# Analyst

Accepted Manuscript

1    **The Model Adaptive Space Shrinkage (MASS) Approach: A**

2    **New Method for Simultaneous Variable Selection and**

3    **Outlier Detection Based on Model Population Analysis**

4    Ming Wen[a,b], Bai-Chuan Deng[c], Dong-Sheng Cao[a*], Yong-Huan Yun[b], Rui-Han Yang[b],

5    Hong-Mei Lu[b†], Yi-Zeng Liang[b]

6    [a] *School of Pharmaceutical Sciences, Central South University, Changsha 410013, PR China*

7    [b] *College of Chemistry and Chemical Engineering, Central South University, Changsha 410083,*

8    *PR China*

9    [c] *College of Animal Science, South China Agricultural University, Guangzhou 510642, P.R. China*

10    **Abstract**

11    Variable selection and outlier detection are important processes in chemical modeling.

12    Usually, they affect each other. Their performing orders also strongly affect the

13    modeling result. Currently, many studies perform them separately and in different

14    orders. In this study, we discussed the interaction between outliers and variables, and

15    compared the modeling procedures performed in different variable selection and

16    outlier detection orders. Because the order of outlier detection and variable selection

17    can affect the interpretation of the model, it is hard to decide which order is better

18    when the predictability (prediction error) of different orders is relatively close. To

19    handle this problem, a simultaneous variable selection and outlier detection approach

---

\*   Corresponding author. *E-mail address:* oriental-cds@163.com (Dong-Sheng Cao)

†   Corresponding author. *E-mail address:* hongmeilu@csu.edu.cn (Hong-Mei Lu)

1

20  called Model Adaptive Space Shrinkage (MASS) was developed. This proposed

21  approach is based on model population analysis (MPA). Through weighted binary

22  matrix sampling (WBMS) from model space, a large number of partial least square

23  (PLS) regression models were built, and the elite part of models were selected for

24  statistically reassigning the weight of each variable and sample. Then, the whole

25  process repeated until the weights of variables and samples were converged. Finally,

26  MASS adaptively found a high performance model which consisted of the optimized

27  variable subset and sample subset. The combination of these two subsets could be

28  considered as the cleaned dataset used for the chemical modeling. In the proposed

29  approach, the problem of the order of variable selection and outlier detection is

30  avoided. One near infrared spectroscopy (NIR) dataset and one quantitative

31  structure-activity relationship (QSAR) dataset were used to test this approach. The

32  result demonstrated that MASS is a useful method in data cleaning before building a

33  predictive model.

34  **Key words:** outlier detection, variable selection, model population analysis,

35  shrinkage, model space

## 1. Introduction

37  With the development of modern analytical instruments, numerous data which

38  contain a large number of variables and samples can be obtained through

39  high-throughput experimental method. Multivariate regression techniques such as

40  multivariate linear regression (MLR) [1], partial least square regression (PLS) [2], support

41  vector regression (SVR) [3] and random forest (RF) [4] are useful tools to analyze those

2

42　data and have been applied in different fields. However, the applications of a built

43　model are seriously affected by the quality of the model. To build a robust and reliable

44　model, variable selection and outlier detection method have been wildly used to

45　improve the performance of regression models.

46　　　In general, variable selection methods can be divided into three categories. One

47　is classical methods such as forward selection method [5] and backward elimination

48　method [6], without considering the combination effect of variables [7]. One is artificial

49　intelligence-based method like genetic algorithm (GA) method [8], artificial neural

50　network (ANN) method [9] and particle swarm optimization (PSO) method [10] which

51　have been applied to search the optimal subset of variables. One is statistical method

52　such as uninformative variable elimination (UVE) [11], variable iterative space

53　shrinkage approach (VISSA) [8] and iteratively retaining informative variables (IRIV)

54　[12]. They select variables by statistically evaluating some values of a model.

55　　　Detecting outlier is troublesome especially when several outliers coexist.

56　Diagnostics and robust regression are two methods to deal with outliers [13]. In the

57　diagnostic method, outliers are identified first, and the rest of samples are used to

58　build model. Monte-Carlo (MC) method is a typical diagnostic method. It uses

59　Monte-Carlo sampling method to build a large number of models. Each sample is

60　predicted by all models. The standard deviation and mean value of predictive error are

61　calculated. The sample with large standard deviation or large mean value could be

62　considered as outliers. In the robust regression method, a regression model is

63　constructed to fit the majority of the data. Outliers are detected by examining the

3

64 residuals which are predicted by the built model. Some representative methods

65 include least median of squares (LMS) [14], robust principal component regression

66 (RPCR) [15] and robust partial least squares (RPLS) [16] and so on.

67     Before building a model, variable selection and outlier detection must be

68 carefully considered, especially their interactions (i.e., their performing orders). It is

69 worth to note that outlier detection and variable selection can influence each other [17].

70 Different results may be obtained by performing these two tasks in the opposite order.

71 Thus, the order of variable selection and outlier detection will intensively influence

72 the application of a model. It is therefore necessary to consider variable selection and

73 outlier detection simultaneously. Jennifer Hoeting proposes a method for

74 simultaneous variable selection and outlier identification in linear regression, which is

75 an early research on this aspect. The approach is based on posterior model

76 probabilities. A Markov chain Monte Carlo approach is used to approximate the

77 Bayesian model average over the space of all possible variables and outliers under

78 consideration. For more detail information see reference [18]. Later some GA-based

79 methods [17, 19, 20] are proposed for this task and have been applied in different fields. J.

80 Tolvi et al. uses an ordinary genetic algorithm for outlier detection and variable

81 selection in linear regression [17]. Patrick Wiegand combines a robust outlier

82 determination method with a genetic algorithm for variable selection [19]. Rachel Cavill

83 et al. develops a genetic algorithm approach which simultaneously selects sub-sets of

84 samples and spectral regions (variables) in metabonomics data. Their results indicate

85 that simultaneous sample and variable selection method improved model performance

4

86　by over 9% compared with those separated method [20]. Rajiv S. Menjoge gives a

87　diagnostic method for simultaneous feature selection and outlier identification in

88　linear regression [21]. The method performs by adding a dummy variable set to the data

89　matrix and running backward selection on the augmented matrix. The sequences of

90　feature-outlier combinations are identified. Another method proposed by Sung-Soo

91　Kim et al [22] consists of two procedures, first identifying the potential outliers

92　(mean-shift outlier model), then exhaustively searching the possible subset

93　regressions for the mean-shift outlier model. A recent method is Monte-Carlo Outlier

94　and Variable Screening approach (MCOVS) [23]. MCOVS builds a series of

95　sub-regression models and simultaneously evaluates the importance of variables and

96　location of outliers statistically.

97　　Model Population Analysis (MPA) [24], proposed by Li et al., is a general

98　framework for designing new types of chemometrics and bioinformatics algorithms[24].

99　In MPA, firstly, randomly produce N sub-training datasets using sampling methods

100　from the original dataset. Secondly, establish a sub-regression model on each

101　sub-training dataset. Finally, statistically analyze interesting outputs of all established

102　N sub-regression models. Many methods such as MCUVE, VISSA, and IRVR are

103　developed based on MPA.

104　　Here, we proposed a strategy based on MPA called Model Adaptive Space

105　Shrinkage (MASS). It was applied to select variables and remove outliers

106　simultaneously. MASS aims to find a high performance model based on a clean

107　dataset in the model space through a weighted iteration strategy. The variable and

5

108     sample subsets are simultaneously obtained. In addition, MASS considers the outlier

109     masking effect and variable combination effect through its random sample procedures.

110     In this study, MASS coupled with PLS was tested on different data. Comparison with

111     other existing popular methods or method combination showed that MASS is a useful

112     method to select variables and outliers simultaneously. It should be noted that MASS

113     can also be coupled to other modeling methods such as artificial neural network

114     (ANN), support vector regression (SVR).

## 115    2. Theory and method

### 116    2.1   Definition of Model Space

117     After obtaining a data with $N$ samples and $P$ variables, a model space is defined

118     as a set of models which are constructed by all possible combinations of samples and

119     variables. Fig. 1 is the sketch of model space. The combination of a variable subset

120     and a sample subset forms a sub-training dataset, and the sub-training dataset is used

121     to building a regression model. The built model is a member of model space. In Fig. 1,

122     #1 and #2 are two models (members) in the model space. The number of all the

123     possible combinations for variables is $2^P$-1 (variable space) and for samples is

124     $2^N$-1(sample space). The model space is the combination of variable space and sample

125     space. It has $(2^P$-1$) \times (2^N$-1$)$ models (combinations).

126

127                            **(Insert Figure 1)**

128

6

129    **2.2    The interaction between variables and outliers**

130        Outliers depend on the variables used for characterization [23]. A sample can be

131    seen as an outlier when its location represented by variables is far away from the bulk

132    of samples. As is shown in Fig. 2a, all samples can be well fitted only using one

133    variable x1. But in Fig. 2b, when added a variable x2, sample 1 turns into an outlier

134    since its location is far away from other samples. In addition, in the dataset with

135    outliers, more variables are needed to reduce the influence of outliers. As is shown in

136    Fig. 2a, with the outlier (sample 1) in dataset, variable x2 is needed to build a model

137    (the red dotted line) to reduce the impact of this outlier. This explicitly indicates that,

138    on the one hand, different variables can lead to different outliers in the sample set; on

139    the other hand, different samples need different variables to build the best model.

140    Building a high performance model not only needs to consider the effects from

141    variable selection and outlier detection separately but also needs to consider their

142    interactions.

143

144                                **(Insert Figure 2)**

145

146    **2.3    BMS and WBMS**

147        Binary matrix sampling (BMS) is a new strategy for random sampling which is

148    proposed by Yun and Deng et.al [8, 12, 25, 26]. It can ensure that all the variables have the

149    same overall frequency of sampling in the sub-regression models. A final sampling

150    matrix with a special variable frequency is consisted of a number of sub-binary

7

151    matrices which have different variable frequency. Weighted binary matrix sampling

152    method (WBMS) [8, 25] is a modified BMS which ensure important variables and

153    samples to have high selected probability in each iteration. In the BMS strategy

154    sample sampling ratio and variable sampling ratio should be manually set according

155    to the real problem.

156    **2.4    Model Adaptive Space Shrinkage (MASS)**

157        By combining MPA and WBMS, a novel method called Model Adaptive Space

158    Shrinkage (MASS) was proposed to select variables and detect outliers

159    simultaneously. The flowchart of MASS is depicted in Fig. 3.

160

161                                **(Insert Figure 3)**

162

163        Firstly, through BMS, a number of sub-training datasets was sampled from the

164    original dataset. That is to say, the samples with specific sampling ratio (e.g., 0.95)

165    and the variables with specific sampling ratio (e.g., 0.5) were randomly selected to

166    construct one sub-training dataset from the original dataset. Initially (i.e., in the first

167    iteration), the frequency of each variable or each sample appearing in these models is

168    somewhat equal according to their sampling ratio. For example, for each

169    sub-regression model, the sub-training dataset is consisted of 95% samples and 50%

170    variables. Thus, these sub-training datasets were used to build sub-regression models

171    which are evenly distributed in model space. Then, these models were sorted by the

172    coefficient of determination of cross-validation $Q_{CV}^2$ (eq. 1) [27, 28].

8

173

$$Q_{CV}^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y}_i)^2} \tag{1}$$

175

176    where $n$ is the number of samples in the model, $\hat{y}_i$ is the prediction property value of

177    the $i$th sample, and $\bar{y}_i$ is the mean property value of sub-training dataset. The models

178    with large $Q_{CV}^2$ were extracted. Then, the frequency of each variable and each sample

179    in the selected models were counted. The weight ($\omega$) of variable $i$ and sample $j$ were

180    obtained by eq.2 and eq.3, respectively.

181

$$\omega_i = \frac{\rho_i}{K_{best}} \tag{2}$$

183

$$\omega_j = \frac{\rho_j}{K_{best}} \tag{3}$$

185

186    Where $\rho_i$ and $\rho_j$ are the frequency of variable $i$ and sample $j$ in the selected

187    sub-regression models respectively, $K_{best}$ is the number of extracted models, and $\omega$

188    is a number between 0 and 1, which represents the ratio of a sub-regression model that

189    contain variable $i$ or sample $j$ in the next iteration. In other words, large $\omega_i$ and $\omega_j$

190    indicate that the variable $i$ and sample $j$ are more important and have more chance to

191    appear in sub-regression models. So far, the first iteration finished.

192        In the next step, WBMS was used to build a number of new sub-regression

193    models by using the weight of variables ($\omega_i$) and sample ($\omega_j$) obtained from last

194    iteration. Unlike the models evenly distributed in model space in the first iteration, the

195 models gradually focus on the high performance model in the next iteration, and the

196 model space is also gradually shrinkage. The procedure for obtaining new weights of

197 variables and samples was repeated until the weights of all variables and samples

198 were constant (either 1 or 0). Thus the best model was obtained; the variables and

199 samples that constructed the best model are simultaneously selected.

200 MASS aims to find a high performance model in the model space through a

201 continuous model space shrink procedure. In the beginning, all variables and samples

202 have the same weight. In each WBMS step, the sampling method focuses on the

203 variables and samples with larger weight until the weight is up to 1. Thus, the extent

204 of the best model space shrinks continuously until we find the best model. The

205 MATLAB codes for implementing MASS are freely available at the Supporting

206 Information.

## 207 **3. Datasets**

208 To illustrate the performance of our proposed method, two online available

209 datasets were used to evaluate the MASS approach.

### 210 **3.1 Wheat kernel dataset**

211 This dataset represents 43 different varieties or variety mixtures from two

212 different locations, and consists of 415 samples and 100 variables. Each sample was

213 analyzed at the range of 850-1050 nm, and 100 wavelengths were recorded as

214 variables. This data is freely available at http://www.models.life.ku.dk/wheat_kernels.

### 215 **3.2 ACE dataset**

10

216    This is a commonly used real QSAR dataset for testing the proposed approach.

217 This dataset consists of 114 angiotensin converting enzyme (ACE) inhibitors

218 originally taken from the work of Depriest et al and 56 descriptors. Activities are

219 spread over a wide range, with each inhibitor pIC50 values ranging from 2.1 to 9.9 [29].

## 4.  Results and discussion

220

### 4.1 The comparison of Wheat kernel and ACE dataset on different methods.

221

222    In this study, for comparison of different approaches, Monte-Carlo sampling

223 (MCS) method and variable iterative space shrinkage approach (VISSA) were used to

224 detect the outliers and select compact subset of variables, respectively. MCS method[13]

225 is an outlier detection method based on MPA. It inherently provides a feasible way to

226 detect different kinds of outliers by establishment of many cross-predictive models.

227 MCS has been demonstrated as a practical outlier detection method by a series of

228 works [30-32]. VISSA, proposed by our group, is a new variable selection method based

229 on MPA. Unlike most of the existing optimization approaches for variable selection,

230 VISSA statistically evaluates the performance of each model and makes full use of the

231 information obtained in each model to iteratively find the best subset of variable. Its

232 acceptability has been proved by comparing with other popular methods [8].

233 Furthermore, the combination of MCS and VISSA were employed to improve the

234 prediction power of the model. Two strategies were considered: removing outliers

235 with MCS followed by variable selection with VISSA (MCS + VISSA) and selecting

236 variables with VISSA followed by outlier detection with MCS (VISSA + MCS).

237 Finally, MASS was compared with PLS, VISSA, MCS, VISSA+MCS, MCS+VISSA.

11

238    All these methods were executed 20 times and used the same parameters to build

239    models: the optimal number of PLS component was obtained by five-fold cross

240    validation and was used for building models. The sampling number used in VISSA

241    and MASS was 2000, and the ratio of selected best sub-regression models was 0.05

242    (that is 100 models). The initial weight of variables in VISSA is 0.5. In MASS, the

243    initial weight of variables is 0.5 and the initial weight of samples is 0.95. In addition,

244    all data were pretreated by mean-center method before modeling. The coefficient of

245    determination of calibration set ($R^2$) and coefficient of determination of cross

246    validation ($Q_{CV}^2$) were used to assess model performance. The number of selected

247    samples and variables was recorded as Sam and Var. The number of optimal latent

248    variables (optPC) was also recorded.

249    The results of wheat kernel dataset and ACE dataset performed by PLS, VISSA,

250    MCS, VISSA +MCS, MCS+VISSA and MASS were listed in Table 1 and Table 2,

251    respectively. As shown in Table 1 and 2, PLS has the worst prediction performance

252    among all these approaches. It gives $R^2$ value of 0.880 and $Q_{CV}^2$ value of 0.869 for

253    wheat kernel dataset and gives $R^2$ value of 0.745 and $Q_{CV}^2$ value of 0.623 for ACE

254    dataset. MCS (with $R^2$ value of 0.899 and $Q_{CV}^2$ value of 0.889 for wheat kernel

255    dataset and $R^2$ value of 0.819 and $Q_{CV}^2$ value of 0.729 for ACE dataset) and VISSA

256    (with $R^2$ value of 0.894 and $Q_{CV}^2$ value of 0.886 and $R^2$ value of 0.775 and $Q_{CV}^2$ value

257    of 0.694 for ACE dataset) yield better prediction accuracy than original PLS model,

258    which indicates that PLS is strongly sensitive to outliers and uninformative variables.

259    Furthermore, the two combination approaches, VISSA+MCS and MCS+VISSA,

12

260   obtained similar prediction accuracy. The results are better than those obtained from

261   single MCS and single VISSA approach. This indicates that variable selection and

262   outlier detection method are two interactively promoted methods and are

263   indispensable in data modeling process. As seen in Table 1 and 2, MASS achieves the

264   best prediction accuracy. It gives $R^2$ value of 0.921 and $Q_{CV}^2$ value of 0.913 for wheat

265   kernel dataset and gives $R^2$ value of 0.865 and $Q_{CV}^2$ value of 0.823 for ACE dataset.

266   Compared with PLS which building the model with all the samples and variables, the

267   $R^2$ and $Q_{CV}^2$ of MASS increased 4.51% and 5.18% for wheat kernel dataset and 16.1%

268   and 32.1% for ACE dataset (P value < 0.05 MASS versus PLS), respectively.

269   Compared with other methods, the $R^2$ and $Q_{CV}^2$ of MASS for both datasets are also

270   increased considerably (P value < 1×0.05 MASS versus MCS, P value < 1×10e-3

271   MASS versus MCS+VISSA, P value < 1×10e-3 MASS versus MCS+VISSA, P value

272   < 1×10e-5 MASS versus VISSA+MCS).

273

274                                 **(Insert Table 1)**

275        **(Insert Table 2)**

276        The accuracies of different orders of variable selection and outlier detection were

277   similar, but the outliers detected and variables selected by different orders varied

278   dramatically. Fig. 4 is the outlier detection plot of wheat kernel dataset. It was

279   detected by MCS method (MCS+VISSA), two blue dash lines separate the picture

280   into 4 areas, the samples located in the lower left are normal samples, the samples

281   located in other areas are outliers [13]. The locations of the dash line are determined by

13

282  3 times of the standard deviation of the mean error Mean and mean error STD [33]. In

283  addition, MCS is a very robust outlier detection method and the outlier detection plots

284  are near the same in 20 times execution. Fig. 5 is the frequency of a wheat kernel

285  sample located in the outlier area in the VISSA+MCS order in 20 times. As is shown

286  in Fig. 4, the samples enclosed by red ellipse (sample number 38, 58, 157, 158 and

287  404) are located in the lower left area. These samples are normal samples. However,

288  as is shown in Fig. 5, they turned into outliers after variables selection. As is shown in

289  Fig. 4, the samples enclosed by green ellipses (sample number 18, 25, 104 and 408)

290  are located in the lower right. These samples are outliers. However, as is shown in Fig.

291  5, they become normal samples after variables selection. Similarly, for ACE dataset,

292  Fig. 6 is the outlier detection plot detected by MCS method (MCS+VISSA).   Fig. 7

293  is the frequency of outlier detected in the VISSA+MCS order in 20 times. Sample

294  number 18, 48, 63, 64 73 and 81 (enclosed by red ellipse) are normal samples in

295  MCS+VISSA order whereas they became outliers in VISSA+MCS order. Sample

296  number 12, 13, 15, 22, 26 and 52 are outliers in MCS+VISSA order whereas they

297  turned into normal samples in VISSA+MCS order. This indicates that different

298  variables can lead to different outliers in sample set. These two different process

299  orders are acceptable if just considering the results of built models. If considering the

300  interpretation of built model, these is a puzzle to decide the final variable selection

301  and outlier detection order. Thus, when dealing with datasets with redundant variables

302  and outliers, it is important to select variables and detect outliers simultaneously.

303

14

304                              **(Insert Figure 4)**

305                              **(Insert Figure 5)**

306                              **(Insert Figure 6)**

307                              **(Insert Figure 7)**

308                              **(Insert Figure 8)**

309                              **(Insert Figure 9)**

310   **4.2 The visualization of the interaction between variables and outliers**

311      Fig. 8 and Fig. 9 can fully explain the interaction between variables and samples

312   in MASS iteration process. Fig. 8 and Fig. 9 are the plots of sample and variable

313   weight against MASS iteration. The weight reveals the trend of sample and variable in

314   the iteration. The weight is the probability of a sample or a variable to be selected to

315   build a model. In other words, large weight indicates that the variable and sample are

316   more important and have more chance to appear in sub-regression models. As shown

317   in Fig. 8 and Fig. 9, each line represents the weight variation of a sample or a variable.

318   There are three different weight variation types: 1), the lines which go down all the

319   time till the weights reach to 0. This kind of variation indicates that these variables (or

320   samples) are uninformative variables (or outliers) and there is no strong interaction

321   between these variables and outliers, these outliers and variables can be easily

322   detected and removed; 2), the lines which go up all the time till the weights reach to 1.

323   These variables (or samples) are informative variable (or normal samples) and should

324   be selected to build model; 3), the lines which go up at first, then go down; or go

325   down at first, then go up till the weights reach to 1 or 0. This kind of variation

15

326    indicates there is strong interaction between these variables and outliers. When the

327    line goes down, it means that the variable (or sample) may be an uninformative

328    variable (or outliers) with current samples (or variables). After several iterations,

329    some outliers (variables) are removed, then the line goes up and the variable (or

330    sample) became important with current samples (or variables).

331        With regard to wheat kernel dataset and the ACE dataset, MASS was converged

332    after 31 and 30 iterations, respectively. The whole MASS iteration process can be

333    separate into 3 parts: 1), the early iteration period (1-17 iteration for wheat kernel and

334    1-13 iterations for ACE dataset). In this period, the weights of most samples gradually

335    reached to 1 except some weights of samples decreased step by step. At the same time,

336    except the weights of some variables rose up to 1 and the weights of few variables

337    went down to 0, the weights of most variables fluctuated dramatically up and down.

338    However, the weights of samples and variables which went down to 0 at this period

339    are without fluctuation or with small fluctuation. As is shown in Fig. 8(a), for wheat

340    kernel dataset, sample number 199, 3, 25, 363, 158 and 71 were detected as outliers in

341    this period. As is shown in Fig. 4 and Fig. 5, sample 199, 3, 363, 158 and 71 were also

342    detected as outliers in both MCS+VISSA order and MCS+VISSA order. In Fig. 9(a),

343    for ACE dataset, sample number 19, 53, 48, 91, 34 and 108 were detected as outliers

344    in this period. As is shown in Fig. 6 and Figure 7, sample 19, 53, 91, 34 and 108 were

345    also detected as outliers in both MCS+VISSA order and MCS+VISSA order. These

346    outliers can be detected by both separate and simultaneous methods. This means that

347    these samples are essentially far away from the main part of sample. These outliers

16

348    are not affected by variables. One can easy detect these outliers without considering

349    the impact of variables. 2), the middle iteration period (18-25 iteration for wheat

350    kernel dataset and 14-25 iteration for ACE dataset). In this period, the variables and

351    samples went ahead along with the tendency in the early period and most variables

352    and samples arrived 0 or 1. The weights of samples and variables which went down to

353    0 at this period varied dramatically. As is shown in Fig. 9(a), 9 samples were detected

354    as outlier and located in this period. Among them, none of them were detected in

355    MCS (Fig. 6), while 4 of them detected in VISSA+MCS (Fig. 7). This means that

356    these samples and outliers strongly affect each other in this period. When more

357    variables are selected, some samples with high weights may not proper for current

358    variables and result in weights decreasing, and vice versa; 3), in the ending iteration

359    period (26-31 iteration for wheat kernel dataset and 26-30 for ACE dataset), the

360    weights of the entire sample kept constant and all the outliers were detected and

361    removed. The rest is to optimize the variable subset which could wonderfully support

362    current selected samples. Finally, MASS converged to and found out the best model in

363    the model space through this model space shrink iteration procedure.

364    **4.3 Comparison with other methods**

365        The comparison with the standard procedures for outlier detection and variable

366    selection such as Williams plot of leverage values versus abnormal residuals (WP)[34],

367    Variable importance in projection (VIP)[35] and genetic algorithm (GA)[36] were also

368    performed, which were listed in Table 3. One can clearly see that the variable

369    selection method based on VIP on these two dataset obtained the poor prediction

17

370    statistics, and its results were similar [35]to those from original PLS models. VIP could

371    effectively eliminate some uninformative or noisy variables and therefore obtained a

372    relatively easy-to-interpret model. Compared to original PLS method without any

373    variable selection, VIP only selected 59 variables for Wheat kernel dataset and 31

374    variables for ACE dataset, respectively. Compared to original PLS and VIP method,

375    GA yields the better prediction results, and obtains the similar prediction performance

376    to VISSA and our proposed MASS. The final variable number used in the regression

377    model is also sharply reduced for these three methods. Compared to the commonly

378    used outlier detection method WP, from the Williams plot of leverage versus

379    abnormal residuals, there are 22 outliers (3, 25, 52, 71, 83, 104, 114, 199, 18, 33, 158,

380    208, 221, 231, 341, 363, 371, 397, 406, 408, 409, 411) in wheat kernel dataset and 4

381    outliers (8, 19, 53, 91) in ACE dataset. All the outliers detected in the first period of

382    MASS (Fig. 8 and Fig. 9) in wheat kernel dataset and ACE dataset were also detected

383    by WP, whereas the remaining outliers detected by these methods vary dramatically.

384    In ACE dataset, besides these four outliers no more outliers where detected by WP but

385    another 11 outliers were detected by MASS. In wheat kernel dataset, the outliers

386    detected by WP embodied most outliers in MASS (12 out of 16) and other 10 outliers.

387    After removing outliers, the $R^2$ and $Q^2_{CV}$ values WP were listed in Table 5. From

388    Table 5, compared with WP, MASS is also increased considerably for both datasets.

389    **4.4 The effect of MASS parameters**

390         In our proposed MASS method, three important parameters related to

391    Monte-Carlo sampling need to be set. The number of Monte-Carlo experiments seems

392    an important parameter which affects the quality of the distribution. Theoretically, the

18

393    fewer samples are selected randomly from the calibration samples, the more repeats

394    are needed. Whereas, it has been proven that the number of number of Monte-Carlo

395    experiments equal to $n^2$ ($n$ is the number of the total samples) is generally enough to

396    make Monte-Carlo strategy better performance. Larger sampling number tends to

397    generate more accurate and stable results. However, the accuracy improvement is very

398    small, indicating that MASS is insensitive to this parameter. To save the computing

399    source, in practice, the number of Monte-Carlo sampling is manually set to 2000. By

400    means of Monte-Carlo method, the computational complexity could be reduced

401    substantially. Similar to VISSA, the initial weights of variables (i.e., the variable

402    sampling ratio) were set to 0.5 in MASS, that is, each variable has 50% probability to

403    be selected in one sub-model in the beginning. Under the circumstance without any

404    prior information, it is a relatively natural choice to set the initial weights of variables

405    to 0.5. Another important parameter is the initial weight of samples (i.e., the sample

406    sampling ratio). To evaluate the influence of the initial weight of samples on MASS,

407    different experiments were carried out on the wheat kernels dataset. These results are

408    shown in Fig.10 and listed in Table 5. From Table 5, one can see that with the

409    decreasing of the initial samples weights, the number of iteration tends to slightly

410    increase and the number of outliers and variables changed a little (when the initial

411    weight is very small such as 0.5, the number of iterations, outliers or variables may

412    vary relatively large). The accuracy of MASS decreased a little when the initial

413    weights of samples decreased. Moreover, from Fig. 10, with different initial weight,

414    the variations of sample weight are similar. This indicates that MASS is insensitive to

415    the initial weights of samples to some extent. Given a dataset without any prior

416    information, we could assume that the main parts of samples are normal and only

417    small parts (e.g., 5%) of samples are outliers. Considering that, in my opinion, it is a

418    relatively suitable choice to set the sample initial weights to 0.95.

419    **4.5 The effect of variable and sample combinations**

420        Given a dataset, the number of combinations of variables and samples is

421    extremely high. Assume that we have a data set with n = 100 and p = 60, all

19

422 combinations from model space will be (2100-1)×(260-1), and this will be an

423 extremely high number for computer simulation. Therefore, it is almost impossible for

424 current computer simulation to enumerate all model combinations. Alternatively, were

425 randomly chosen some combinations from all possible combinations by Monte-Carlo

426 sampling strategy and then use the best part of generated models to represent the

427 distribution of important variables and samples. In general, Monte-Carlo approach can

428 be used to generate such a distribution of some statistic of interest by repeatedly

429 calculating that statistic randomly selected portions of the data because of its good

430 asymptotic properties. Through this sampling procedure, though the model

431 combination is usually high, only part of combinations was used which can

432 dramatically reduce the modeling time. Take a wheat kernel dataset for an example,

433 we only used about sixty thousands (2000×31 (31 iterations)) combinations to shrink

434 to a relatively good solution. We calculated the elapsed time of MASS on this dataset,

435 which is listed in Table 4. Although computation time of MASS is slightly higher than

436 those from other variable selection or outlier detection programs including MCS and

437 VISSA, it is worthy to waste somewhat more time to obtain a clean dataset and a

438 higher performance model. Simultaneous variables selection and outlier detection is

439 usually a hard task. We applied a computing-intensive method and therefor a little

440 more time was required. After MASS were performed, no additional codes between

441 variables selection and outlier detection was needed. From an overall perspective,

442 MASS takes less time in the model building process.

443 **Conclusion**

20

444   In this study, we proposed MASS to simultaneously detect outliers and select

445   variables before building a final prediction model. The proposed method is based on

446   MPA which iteratively and smoothly shrinks the model space to obtain the best model.

447   MASS is a mild stepwise optimization method. The model space shrinks smoothly

448   which reduce the risk of eliminating informative variables and normal samples. The

449   weights variation of variables and outliers illustrate the cross interaction between

450   variables and outliers: if the weights of variables and samples go down to 0 in the first

451   period, these variables and outliers do not interact with each other and they can be

452   easily identified. If the weights of variable and samples go down to 0 in the middle

453   period, these samples and outliers strongly affect each other. In the last period, the

454   weights of samples were constants, and the rest is to optimize the variable subset

455   which can wonderfully support current selected samples. The performance of the new

456   algorithm was compared with several other outlier detection and variable selection

457   methods and methods combination. The results clearly indicate that: when outlier

458   detection and variable selection performed separately, there is a great opportunity to

459   obtain a wrong model that fails to reflect the true relationship between variables and

460   outliers. To avoid this failure, it is recommended to do these tasks simultaneously. The

461   results demonstrated that MASS is a useful method in data cleaning before building a

462   predictive model.

463   ## Acknowledgement

21

469  # References

470  1.  G. C. Reinsel and R. P. Velu, in *Multivariate Reduced-Rank Regression*, Springer, 1998, pp.
471      1-14.
472  2.  W. W. Chin, *Modern methods for business research*, 1998, **295**, 295-336.
473  3.  A. J. Smola and B. Schölkopf, *Statistics and computing*, 2004, **14**, 199-222.
474  4.  V. Svetnik, A. Liaw, C. Tong, J. C. Culberson, R. P. Sheridan and B. P. Feuston, *Journal of*
475      *chemical information and computer sciences*, 2003, **43**, 1947-1958.
476  5.  F. G. Blanchet, P. Legendre and D. Borcard, *Ecology*, 2008, **89**, 2623-2632.
477  6.  J. M. Sutter and J. H. Kalivas, *Microchemical journal*, 1993, **47**, 60-66.
478  7.  M. Shahlaei, *Chemical reviews*, 2013, **113**, 8093-8103.
479  8.  B.-c. Deng, Y.-h. Yun, Y.-z. Liang and L.-z. Yi, *Analyst*, 2014, **139**, 4836-4845.
480  9.  M. Gevrey, I. Dimopoulos and S. Lek, *Ecological Modelling*, 2003, **160**, 249-264.
481  10. R. C. Eberhart and Y. Shi, 2001.
482  11. W. Cai, Y. Li and X. Shao, *Chemometrics and intelligent laboratory systems*, 2008, **90**,
483      188-194.
484  12. 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,
485      *Analytica chimica acta*, 2014, **807**, 36-43.
486  13. D. S. Cao, Y. Z. Liang, Q. S. Xu, H. D. Li and X. Chen, *Journal of computational chemistry*,
487      2010, **31**, 592-602.
488  14. P. J. Rousseeuw, *Journal of the American statistical association*, 1984, **79**, 871-880.
489  15. P. Filzmoser, *Computer data analysis and modeling. Robust and computer intensive methods.*
490      *Belarusian State University, Minsk*, 2001, 132-137.
491  16. J. A. Gil and R. Romera, *Journal of chemometrics*, 1998, **12**, 365-378.
492  17. J. Tolvi, *Soft Computing*, 2004, **8**, 527-533.
493  18. J. Hoeting, A. E. Raftery and D. Madigan, *Computational Statistics & Data Analysis*, 1996, **22**,
494      251-270.
495  19. P. Wiegand, R. Pell and E. Comas, *Chemometrics and intelligent laboratory systems*, 2009, **98**,
496      108-114.
497  20. R. Cavill, H. C. Keun, E. Holmes, J. C. Lindon, J. K. Nicholson and T. M. Ebbels,
498      *Bioinformatics*, 2009, **25**, 112-118.
499  21. R. S. Menjoge and R. E. Welsch, *Computational Statistics & Data Analysis*, 2010, **54**,
500      3181-3193.
501  22. S.-S. Kim, S. H. Park and W. Krzanowski, *Journal of Applied Statistics*, 2008, **35**, 283-291.
502  23. D. Cao, Y. Liang, Q. Xu, Y. Yun and H. Li, *Journal of computer-aided molecular design*, 2011,
503      **25**, 67-80.
504  24. H.-D. Li, Y.-Z. Liang, D.-S. Cao and Q.-S. Xu, *TrAC Trends in Analytical Chemistry*, 2012, **38**,

22

154-162.

25. B.-C. Deng, Y.-H. Yun, P. Ma, C.-C. Lin, D.-B. Ren and Y.-Z. Liang, *Analyst*, 2015, **140**, 1876-1885.

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, *Analytica chimica acta*, 2015, **862**, 14-23.

27. R. Kohavi, 1995.

28. N. J. Nagelkerke, *Biometrika*, 1991, **78**, 691-692.

29. J. J. Sutherland, L. A. O'Brien and D. F. Weaver, *Journal of Medicinal Chemistry*, 2004, **47**, 5541-5554.

30. J. B. Wang, D. S. Cao, M. F. Zhu, Y. H. Yun, N. Xiao and Y. Z. Liang, *Journal of Chemometrics*, 2015.

31. D. S. Cao, Q. S. Xu, Y. Z. Liang, X. Chen and H. D. Li, *Journal of Chemometrics*, 2010, **24**, 584-595.

32. E. Pourbasheer, S. Shokouhi Tabar, V. Masand, R. Aalizadeh and M. Ganjali, *SAR and QSAR in Environmental Research*, 2015, **26**, 461-477.

33. N. Xiao, D.-S. Cao and Q.-S. Xu, 2015.

34. B. S. Everitt, *American Mathematical Monthly*, 1998, 387-388.

35. L. Eriksson, E. Johansson, H. Antti and E. Holmes, *Multi- and Megavariate Data Analysis*, 2005.

36. R. Leardi and A. L. González, *Chemometrics & Intelligent Laboratory Systems*, 1998, **41**, 195-207.

23

528 **Tables:**

529 Table 1    The results of wheat kernel dataset performed on different methods.

| Method | Sam | Var | $R^2$ | $Q^2_{CV}$ | optPC | Iteration |
|---|---|---|---|---|---|---|
| **PLS** | 415 | 100 | 0.880 | $0.868\pm0.005$ | 10 | |
| **VISSA** | 415 | $32\pm2$ | $0.894\pm0.001$ | $0.886\pm0.003$ | 9 | $13\pm2$ |
| **MCS** | $402\pm0$ | 100 | $0.899\pm0$ | $0.889\pm0.003$ | 10 | |
| **VISSA + MCS** | $404\pm2$ | $32\pm2$ | $0.909\pm0.002$ | $0.902\pm0.002$ | 9 | $13\pm2$ |
| **MCS+ VISSA** | $402\pm0$ | $31\pm3$ | $0.911\pm0.001$ | $0.904\pm0.003$ | 9 | $13\pm2$ |
| **MASS** | $398\pm2$ | $31\pm5$ | $0.921\pm0.003$ | $0.913\pm0.005$ | 10 | $31\pm4$ |

530

531 Table 2    The results of ACE dataset performed on different methods.

| Method | Sam | Var | $R^2$ | $Q^2_{CV}$ | optPC | Iteration |
|---|---|---|---|---|---|---|
| **PLS** | 114 | 56 | $0.745\pm0$ | $0.623\pm0.038$ | 10 | - |
| **VISSA** | 114 | $23\pm10$ | $0.775\pm0.023$ | $0694\pm0.036$ | 10 | $12\pm3$ |
| **MCS** | 102 | 56 | 0.819 | $0.729\pm0.033$ | 10 | - |
| **VISSA + MCS** | $106\pm3$ | $23\pm10$ | $0.837\pm0.017$ | $0.772\pm0.044$ | 10 | $12\pm3$ |
| **MCS+ VISSA** | 102 | $30\pm13$ | $0.841\pm0.017$ | $0.775\pm0.031$ | 10 | $12\pm3$ |
| **MASS** | $102\pm3$ | $26\pm10$ | $0.865\pm0.021$ | $0.823\pm0.027$ | 10 | $24\pm5$ |

532

533

24

534    Table 3    The results of Wheat kernel and ACE dataset performed on VIP and GA.

| Dataset | Methods | Sam | *Var* | $R^2$ | $Q^2_{CV}$ |
|---|---|---|---|---|---|
| | WP | 393 | 100 | 0.913 | 0.904 |
| Wheat kernel dataset | VIP | 415 | 59 | 0.877 | 0.861 |
| | GA | 415 | 34 | 0.891 | 0.881 |
| | WP | 110 | 56 | 0.774 | 0.691 |
| ACE dataset | VIP | 114 | 31 | 0.728 | 0.624 |
| | GA | 114 | 13 | 0.772 | 0.703 |

535

536    Table 4    The elapsed time of MCS, VISSA and MASS.

| | Wheat kernel | | ACE | |
|---|---|---|---|---|
| Methods | Time (second) | $Q^2_{CV}$ | Time (second) | $Q^2_{CV}$ |
| MCS | 30 | 0.889 | 12 | 0.729 |
| VISSA | 810 | 0.889 | 688 | 0694 |
| MASS | 1260 | 0.917 | 981 | 0.823 |

537

538    Table 5    The performance of MASS with different sample initial weight.

| Initial weight | 0.95 | 0.9 | 0.8 | 0.7 | 0.6 | 0.5 |
|---|---|---|---|---|---|---|
| Number of iterations | 34 | 40 | 37 | 49 | 39 | 54 |
| Number of outliers | 17 | 13 | 14 | 18 | 13 | 29 |
| Number of variables | 34 | 28 | 36 | 36 | 13 | 17 |
| $Q^2_{CV}$ | 0.9173 | 0.9109 | 0.9111 | 0.9155 | 0.8783 | 0.9168 |

539

25
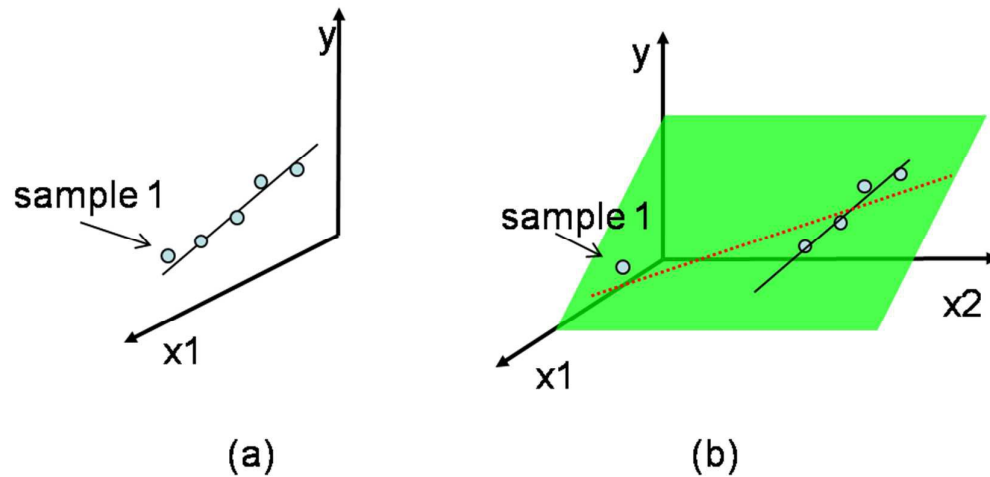
540    **Figure captions:**

541



542

543    Fig. 1 The sketch of model space. A model is constructed by the combination of some variables

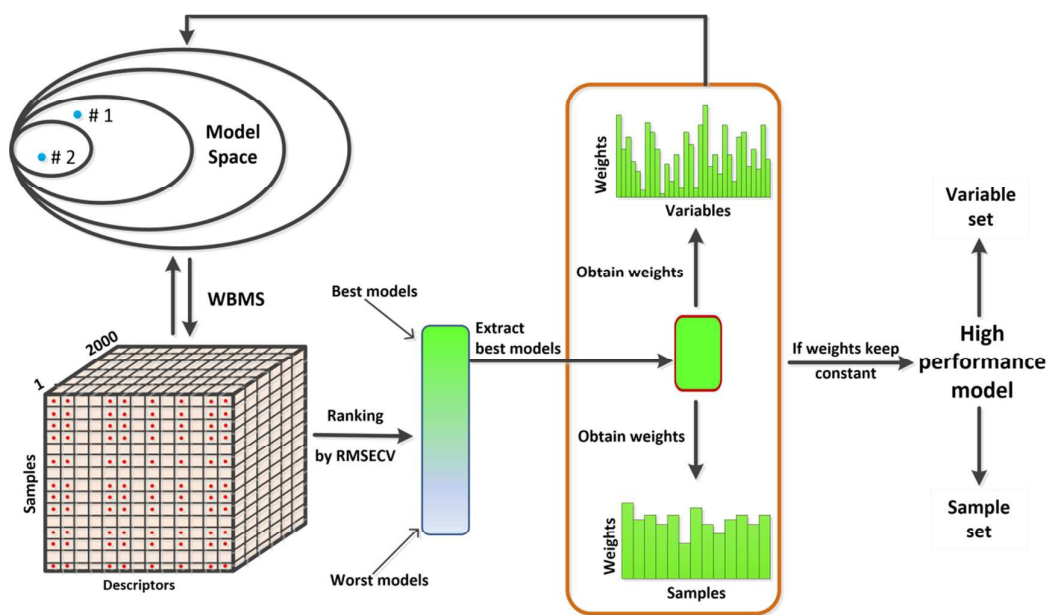544    and samples like #1 and #2, all the combinations make up the model space.

545



546

547    Fig. 2 The interactions between variables and outliers. (a) With only one variable $x1$, all samples

548    (including sample 1) can be well fitted. (b) When variable $x2$ was added, sample 1 turns into an

549    outlier.

550

551

26

552

553    Fig. 3 The framework of MASS. Firstly, a number of sub-training datasets were sampled from the

554    original dataset and build sub-regression models. The frequency of each variable and each sample

555    in the best part models were counted. Then, a number of new sub-training datasets were sampled

556    using the weight of variables and sample obtained from last iteration. The procedure for obtaining

557    new weights of variable and sample was repeated until the weights of all variables and samples

558    were constant (either 1 or 0). Thus the best model was obtained; the variables and samples that

559    constructed the best model are simultaneously selected.

560

561

562

563    Fig. 4 Outlier detection plot of wheat kernel dataset detected by MCS+VISSA order. The samples

564    enclosed by red ellipse located in lower left area and they are normal samples. Whereas they

565    turned into outliers after variables selection (see Fig. 5). The samples enclosed by green ellipses

566    are outliers but they become normal samples after variables selection (see Fig. 5)
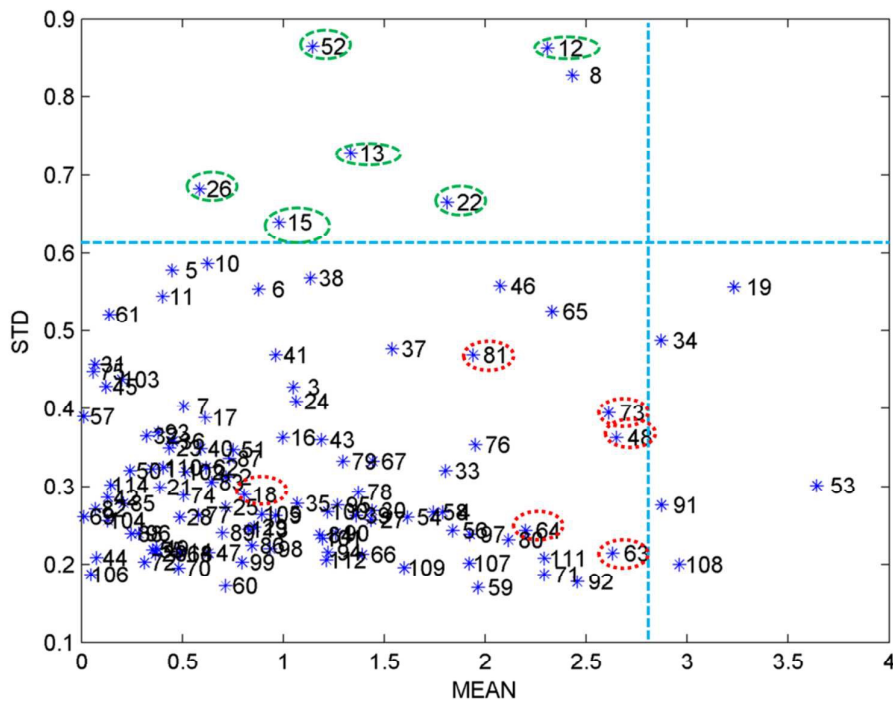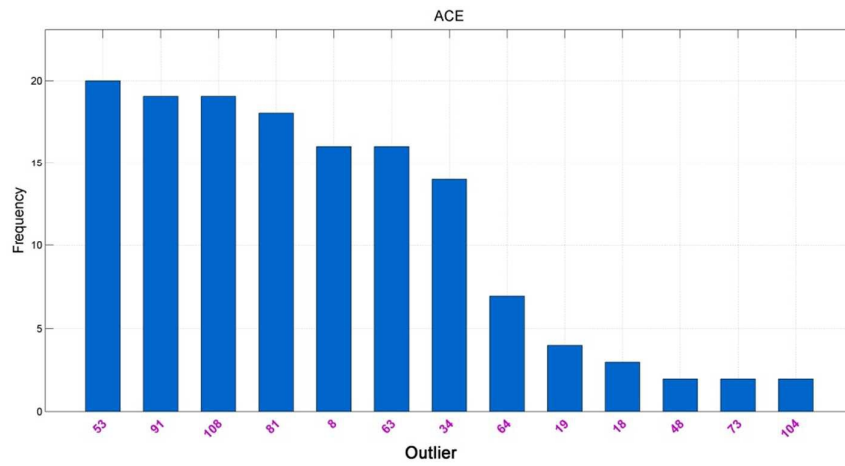
567

28

568

569    Fig. 5 Frequencies of outliers detected in wheat kernel dataset in 20 times in VISSA+MCS order.

570



571

572    Fig. 6 Outlier detection plot of ACE dataset detected by MCS+VISSA order. The samples

573    enclosed by red ellipse located in lower left area and they are normal samples. Whereas they

574    turned into outliers after variables selection (see Fig. 7). The samples enclosed by green ellipses

575    are outliers but they become normal samples after variables selection (see Fig. 7)
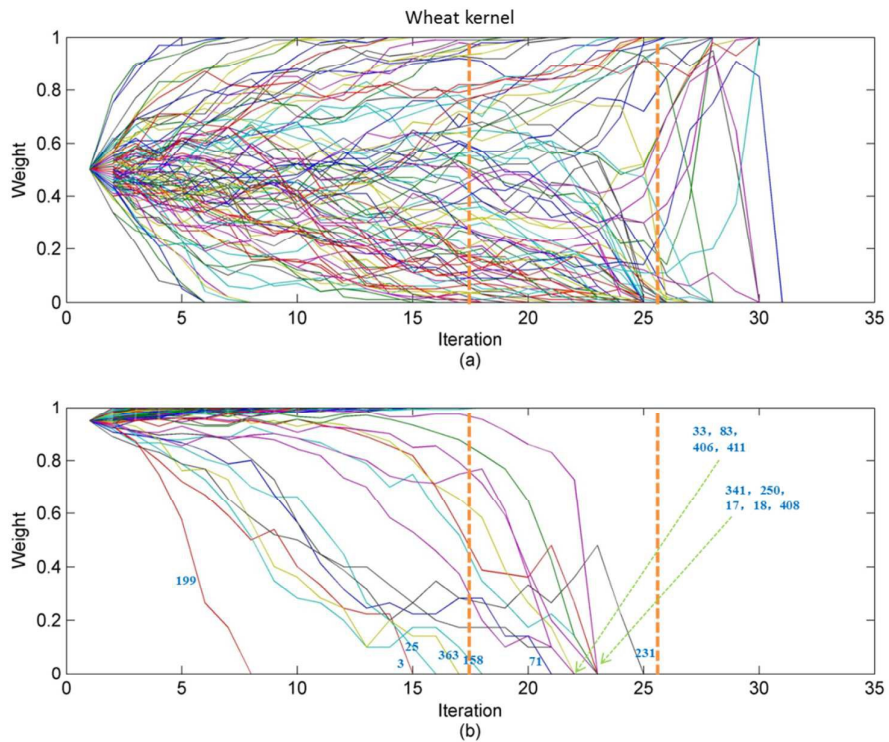
29

576



577

578    Fig. 7 Frequencies of outliers detected in ACE dataset in 20 times in VISSA+MCS order.
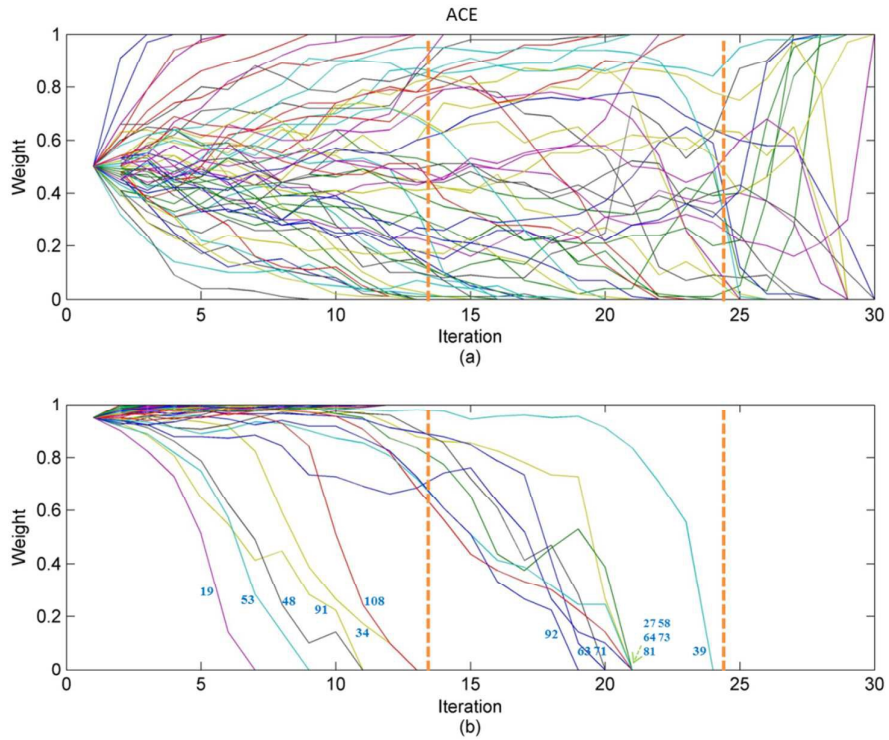
579



580

581    Fig. 8 The weight variation of (a) variables and (b) samples of wheat kernel dataset. Each line

582    represents the weight variation of a sample or a variable.

30

1
2
3
4          583
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
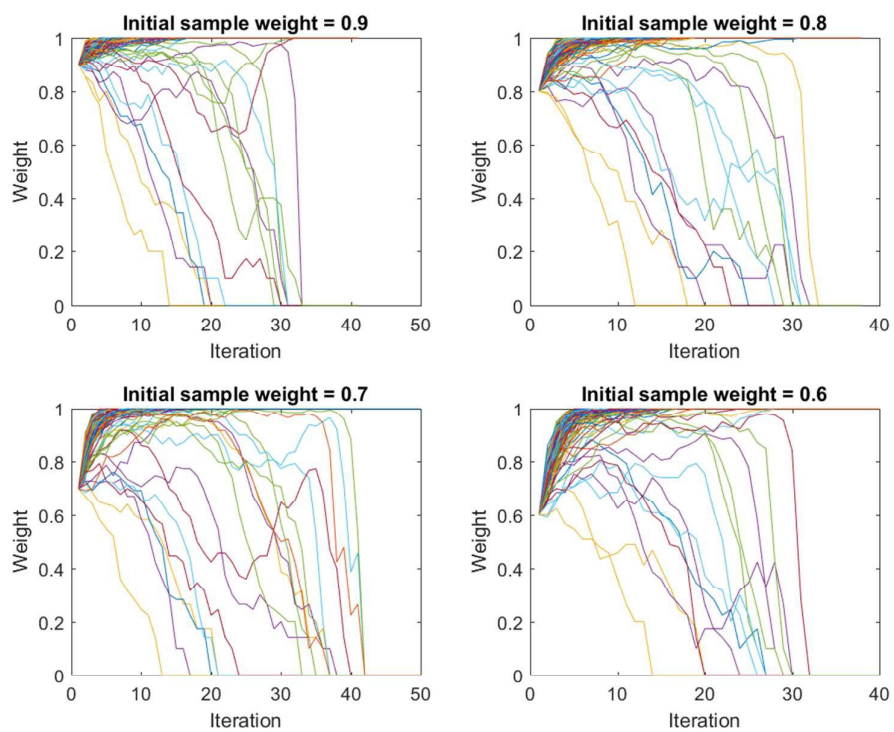27
28
29
30          584
31
32
33          585    Fig. 9 The weight variation of (a) variables and (b) samples of ACE data set. Each line represents
34
35          586    the weight variation of a sample or a variable.
36
37          587
38

31

588

589

590     Fig. 10     The weight variation of samples of wheat kernel data set with different sample initial

591                 weight. Each line represents the weight variation of a sample.

592

593

32