlo sampling runs,
70 one variable is evaluated according to a criterion which is equal to the ratio of the mean
71 of the regression coefficients and its standard deviation. The variables with small
72 criteria are eliminated. Unlike MC-UVE, in each iterative round, CARS removes the
73 variables with small means of regression coefficients by the exponentially decreasing
74 function (EDF) by force and adaptive reweighted sampling (ARS) competitively.
75 However, it is the full spectrum in MC-UVE that is used to establish sub-models, which
76 will lead to that the regression coefficients of the informative variables can be
77 influenced by the uninformative variables. With regard to CARS, the enforced
78 elimination of variables by EDF may lose important variables and further result in
79 instability. Hence, in most cases the results achieved by MC-UVE and CARS are not
80 satisfied enough.

81 In this study, a novel strategy for variable selection based on regression coefficient
82 is proposed, called iteratively variable subset optimization (IVSO). At first, we
83 introduce a new random sampling method, named weighted binary matrix sampling
84 (WBMS),^{31, 32} which is an improvement of the binary matrix sampling (BMS).^{25, 33}
85 Giving different weights to different variables, WBMS aims to make variables with
86 larger weights more likely to be chosen. On the contrary, if the weight of one variable is
87 small, it will be selected with little or even no possibility. Furthermore, combining
88 WBMS with another strategy called sequentially addition, the variables with small
89 criteria are deleted and a new variable subset is yielded. After N WBMS runs, N
90 different variable subsets are obtained and the root mean squares error of

91 cross-validation (RMSECV) is used as the objective function to search for the best
92 variable subset. In addition, the regression coefficients of one variable in all
93 sub-models are summed to evaluate its importance level. This data fusion step is a
94 good option, as the noise cancels out and the systematic information accumulates.
95 However, we find that the optimal number of latent variables generated by
96 cross-validation will make a great difference to the regression coefficients, which is
97 consistent with the viewpoint in Reference.³⁴ Thus, the regression coefficients of the
98 same variable in different sub-models cannot be calculated or compared directly due to
99 the great difference. The strategy of normalization is applied to eliminate the influence.
100 Tested on a simulated dataset and two NIR datasets, IVSO coupled with partial least
101 squares (PLS) demonstrates better prediction ability and higher stability compared to
102 the two outstanding methods above, namely MC-UVE and CARS. The result
103 demonstrates that IVSO has the ability to eliminate uninformative variables gradually
104 and gently in a competitive way, which can avoid those two problems of MC-UVE and
105 CARS discussed above. It proves that IVSO is an efficient method for variable
106 selection in multivariate calibration.

107 Additionally, it should be noted that IVSO is just evaluated by NIR datasets with
108 PLS in this study, but it is a general strategy and can be combined with other regression
109 and pattern recognition methods, and applied to other kinds of datasets, such as
110 metabolomic and quantitative structure-activity relationship (QSAR).

111

112 **2. Theory and algorithms**

113 2.1. Notation

114 In this study, the matrix \mathbf{X} with dimensionality $K \times P$ represents the observation
115 matrix, in which K stands for the number of samples in rows and P is the number of
116 variables in columns. Vector \mathbf{y} with dimensionality $P \times 1$ denotes the measured
117 property of interest, for example the concentration. In addition, the number of WBMS
118 runs is set to N .

119 2.2. Weighted binary matrix sampling (WBMS)

120 In IVSO, a new method called WBMS is introduced for randomly sampling and
121 further eliminating a part of uninformative variables, which is an improvement of
122 binary matrix sampling (BMS). If the weight of one variable is small, the variable will
123 be selected with little or even no possibility. Therefore, WBMS can eliminate variables
124 competitively.

125 It works as follows: assume that the weight of the i th variable is w_i . At first, a
126 binary matrix \mathbf{M} with dimensionality $N \times P$ is generated, which contains either '1' or
127 '0'. '1' represents that the responding variable is included for modeling, while '0'
128 represents non-sampling for the variable. In each column, there are Nw_i '1' and the left
129 ones are all '0'. The procedure is displayed in Fig. 1, where the row of \mathbf{M} is set to 5 and
130 the column is 7 for simplicity. When sampling, the weights of some variables are too
131 small to be sampled in any column. The first and second columns in Fig. 1 can
132 represent this case. As the last column shows, if the weight of one variable is equal to
133 1, it will be sampled in each iterative round. Next, permuting each column in \mathbf{M}
134 generates a new binary matrix \mathbf{NM} . Remarkably, after the permutation, the number of
135 '1' or '0' in each column is kept unchanged.

136 In the matrix of \mathbf{NM} , each row represents a sampling process for building a
137 sub-model. It can be summarized that when Nw_i of one variable is less than 1, it will be

138 eliminated.

139 (Insert Figure 1)

140 **2.3. Normalizing PLS regression coefficients**

141 PLS is one of the most widely used methods for establishing the relationship
 142 between the observation matrix \mathbf{X} and the properties of interest \mathbf{y} . The scores matrix \mathbf{T}
 143 is a linear combination of \mathbf{X} with the combination coefficients \mathbf{W} , and \mathbf{c} is the
 144 regression coefficient vector of \mathbf{y} against \mathbf{T} by least squares.^{1,35} PLS can be expressed
 145 by by the following formulas:

$$146 \quad \mathbf{T} = \mathbf{XW} \quad (1)$$

$$147 \quad \mathbf{y} = \mathbf{Tc} + \mathbf{e} = \mathbf{XWc} + \mathbf{e} = \mathbf{Xb} + \mathbf{e} \quad (2)$$

148 where $\mathbf{b} = \mathbf{Wc}$ is the vector of PLS regression coefficients and \mathbf{e} is the vector of
 149 residuals that cannot be explained by the model.

150 In addition, the matrix \mathbf{X} needs to be autoscaled to guarantee that each variable has
 151 the same variance before modeling. It should be noted that the regression coefficients
 152 mentioned in this study have been changed into the absolute value before calculating.
 153 Afterwards, the larger the regression coefficient is, the more important the variable is.

154 Moreover, it is found that the optimal number of latent variables generated by
 155 cross-validation will make a great difference to the regression coefficients and
 156 sometimes the differences can even vary by several orders of magnitude.³⁴ Thus, the
 157 regression coefficients of the same variable in different sub-models may change a great
 158 deal and they cannot be calculated or compared directly. In this study, we employ the
 159 strategy of normalization to remove this influence. Assume that after building N
 160 sub-models, a regression coefficient matrix \mathbf{B} ($\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_N]^T$) is generated. The j th

161 row vector in \mathbf{B} , denoted by \mathbf{b}_j ($1 \leq j \leq N$) records the regression coefficients of the j th
 162 sub-model. The elements in the matrix \mathbf{B} will be set to 0 if the responding variables are
 163 not included into the sub-models. The regression coefficient b_{ij} of the i th variable in the
 164 j th sub-model is normalized as follow:

$$165 \quad c_{ij} = b_{ij} / \max(\mathbf{b}_j) \quad (3)$$

166 where $\max(\mathbf{b}_j)$ stands for the maximum of the row vector \mathbf{b}_j . The normalized regression
 167 coefficient matrix composed by all c_{ij} is denoted by \mathbf{C} . The elements in \mathbf{C} range from 0
 168 to 1.

169 **2.4. The criteria and weights of variables**

170 For CARS, it is the mean of the regression coefficients of one variable that is
 171 considered as the criterion to determine its importance level. In IVSO, the normalized
 172 regression coefficient matrix \mathbf{C} is summed in columns to generate a row vector \mathbf{s} ($\mathbf{s} = [s_1,$
 173 $s_2, \dots, s_p]$), where s_i stands for the sum of the regression coefficients of the i th variable in
 174 all N sub-models. The sum s_i of the i th variable is regarded as its criterion. By this data
 175 fusion step, the noise can be cancelled out and the systematic information can be
 176 accumulated. In this way the difference between variables will become larger than that
 177 in CARS, which accelerating the iteration.

178 In each iterative round, the weight of the i th variable is defined as:

$$179 \quad w_i = s_i / \max(\mathbf{s}), \quad i = 1, 2, \dots, P \quad (4)$$

180 where $\max(\mathbf{s})$ is the maximum of the vector \mathbf{s} . The weights of the variables having been
 181 eliminated are set to zero automatically so that the weight vector \mathbf{w} is always
 182 p -dimensional. Moreover, it should be mentioned that the weight vector only work for
 183 sampling by WBMS in the next iterative round.

184 **2.5. Sequentially addition**

185 In each iterative round, we combine WBMS with another strategy called
186 sequentially addition to optimize the variable space. Firstly we use WBMS to eliminate
187 variables in a competitive way. Denote L_1 as the number of the variables which can be
188 sampled by WBMS. Then the L_1 variables are ranked based on their criteria. The
189 variable space is further shrunk by sequentially addition. The L_1 variables are
190 sequentially added step by step to establish L_1 PLS sub-models according to the rank
191 and the performance of the sub-models is estimated by cross-validation. The first
192 sub-model consists of only the first variable in the rank, and the second sub-model
193 contains the first two variables, and the i th sub-model contains the first i variables.
194 Repeat this process until the L_1 variables are all included into the last sub-model. When
195 the RMSECV value of the sub-model reaches minimum with addition one by one, the
196 corresponding variable subset in this best sub-model is chosen. The number of
197 variables in this variable subset is denoted by L_2 . The iterative process is continued with
198 L_1 getting closer to L_2 , until both L_1 and L_2 reach an equal value. One variable subset is
199 yielded in one iterative round and finally many different variable subsets are generated.
200 The RMSECV value is employed as the objective function to search for the best
201 variable subset.

202 In each iterative round, sequentially addition can select a variable subset which
203 contains informative variables. Thus, if some important variables are lost by WBMS,
204 they still can be retained in the variable subset in the previous rounds by sequentially
205 addition. When selecting the best variable subset among all iterative rounds, these lost

206 variables still have the opportunity to be included into the ultimate variable subset. In
207 this way, no loss of important variables can be assured. For the same reason, IVSO
208 possesses high stability. Overall, IVSO has the ability to eliminate variables gradually
209 and gently in a competitive way and reduce the risk of losing important variables.

210 **2.6. General description of IVSO**

211 Fig. 2 shows the scheme of IVSO algorithm. The initial value of the weight of
212 each variable is set to 1. It should be mentioned that the weights for sampling by
213 WBMS are obtained from the previous iterative round. The detailed algorithm of IVSO
214 is described as follows:

215 Step 1: Creating a binary matrix \mathbf{NM} with dimensionality $N \times P$ for sampling by
216 WBMS gives N sampling runs. In each column of \mathbf{NM} , there are Nw_i '1' and the left
217 ones are all '0'. If the Nw_i of one variable is less than 1, it will not be sampled in any row.
218 Record the number of variables which can be sampled by WBMS, namely L_1 .

219 Step 2: Build N PLS sub-models to calculate the regression coefficient matrix \mathbf{B} .
220 Each row of the matrix \mathbf{B} is normalized to generate the matrix \mathbf{C} , as Formula 2 did.

221 Step 3: Each column of the matrix \mathbf{C} is summed as the criterion of the
222 corresponding variable, denoted by the vector \mathbf{s} . Rank the L_1 variables based on their
223 criteria.

224 Step 4: Build L_1 sub-models through sequentially addition according to the rank.
225 Take the variable subset in the sub-model with the lowest RMSECV value as the
226 objective variable subset of this iterative round. Record this RMSECV value R and the
227 length of this variable subset L_2 .

228 Step 5: The vector \mathbf{s} is normalized to calculate weights as Formula 3. The weights
229 in this iterative round only work in the sampling of the next iterative round.

230 Step 6: Repeat the steps 1-5 until L_1 is equal to L_2 , then many variable subsets are
231 obtained. The variable subset with the lowest R value is chosen as the ultimate variable
232 subset of the algorithm.

233 (Insert Figure 2)

234

235 **3. Datasets and Software**

236 **3.1. Simulated dataset**

237 This dataset, called SIMUIN, is simulated as described in Reference 18. SIMUIN
238 contains 100 samples in rows and 200 variables in columns with exactly five latent
239 variables. The relative eigenvalues by principal component analysis on the autoscaled
240 data are 21.29%, 20.30%, 19.84%, 19.61%, 18.96%. The first 100 variables of
241 SIMUIN are linearly relative with \mathbf{y} but the last 100 ones are random numbers from 0 to
242 1, regarded as uninformative variables. The noises with normal distribution in the range
243 from 0 to 0.005 are added to SIMUIN .

244 **3.2. Corn moisture dataset**

245 The corn dataset is available in the website: <http://www.eigenvector.com/data/Corn/index.html>. This dataset contains 80 samples of corn measured on three
246 different NIR spectrometers. The spectrum has been recorded from 1100 - 2498 nm
247 with 700 spectral points at intervals of 2 nm. In this study, the NIR spectrum of 80 corn
248 samples measured on m5 instrument is used and the moisture value is considered as
249 property of interest \mathbf{y} .

251 **3.3. Wheat dataset**

252 This NIR dataset⁴ consists of 100 wheat samples and the protein value is
253 considered as property of interest \mathbf{y} . The spectrum has been recorded from 1100 - 2500
254 nm with 701 spectral points at intervals of 2 nm. Due to the ‘large p , small n ’ problem,^{36,}

255 ³⁷ the original spectrum is compressed into a maximum of 200 points by an appropriate
256 window size as did by Leardi.¹⁴ Setting the window size to 4, this dataset is reduced to
257 175 variables with the average of every four original variables.

258 **3.4. Software**

259 All the computations are achieved in MATLAB on an ordinary computer
260 configured to Intel Core i5 3.2 GHz CPU, 3G RAM, WIN7 Ultimate. The MATLAB
261 code for implementing IVSO is available in the supplemental materials.

262

263 **4. Results and Discussion**

264 In this study, all the datasets are split into calibration set (80% of the dataset) and
265 independent test set (20% of the dataset) based on Kennard - Stone (KS) method.³⁸ KS
266 method aims at covering the multidimensional space by maximizing the Euclidean
267 distances between each pair of the selected samples. The calibration set is used for
268 variable selection and goodness of fit, and the independent test set is used for
269 validation of the calibration model for prediction. When conducting variable selection
270 on the calibration set, cross-validation is conducted. Furthermore, in order to evaluate
271 the performance of IVSO, we compare it with another two outstanding methods based
272 on the regression coefficient, namely MC-UVE and CARS. For MC-UVE, the number
273 of Monte Carlo sampling runs is set to 1000, and 80% samples are randomly chosen for
274 modeling in each sampling run. As to CARS, the number of Monte Carlo sampling runs
275 is set to 100. For all methods, the maximum latent variable is limited to 10 and the
276 number of latent variables is determined by 10-fold cross-validation. Each dataset is

277 autoscaled to have zero mean and unit variance before modeling. Besides, the root
278 mean square error of calibration set (RMSEC) and the root mean square error of
279 prediction of test set (RMSEP) are employed to assess the performance of the three
280 methods. In addition, because of the random sampling, these methods are conducted 50
281 times to obtain statistical results and compare the three methods fairly.

282 **4.1. The influence of the number of sampling by WBMS**

283 To investigate the influence of the number of sampling runs by WBMS, we
284 discuss four cases about the performance of IVSO, as shown in Fig. 3. The number of
285 sampling runs is set to 3000, 5000, 8000 and 10000, respectively, in the three datasets.
286 For these three datasets, their RMSEP values generated by full spectrum are 0.4043,
287 0.0237 and 0.2382, respectively. All of the results of the three datasets have good
288 stability. Overall, no significant influence on the results of IVSO has been found among
289 these four cases. For the dataset of wheat protein, the median values of the four RMSEP
290 values are the same, but the results with the parameter of 8000 shows the best stability.
291 Thus, the number of sampling by WBMS is set to 8000 in this study.

292 (Insert Figure 3)

293 **4.2. Simulated dataset**

294 This dataset is simulated to assess the ability of IVSO to select appropriate
295 variable subset. The first 100 variables are linearly relative with y and regarded as
296 informative ones, but the last 100 ones are noisy. The results obtained by conducting
297 different methods within 50 times are discussed in detail.

298 Table 1 includes the results of the three methods on the three datasets. The mean
299 and 95% confidence interval are given as well. The simulated dataset is investigated
300 by comparing with MC-UVE, CARS, the full spectrum and the first 100 informative

301 variables. Compared with the full spectrum, the RMSEC value of only the first 100
302 variables drops from 0.0644 to 0.0091 and the RMSEP value drops from 0.4043 to
303 0.0135, even using a smaller number of latent variable, 6. It demonstrates the
304 importance and necessity for variable selection in multivariate calibration.

305 The statistical features of the results by different methods can be observed visually
306 in the boxplot of Fig. 4, in which Fig. 4A displays the results of the simulated dataset.
307 As it can be seen from Table 1 and Fig. 4A, IVSO shows the best performance with
308 regard to the improving the prediction ability of the model and good stability. The 95 %
309 confidence interval of RMSEC and RMSEP results for each method shows that IVSO
310 has no overlap with other methods. In addition, the selected latent variable of IVSO is 6,
311 which is much smaller than that of the full spectrum.

312 The frequency distribution of selected variables within 50 times is displayed in Fig.
313 5. For different methods the selected variables all concentrate in the first 100 variables.
314 Both IVSO and MC-UVE can select variables with high frequencies. However, the
315 selected variables by CARS are of low frequencies and even no one variable can be
316 selected by all 50 times, which reveals its instability. The fact is just consistent with its
317 large confidence interval in Table 1 and standard deviation in Fig. 4A.

318 Fig. 6A and Fig. 6B show the changing trend of the number of variables sampled
319 by IVSO and CARS respectively. The arrow indicates the point reaching the optimal
320 variable subset. As to MC-UVE, it is the full spectrum that is used to establish the
321 sub-models, so no iterative round has ever occurred. For the simulated dataset in Fig.
322 6A, the number of sampled variables decreases to 100 in the 3th and 4th iterative
323 rounds, then the curve drops much more slowly. In stark contrast, the number of
324 sampled variables of the simulated dataset in Fig. 6B varies tremendously in the front
325 section of the curve. It is in the first iterative round that the number decreases to 97,

326 which means that just the first iterative round can removes not only uninformative but
327 also informative variables. In the 11th iterative round CARS achieves its optimal
328 variables subset containing only 26 variables. It can be concluded that CARS
329 eliminates variables too quickly and thus lose informative variables, which may result
330 from the enforced elimination of variables by EDF. On the contrary, IVSO has the
331 ability to eliminate uninformative variables gradually and gently, and achieve much
332 higher stability.

333 Fig. 7 shows the RMSECV value of the variable subset chosen by sequentially
334 addition in each iterative round, which is corresponding to the sampling curve of the
335 simulated dataset in Fig. 6A. The '0' iterative round stands for the process during which
336 the weights of all variables are set to '1' as initial values but all these variables are
337 conducted by sequentially addition. From Fig. 7A, the RMSECV value of the
338 simulated dataset drops first because of the existing of some uninformative variables
339 and begins to rise again due to the missing of some informative variables. It
340 demonstrates the good ability of IVSO to eliminate uninformative variables and keep
341 informative ones.

342 (Insert Table 1)

343 (Insert Figure 4)

344 (Insert Figure 5)

345 (Insert Figure 6)

346 (Insert Figure 7)

347 **4.3. Corn moisture dataset**

348 The results obtained by repeating the three different methods 50 times are reported
349 in Table 1 and Fig. 4B. In Table 1, compared with the full spectrum, the RMSEC and
350 RMSEP values of IVSO decrease by 98.4% and 98.6%, respectively. Clearly, IVSO

351 has highly improved the prediction performance. IVSO exhibits not only the best
352 prediction ability in terms of the RMSEC and RMSEP values but also owns the best
353 stability based on confidence interval. The number of latent variable is also the
354 smallest, which means that it can generate more parsimony model.

355 The frequencies of variables selected by these methods are displayed in Fig. 8.
356 Both CARS and IVSO mainly select variables of 1908nm and 2108nm, which have
357 been discussed and proven as the key wavelengths by the literature of CARS
358 experimentally and theoretically²⁷. These two wavelengths are relative with the water
359 absorption and the combination of O-H bond.¹³ For CARS, it cannot select the key
360 wavelength of 2108nm in every iterative round. However, except for these two key
361 variables, MC-UVE selects too many other variables with high frequencies.

362 From the corn moisture dataset in Fig. 6A and Fig. 6B, we also can see that the
363 number of variables sampled by IVSO in the previous rounds drops much more
364 gradually and gently than that of CARS. In the latter rounds, though this number of
365 IVSO changes more quickly, the key variables of 1908nm and 2108nm still can be
366 retained in every iterative round due to sequentially addition. But CARS cannot do it.

367 In Fig. 7B, it reaches the optimal variable subset with the two key variables firstly
368 in the 4th iterative round. Then the optimal variable subset keeps unchanged, so the
369 RMSECV value is stable. From the RMSECV values, we can summary that the strategy
370 of sequentially addition used in every iterative round makes the result stable.

371 (Insert Figure 8)

372 **4.4. Wheat protein dataset**

373 In Table 1 and Fig. 4C, IVSO can achieve better results with the smallest number
374 of latent variable than the full spectrum. But after selecting variables, the RMSEP
375 values of both MC-UVE and CARS get much worse. Fig. 9 displays the frequencies of

376 variables selected by these methods. The variables around 1144-1296nm can be
377 selected by all methods, which is the same as the result of GA-PLS.¹⁴ IVSO can select
378 variables with quite high frequencies. As to MC-UVE and CARS, it selects many
379 variables in other regions. Moreover, the frequencies of variables selected by CARS are
380 not high and even quite a number of variables are selected by less than five times. From
381 the wheat protein dataset in Fig. 6A and Fig. 6B, we also can see that the number of
382 variables sampled by IVSO decreases much more slowly than that of CARS. The
383 RMSECV value of the variable subset in Fig. 7C goes down at first with the decrease of
384 the uninformative variables and then goes up because of increasingly deleting the
385 informative variables.

386 (Insert Figure 9)

387

388 5. Conclusion

389 This paper presents a new method for variable selection based on the regression
390 coefficient, called iteratively variable subset optimization (IVSO). Investigated by one
391 simulated dataset and two NIR datasets, IVSO is proven to be a better variable
392 selection method than another two methods, namely Monte Carlo uninformative
393 variable elimination (MC-UVE) and competitive adaptive reweighted sampling
394 (CARS). IVSO can eliminate uninformative variables gradually and gently, and
395 achieve good prediction and stability. The outstanding performance of IVSO indicates
396 that it is an efficient procedure and an alternative for variable selection.

397 Although IVSO is worked with partial least squares (PLS) to select variables in
398 this study, it also can be coupled with other modeling methods on regression or

399 pattern recognition. Our future work will focus on investigating the application of
400 IVSO in other fields, such as metabolomic and quantitative structure-activity
401 relationship (QSAR).

402 **Acknowledgments**

403 This work is financially supported by the National Nature Foundation Committee
404 of P.R. China (grants no. 21275164 and 21465016) and also supported by the
405 Fundamental Research Funds for the Central Universities of Central South University
406 (grants no. 2014zzts014). The studies meet with the approval of the university's review
407 board.

408

409

410

411

412

413

414

415

416

417

418

419

420

421 **References**

- 422 1. S. Wold, M. Sjöström and L. Eriksson, *Chemometr. Intell. Lab.*, 2001, **58**, 109-130.
- 423 2. C. M. Andersen and R. Bro, *J. Chemometr.*, 2010, **24**, 728-737.
- 424 3. D. Jouan-Rimbaud, B. Walczak, D. L. Massart, I. R. Last and K. A. Prebble, *Anal. Chim. Acta.*,
425 1995, **304**, 285-295.
- 426 4. J. H. Kalivas, *Chemometr. Intell. Lab.*, 1997, **37**, 255-259.
- 427 5. C. H. Spiegelman, M. J. McShane, M. J. Goetz, M. Motamedi, Q. L. Yue and G. L. Coté, *Anal.*
428 *Chem.*, 1998, **70**, 35-44.
- 429 6. Y. H. Yun, Y. Z. Liang, G. X. Xie, H. D. Li, D. S. Cao and Q. S. Xu, *Analyst*, 2013, **138**, 6412-6421.
- 430 7. I. Guyon and A. Elisseeff, *J. Mach. Learn. Res.*, 2003, **3**, 1157-1182.
- 431 8. V. Centner, D. L. Massart, O. E. de Noord, S. de Jong, B. M. Vandeginste and C. Sterna, *Anal.*
432 *Chem.*, 1996, **68**, 3851-3858.
- 433 9. W. Cai, Y. Li and X. Shao, *Chemometr. Intell. Lab.*, 2008, **90**, 188-194.
- 434 10. Q. J. Han, H. L. Wu, C. B. Cai, L. Xu and R. Q. Yu, *Anal. Chim. Acta.*, 2008, **612**, 121-125.
- 435 11. S. Favilla, C. Durante, M. L. Vigni and M. Cocchi, *Chemometr. Intell. Lab.*, 2013, **129**, 76-86.
- 436 12. O. M. Kvalheim, *J. Chemometr.*, 2010, **24**, 496-504.
- 437 13. J. H. Jiang, R. J. Berry, H. W. Siesler and Y. Ozaki, *Anal. Chem.*, 2002, **74**, 3555-3565.
- 438 14. R. Leardi, *J. Chemometr.*, 2000, **14**, 643-655.
- 439 15. R. Leardi and A. Lupiáñez González, *Chemometr. Intell. Lab.*, 1998, **41**, 195-207.
- 440 16. Y. H. Yun, D. S. Cao, M. L. Tan, J. Yan, D. B. Ren, Q. S. Xu, L. Yu and Y. Z. Liang, *Chemometr. Intell.*
441 *Lab.*, 2014, **130**, 76-83.
- 442 17. J. Kennedy, in *Encyclopedia of Machine Learning*, eds. C. Sammut and G. Webb, Springer US,
443 2010, ch. 630, pp. 760-766.
- 444 18. M. Goodarzi and L. dos Santos Coelho, *Anal. Chim. Acta.*, 2014, **852**, 20-27.
- 445 19. F. Allegrini and A. C. Olivieri, *Anal. Chim. Acta.*, 2011, **699**, 18-25.
- 446 20. S. Tabakhi and P. Moradi, *Pattern. Recogn.*, 2015, **48**, 2798-2811.
- 447 21. J. Xiang, X. Han, F. Duan, Y. Qiang, X. Xiong, Y. Lan and H. Chai, *Appl. Soft. Comput.*, 2015, **31**,
448 293-307.
- 449 22. J. H. Kalivas, N. Roberts and J. M. Sutter, *Anal. Chem.*, 1989, **61**, 2024-2030.
- 450 23. L. Breiman, *Mach. Learn.*, 2001, **45**, 5-32.
- 451 24. M. C. U. Araújo, T. C. B. Saldanha, R. K. H. Galvão, T. Yoneyama, H. C. Chame and V. Visani,
452 *Chemometr. Intell. Lab.*, 2001, **57**, 65-73.
- 453 25. Y. H. Yun, W. T. Wang, M. L. Tan, Y. Z. Liang, H. D. Li, D. S. Cao, H. M. Lu and Q. S. Xu, *Anal.*
454 *Chim. Acta.*, 2014, **807**, 36-43.
- 455 26. Y. H. Yun, W. T. Wang, B. C. Deng, G. B. Lai, X. B. Liu, D. B. Ren, Y. Z. Liang, W. Fan and Q. S. Xu,
456 *Anal. Chim. Acta.*, 2015, **862**, 14-23.
- 457 27. H. Li, Y. Liang, Q. Xu and D. Cao, *Anal. Chim. Acta.*, 2009, **648**, 77-84.
- 458 28. 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,
459 *Spectrochim. Acta. A.*, 2013, **111**, 31-36.
- 460 29. A. G. Frenich, D. Jouan-Rimbaud, D. L. Massart, S. Kuttatharmmakul, M. M. Galera and J. L. M.
461 Vidal, *Analyst*, 1995, **120**, 2787-2792.

- 462 30. O. M. Kvalheim, R. Arneberg, O. Bleie, T. Rajalahti, A. K. Smilde and J. A. Westerhuis, *J.*
463 *Chemometr.*, 2014, **28**, 615-622.
- 464 31. B. C. Deng, Y. H. Yun, P. Ma, C. C. Lin, D. B. Ren and Y. Z. Liang, *Analyst*, 2015, **140**, 1876-1885.
- 465 32. B. C. Deng, Y. H. Yun, Y. Z. Liang and L. Z. Yi, *Analyst*, 2014, **139**, 4836-4845.
- 466 33. H. Zhang, H. Wang, Z. Dai, M.-s. Chen and Z. Yuan, *BMC Bioinformatics*, 2012, **13**, 298.
- 467 34. B. C. Deng, Y. H. Yun, Y. Z. Liang, D. S. Cao, Q. S. Xu, L. Z. Yi and X. Huang, *Anal. Chim. Acta.*,
468 2015, **880**, 32-44.
- 469 35. Q. S. Xu, Y. Z. Liang and H. L. Shen, *J. Chemometr.*, 2001, **15**, 135-148.
- 470 36. E. Candes and T. Tao, *Ann. Stat.*, 2007, 2313-2351.
- 471 37. H. Zou and T. Hastie, *J. Roy. Stat. Soc. B.*, 2005, **67**, 301-320.
- 472 38. R. W. Kennard and L. A. Stone, *Technometrics*, 1969, **11**, 137-148.

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489 **Table 1**

490 Results of different methods on the three datasets.

Methods	nVar ^a	nLVs ^b	RMSEC		RMSEP	
			Min-Max	Average \pm CI ^c	Min-Max	Average \pm CI
Simulated						
PLS ^d	200	10		0.0644		0.4043
PLS ^e	100	6		0.0091		0.0135
MC-UVE	78.2 \pm 6.3	6	0.0087-0.0093	0.0089 \pm 8.3 \times 10 ⁻⁵	0.013-0.0135	0.0132 \pm 5.8 \times 10 ⁻⁵
CARS	27.6 \pm 4.8	6.3 \pm 0.7	0.0080-0.0124	0.0100 \pm 3.1 \times 10 ⁻⁴	0.0118-0.0197	0.0155 \pm 5.3 \times 10 ⁻⁴
IVSO	68.1 \pm 2.1	6	0.0090-0.0093	0.0091 \pm 1.5 \times 10 ⁻⁵	0.012-0.0137	0.0125 \pm 8.7 \times 10 ⁻⁵
Corn moisture						
PLS	701	10		0.017		0.0237
MC-UVE	70.4 \pm 2.6	10	2.4 \times 10 ⁻³ -3.0 \times 10 ⁻³	2.7 \times 10 ⁻³ \pm 4.1 \times 10 ⁻⁵	2.8 \times 10 ⁻³ -3.7 \times 10 ⁻³	3.2 \times 10 ⁻³ \pm 4.0 \times 10 ⁻⁵
CARS	3.4 \pm 2.7	3.1 \pm 1.9	2.4 \times 10 ⁻⁴ -2.7 \times 10 ⁻³	4.6 \times 10 ⁻⁴ \pm 1.8 \times 10 ⁻⁴	3.4 \times 10 ⁻⁴ -4.5 \times 10 ⁻³	6.4 \times 10 ⁻⁴ \pm 2.8 \times 10 ⁻⁴
IVSO	2.3 \pm 0.8	2.3 \pm 0.8	2.6 \times 10 ⁻⁴ -2.8 \times 10 ⁻⁴	2.8 \times 10 ⁻⁴ \pm 1.1 \times 10 ⁻⁶	3.3 \times 10 ⁻⁴ -3.6 \times 10 ⁻⁴	3.4 \times 10 ⁻⁴ \pm 1.4 \times 10 ⁻⁶
Wheat protein						
PLS	175	10		0.3923		0.2382
MC-UVE	10.6 \pm 1.3	9.9 \pm 0.3	0.3370-0.3657	0.3475 \pm 0.0012	0.2466-0.2791	0.2532 \pm 0.0027
CARS	9.8 \pm 2.8	8.2 \pm 1.3	0.2501-0.3427	0.2969 \pm 0.0054	0.1818-0.3535	0.2432 \pm 0.0111
IVSO	14.8 \pm 3.0	7.5 \pm 1.0	0.2415-0.2695	0.2641 \pm 0.0030	0.2313-0.2363	0.2339 \pm 0.0009

491 ^a The number of selected variables492 ^b The number of selected latent variables of PLS493 ^c 95% confidence interval (CI)494 ^d Results using the full spectrum with 200 variables by PLS495 ^e Results using only the first 100 informative variables by PLS

496

497

498

499 **Figure Captions**

500 Fig.1. The process of generating binary matrix.

501

502 Fig.2. The scheme of iteratively variable subset optimization (IVSO) algorithm.

503

504 Fig.3. The boxplot for each dataset with the number of sampling runs by WBMS set
505 to 3000, 5000, 8000 and 10000, respectively. (A) simulated dataset; (B) corn moisture
506 dataset; (C) wheat protein dataset. On each box, the central mark is the median, the
507 edges of the box are the 25th and 75th percentile, the whiskers extend to the most
508 extreme data points are the maximum and minimum, the “+” plotted individually
509 represents outliers.

510

511 Fig.4. The boxplot of 50 RMSEP values for the three methods. (A) simulated dataset;
512 (B) corn moisture dataset; (C) wheat protein dataset. On each box, the central mark is
513 the median, the edges of the box are the 25th and 75th percentile, the whiskers extend
514 to the most extreme data points are the maximum and minimum, and the “+” plotted
515 individually represents outliers.

516

517 Fig.5. The frequencies of variables selected by different methods within 50 times on
518 the simulated dataset. (A) MC-UVE; (B) CARS; (C) IVSO.

519

520 Fig.6. The changing trend of the number of sampled variables by IVSO (A) and
521 CARS (B).

522

523 Fig.7. The root mean squares error of cross-validation (RMSECV) of the variable
524 subset chosen by sequentially addition in each iterative round. (A) simulated dataset;
525 (B) corn moisture dataset; (C) wheat protein dataset.

526

527 Fig.8. The frequencies of variables selected by different methods within 50 times on
528 the corn moisture dataset. (A) MC-UVE; (B) CARS; (C) IVSO.

529

530 Fig.9. The frequencies of variables selected by different methods within 50 times on
531 the wheat protein dataset. (A) MC-UVE; (B) CARS; (C) IVSO.

532

533

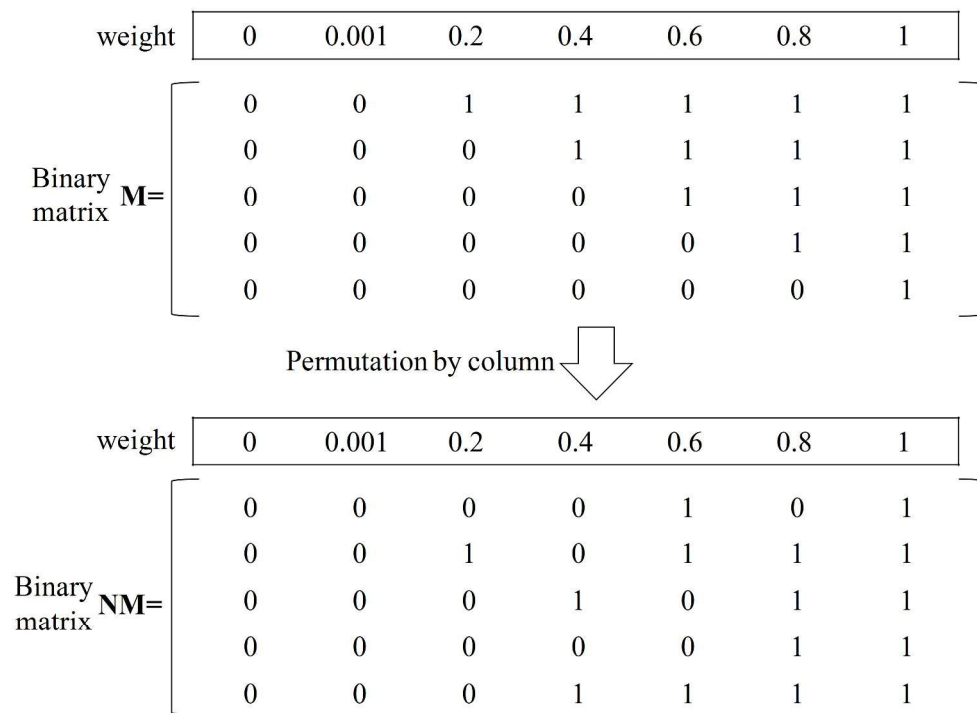


Fig.1. The process of generating binary matrix.
295x216mm (300 x 300 DPI)

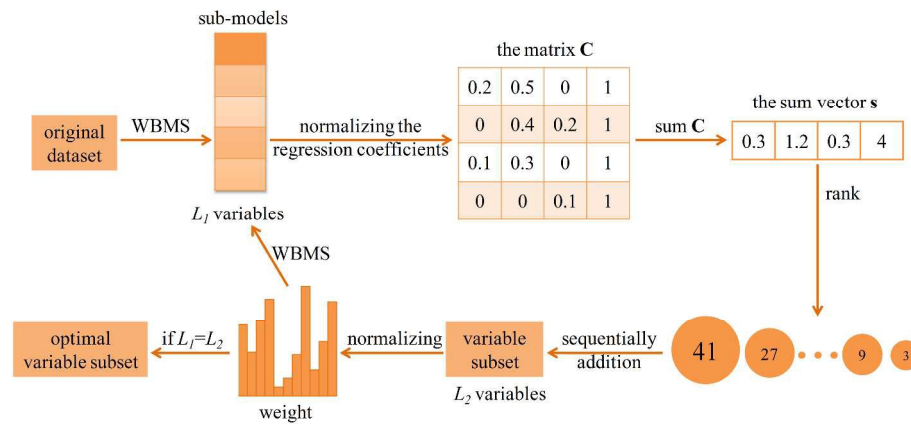


Fig.3. The boxplot for each dataset with the number of sampling runs by WBMS se 438x187mm (300 x 300 DPI)

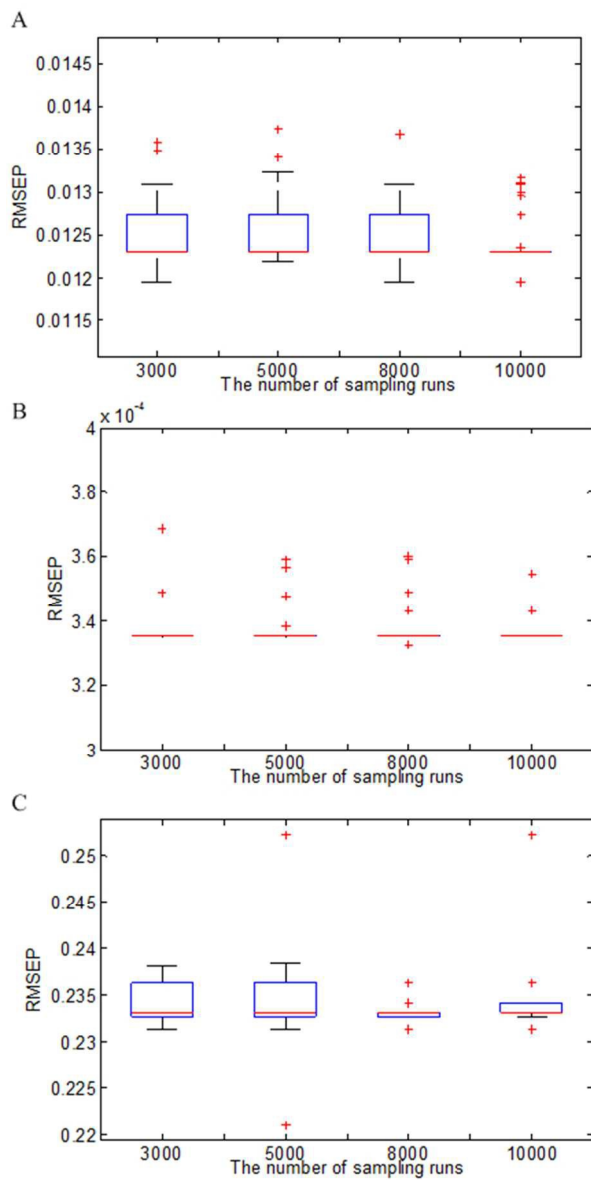


Fig.3. The boxplot for each dataset with the number of sampling runs by WBMS set to 3000, 5000, 8000 and 10000, respectively. (A) simulated dataset; (B) corn moisture dataset; (C) wheat protein dataset. On each box, the central mark is the median, the edges of the box are the 25th and 75th percentile, the whiskers extend to the most extreme data points are the maximum and minimum, the "+" plotted individually represents outliers.

239x445mm (50 x 50 DPI)

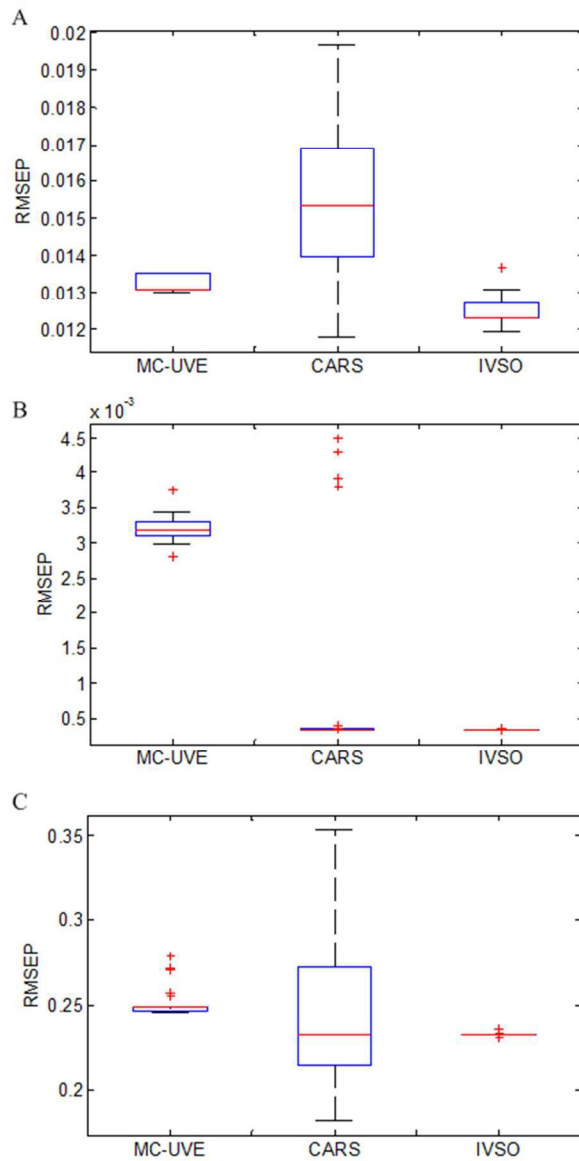


Fig.4. The boxplot of 50 RMSEP values for the three methods. (A) simulated dataset; (B) corn moisture dataset; (C) wheat protein dataset. On each box, the central mark is the median, the edges of the box are the 25th and 75th percentile, the whiskers extend to the most extreme data points are the maximum and minimum, and the "+" plotted individually represents outliers.

236x443mm (50 x 50 DPI)

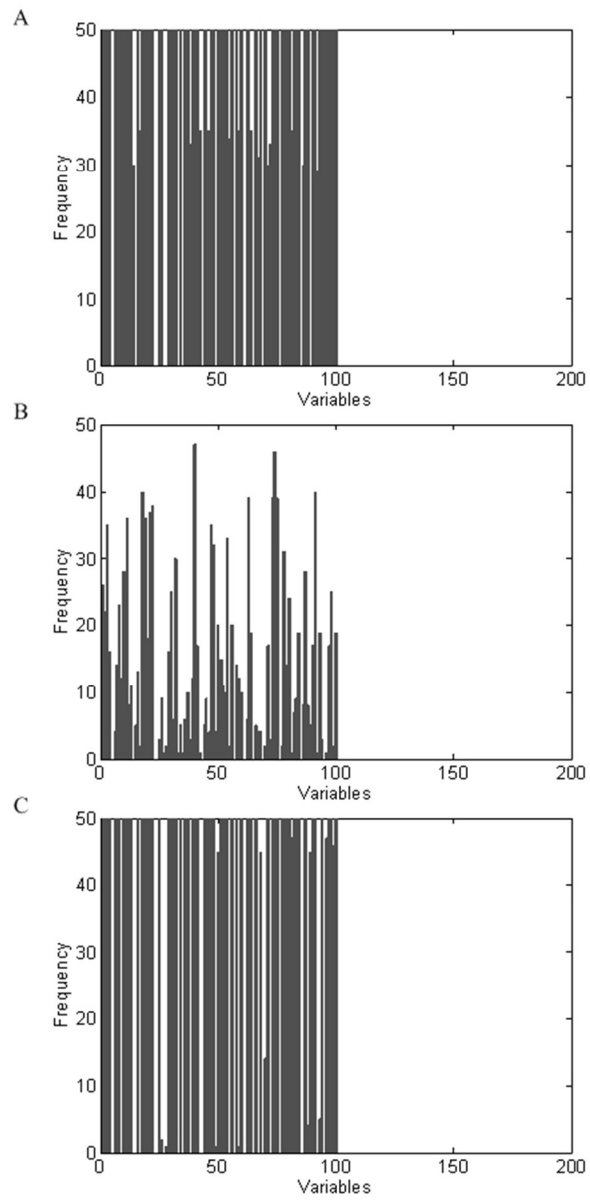


Fig.5. The frequencies of variables selected by different methods within 50 times on the simulated dataset.
(A) MC-UVE; (B) CARS; (C) IVSO.
235x457mm (50 x 50 DPI)

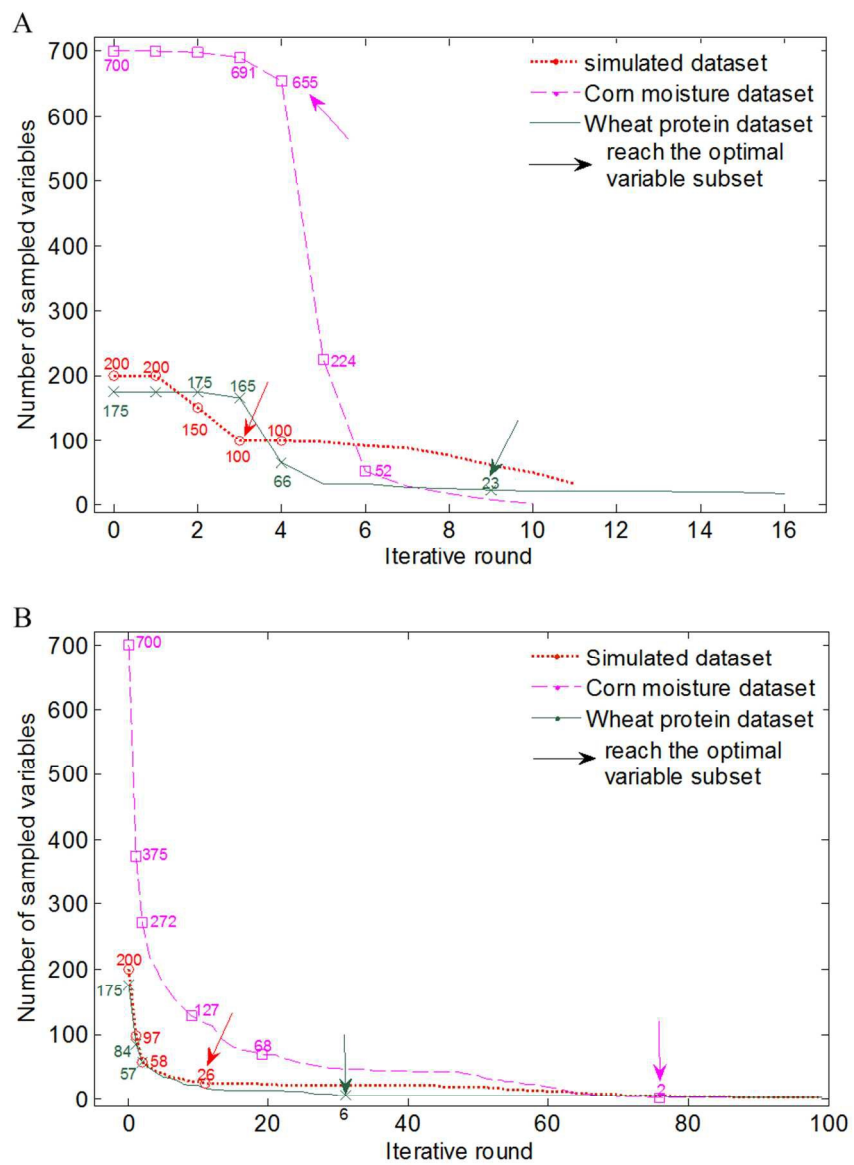


Fig.6. The changing trend of the number of sampled variables by IVSO (A) and CARS (B).
271x343mm (100 x 100 DPI)

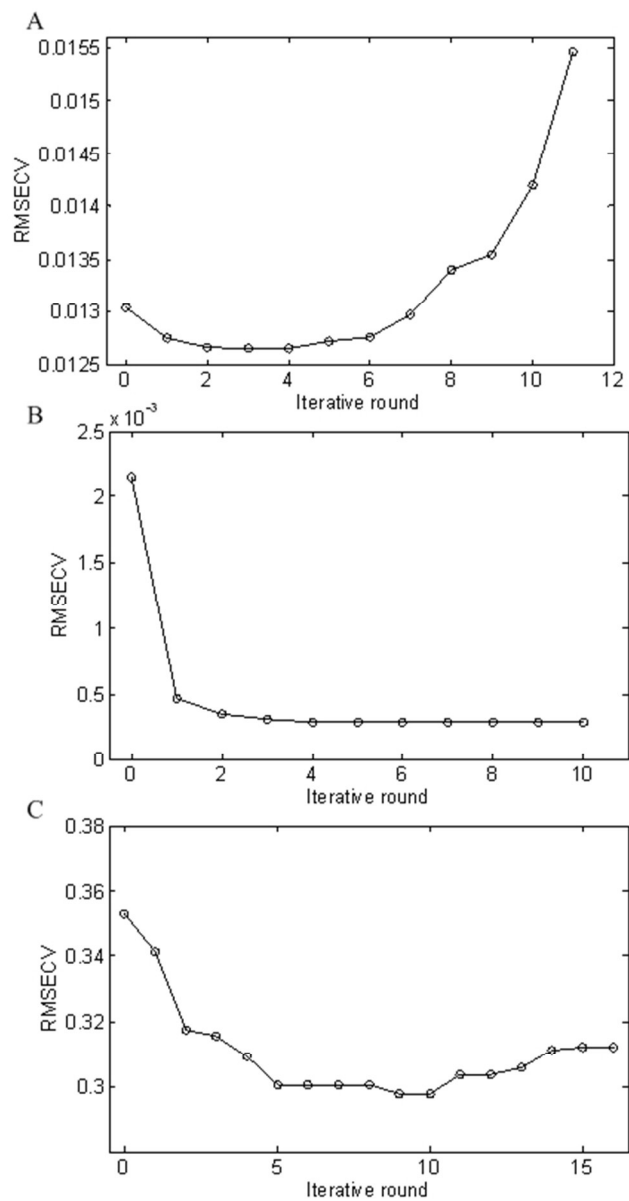


Fig.7. The root mean squares error of cross-validation (RMSECV) of the variable subset chosen by sequentially addition in each iterative round. (A) simulated dataset; (B) corn moisture dataset; (C) wheat protein dataset.
234x411mm (50 x 50 DPI)

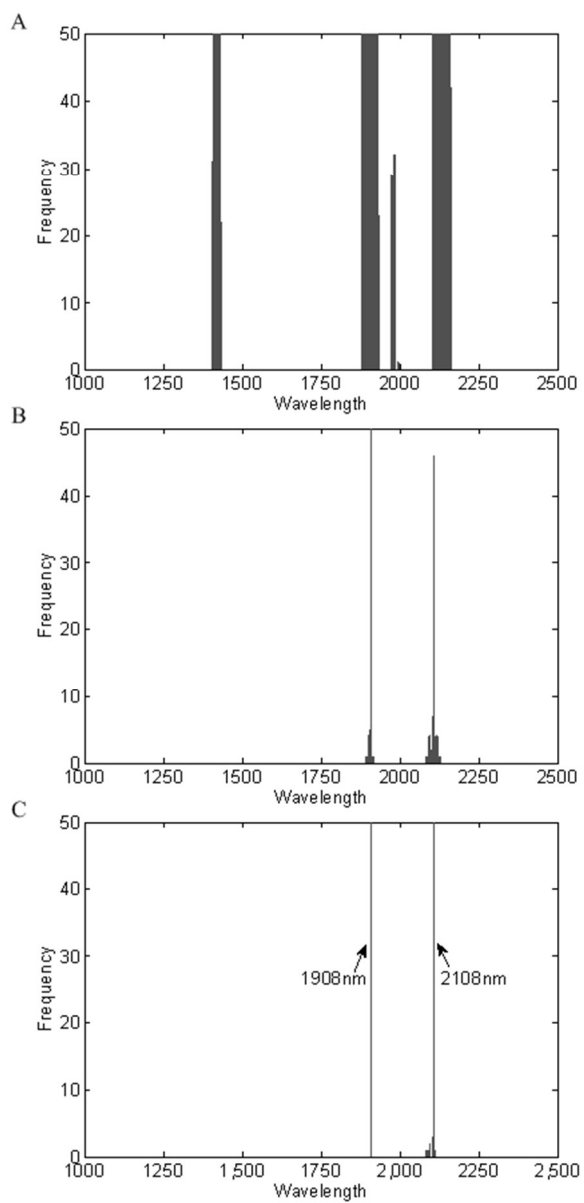


Fig.8. The frequencies of variables selected by different methods within 50 times on the corn moisture dataset. (A) MC-UVE; (B) CARS; (C) IVSO.
231x457mm (50 x 50 DPI)

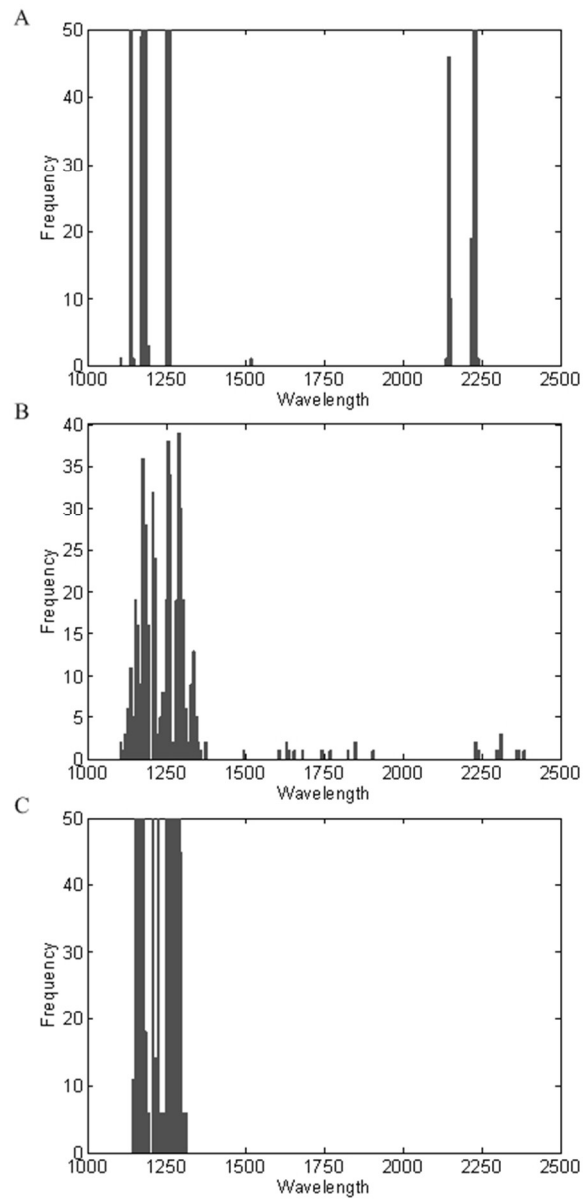


Fig.9. The frequencies of variables selected by different methods within 50 times on the wheat protein dataset. (A) MC-UVE; (B) CARS; (C) IVSO.

231x457mm (50 x 50 DPI)