

Document downloaded from:

<http://hdl.handle.net/10251/52301>

This paper must be cited as:

Kuligowski, J.; Pérez Guaita, D.; Escobar, J.; Guardia, MDL.; Vento, M.; Ferrer Riquelme, AJ.; Quintás, G. (2013). Evaluation of the effect of chance correlations on variable selection using Partial Least Squares -Discriminant Analysis. *Talanta*. 116:835-840. doi:10.1016/j.talanta.2013.07.048.



The final publication is available at

<http://dx.doi.org/10.1016/j.talanta.2013.07.048>

Copyright Elsevier

Evaluation of the effect of chance correlations on variable selection using Partial Least Squares – Discriminant Analysis

Julia Kuligowski^{a,§}, David Pérez-Guaita^{b,§}, Javier Escobar^a, Miguel de la Guardia^b, Máximo Vento^{a,c}, Alberto Ferrer^d, Guillermo Quintás^{*,e}

^aNeonatal Research Centre, Health Research Institute La Fé, 46009 Valencia, Spain

^bDepartment of Analytical Chemistry, University of Valencia, 46100 Burjassot, Spain

^cDivision of Neonatology, University & Polytechnic Hospital La Fé, 46021 Valencia, Spain

^dDepartment of Applied Statistics, Operations Research and Quality, Universidad Politécnica de Valencia, 46071 Valencia, Spain

^eLeitat Technological Center, Bio *In Vitro* Division, 08225 Terrassa, Spain, e-mail: guillermo.r.quintas@uv.es,

Tel: +34 96 354 4838, Fax: +34 96 354 4845

[§] Both authors contributed equally to this work.

Abstract

Variable subset selection is often mandatory in high throughput metabolomics and proteomics. However, depending on the variable to sample ratio there is a significant susceptibility of variable selection towards chance correlations. The evaluation of the predictive capabilities of PLSDA models estimated by cross-validation after feature selection provides overly optimistic results if the selection is performed on the entire set and no external validation set is available. In this work, a simulation of the statistical null hypothesis is proposed to test whether the discrimination capability of a PLSDA model after variable selection estimated by cross model validation is statistically higher than that attributed to the presence of chance correlations in the original data set. Statistical significance of PLSDA CV-figures of merit obtained after variable selection is expressed by means of p-values

27 calculated by using a permutation test that included the variable selection step. The
28 reliability of the approach is evaluated using two variable selection methods on
29 experimental and simulated data sets with and without induced class differences. The
30 proposed approach can be considered as a useful tool when no external validation set is
31 available and provides a straightforward way to evaluate differences between variable
32 selection methods.

33 **KEYWORDS:** metabolomics; chance correlations; variable selection; Partial Least Squares -
34 Discriminant Analysis (PLSDA)

35

36 **1. Introduction**

37 Nowadays, nuclear magnetic resonance (NMR) and the hyphenation of high resolution
38 separation techniques (e.g. gas and liquid chromatography as well as capillary
39 electrophoresis) with mass spectrometry (MS) play leading roles as high throughput
40 analytical tools in comprehensive metabolomics and proteomics. Frequently, studies involve
41 the discriminant analysis of samples under two distinct experimental conditions such as
42 treated vs. untreated or diseased vs. control, for the identification of biomarkers or the
43 calculation of predictive models. This is a challenging task, as the number of detected
44 variables typically largely exceeds the number of samples, and variables are usually
45 correlated. Moreover the unambiguous identification of metabolites or proteins can be
46 highly difficult, and the concentration and response ranges involved normally cover several
47 orders of magnitude. Besides, the majority of the detected variables are frequently irrelevant
48 for the outcome prediction [1-3] and so the predictive precision and accuracy of
49 discriminant models can be improved if uninformative variables are removed in advance
50 [4,5]. Furthermore, feature selection also provides simplified models of easier
51 interpretability in a subsequent statistical or biochemical data analysis.

52 Whereas a wide range of multivariate methods for supervised learning (i.e. pattern
53 recognition) is available, each with its own strengths and weaknesses, the most commonly
54 used multivariate classification technique is Partial Least Squares - Discriminant Analysis
55 (PLSDA) [4]. PLSDA is a multivariate PLS method that extracts a set of latent variables (LVs)
56 that explain the sources of variation in the \mathbf{X} -block correlated to an \mathbf{y} -vector that encodes the
57 class membership [1]. One of the key features of PLSDA is its applicability in situations in
58 which variables far outnumber samples, and correlation among variables exists [2,6]. In a
59 PLSDA model, the relation between the predictors \mathbf{X} ($N \times J$) and the response \mathbf{y} ($N \times 1$) can be
60 described as:

61
$$\mathbf{y} = \mathbf{X}\mathbf{b}^T + \mathbf{e} \text{ (Equation 1)}$$

62 where \mathbf{b} ($1 \times J$) is the vector of regression coefficients, \mathbf{e} ($N \times 1$) is the error vector (i.e.
63 residuals) and N and J are the number of objects (e.g. MS or NMR spectra) and variables (e.g.
64 m/z features or chemical shifts), respectively.

65 In metabolomic and proteomic studies, results should be subjected to thorough statistical
66 and biological validation. Whereas the biological validation determines whether biomarkers
67 are involved in processes related to the stated difference between classes, the statistical
68 validation determines the performance of the biomarker and the probability of a chance
69 result [7]. There are two statistical validation approaches, namely external and cross-
70 validation. While external validation is considered the 'gold standard', cross-validation (CV)
71 can be seen as a sub-optimal approximation to external validation [7] that, in spite of its
72 limitations, still is very useful in case of a limited number of samples. Cross-validation is
73 used for both the selection of the complexity of PLS models and to obtain an estimation of
74 their predictive performance. During CV, a subset of objects from the data set is removed (i.e.
75 validation set) and a PLS model is calculated using the remaining objects (i.e. training set).
76 Then, the calculated model is used for the prediction of the \mathbf{y} values of the validation set, and
77 averaging over several splits yields the CV estimation of the model performance. CV methods
78 are classified according to the procedure employed for the selection of the different subsets.
79 In spite of being widely employed, CV increases the risk of model over-fitting and it also
80 provides overoptimistic internal figures of merit in explorative and predictive PLS analysis
81 [9,10]. Double cross-validation (2CV), also known as cross model validation, is an alternative
82 CV strategy that circumvents these drawbacks providing external figures of merit [9-12]. In
83 2CV a subset of objects is set aside as a test set. The remaining set of objects are again split
84 into training and validation sets, and they are subjected to a standard CV procedure for the
85 selection of the number of latent variables [3]. Besides, non-parametric permutation tests

86 based on random rearrangements of the elements of the y vector of a data set are useful for
87 determining the significance of a statistic [9,13] and its use in combination with 2CV has
88 been repeatedly proposed as a suitable approach to assess the statistical significance of
89 PLSDA figures of merit [2,3,9,10,14].

90 As aforementioned, dimensionality reduction methods are often employed to increase PLS
91 prediction accuracy. If variable selection is performed in advance on the entire data set, it
92 gives overly optimistic CV results. This apparent improvement, however, partly originates
93 from the susceptibility of variable selection towards chance correlations depending on both
94 the variable to sample ratio and the correlation structure of data [15-18].

95 Addressing the aforementioned concerns, a straightforward strategy based on permutation
96 testing and 2CV is proposed to grade the effect of chance correlations on PLSDA model
97 performance during variable selection. For this purpose, the number of misclassified
98 samples (NMC) and the discriminant Q^2 (dQ^2) [19] performance statistics calculated using
99 real class labels were compared to a distribution of the same estimators obtained after class
100 randomization before and after variable selection. Simulated data sets, an experimental MS
101 data set and two variable selection procedures were used to demonstrate the potentials and
102 drawbacks of the approach.

103

104 **2. Material and Methods**

105 **2.1 Software**

106 Data analysis was run under Matlab 7.7.0 from Mathworks (Natick, USA, 2004) using in-
107 house written MATLAB functions and the PLS Toolbox 6.2 from Eigenvector Research Inc.
108 (Wenatchee, WA, USA). Bold capital letters represent matrices, bold italic lowercase
109 characters represent vectors, and italic uppercase letters represent scalars. Both simulated

110 | and experimental data sets are assembled in matrices \mathbf{X} ($N \times J$), where rows (N) and columns
111 | (J) correspond to samples and variables, respectively.

112

113 2.2 Data sets

114 Four 'null' data sets (N_{250} (60×250), N_{540} (60×540), N_{1000} (60×1000), and N_{2000} ($60 \times$
115 2000)) were generated using the *randn* MATLAB function [20] and contained
116 pseudorandom values drawn from standard normal distributions (i.e. mean zero and
117 standard deviation one). The first 30 objects were classified as class A ($Y=1$) and the rest as
118 class B ($Y=0$).

119 Then, a set of simulated data sets (**SIMUIN_5**, **SIMUIN_15** and **SIMUIN_25**) was calculated as
120 described by Centner et al. [21]: **SIM** ($60 \times V$) was a simulated pure (noise free) data matrix
121 generated with an exact dimensionality of 3, only containing informative variables. The
122 **SIMUI** (60×540) data matrix resulted from the attachment of an uninformative variable
123 matrix **UI** ($60 \times (540-V)$) to the **SIM** matrix. **UI** consisted of pseudorandom numbers drawn
124 from a standard normal distribution. **SIMUIN** is the sum of the **SIMUI** matrix and a noise
125 matrix **N** (60×540) containing pseudorandom numbers drawn from a normal distribution
126 with mean zero and standard deviation 0.025. The **SIMUIN_5**, **SIMUIN_15** and **SIMUIN_25**
127 data sets corresponded to $V=5$, 15 and 25, respectively. For each **SIM** matrix, a \mathbf{y} (60×1)
128 vector was calculated as $\mathbf{y} = 3 \cdot \mathbf{t}_1 + 2 \cdot \mathbf{t}_2 + 1 \cdot \mathbf{t}_3$, where \mathbf{t}_a ($a=\{1,2,3\}$), is the vector of scores of
129 the a^{th} principal component. Class assignment of each simulated sample was carried out
130 according to the sign of its calculated y value (see Figure 1). A new \mathbf{y} vector was generated
131 where each class A sample was assigned a value of zero and each class B sample a value of
132 | one. The use of class labels ($\underline{1/0}$) instead of the actual y values was selected to simulate real
133 | situations where samples are typically clustered in two classes in spite of within-class

134 differences among samples. Then, 20 randomly selected samples of each class were removed
135 from each data set for being used as external test sets.

136 An experimental data set was used to test the applicability of the approach. The employed
137 data set (**Gaucher** (40 x 590)) was obtained from the Biosystems Data Analysis Group
138 website (www.bdagroup.nl) and contains Surface Enhanced Laser Desorption Ionization –
139 Time of Flight – Mass Spectrometry (SELDI-TOF-MS) data of 40 serum samples from 20
140 Gauchy patients and 20 healthy controls. Each sample spectrum consists of 590 m/z
141 variables between 1000 and 10000. Y values of 1 and 0 were assigned to the spectra
142 obtained from Gauchy and healthy patients, respectively. Background information on the
143 **Gaucher** data set can be found in a previous work [22] and on the aforementioned website.

144 *PLS modelling*

145 Prior to PLS model calculation, autoscaling was employed to equal the relative importance of
146 all variables. The \mathbf{y} vector containing the class labels was mean centered. Scaling factors
147 were calculated from the calibration subsets. No outlier detection was performed and all
148 samples were used for variable selection and 2CV, employing a maximum of 5 PLS
149 components selected from dQ^2 values calculated by CV.

150

151 **2.3 Variable selection procedures**

152 The following variable selection procedures were considered:

153 *Approach 1. Variance of the \mathbf{b} regression vector (b_{cv} -PLSDA)*

154 This approach uses the set of PLSDA regression vectors obtained after M random K -fold
155 cross-validations ($M=20$ and $K=4$ in this work). The number of LVs included in a PLSDA
156 model was selected from dQ^2 values obtained after each K -fold CV. Then, the mean vector ($\bar{\mathbf{b}}$)
157 and the standard deviation vector (\mathbf{s}_b) of the regression coefficient vectors (\mathbf{b}) were
158 calculated. Variables were selected as informative according to $(|\bar{b}_j| - d s_{b,j}) > 0$ ($j=1, \dots, J$)

159 The value of d should be selected as a compromise between Type I and II errors. In this
160 work $d=4$ was used.

161

162 *Approach 2. Uninformative Variable Elimination (UVE-PLSDA)*

163 The second approach involved UVE-PLSDA [21]. In this procedure, the original data matrix is
164 augmented column-wise by a matrix containing normally distributed artificial random
165 variables of very small magnitude (e.g. 10^9 times lower than the real variables [21]). Then,
166 the standard deviation vector of the regression coefficients (\mathbf{s}_b) is obtained from the
167 variation of the PLS regression coefficients by leave-one-out CV. The obtained values are
168 used to calculate a reliability coefficient (c_j), which is an equivalent to the calculated t-value,
169 for each original variable according to Equation 2:

$$170 \quad c_j = \frac{b_j}{std(b_j)} \quad (\text{Equation 2})$$

171 Different criteria have been proposed [9] to establish the cut-off level for classification of
172 real variables as informative using the reliability values of the artificial (uninformative)
173 variables ($c_{artif,j}$). In this work, the UVE-R approach was employed where $|c_{artif,j}|$ are ranked
174 and a cut-off level corresponding to a defined α -quantile [21,24] is selected. Due to random
175 generation of artificial variables, results found by UVE-PLSDA show certain variability as
176 real variables with c_j values close to the cut-off level might be retained or not in the final
177 model depending on minor differences in $c_{artif,j}$. Daszykowsky et al. [24] improved the
178 performance of UVE-PLSDA by using a Monte Carlo approach in which UVE-PLSDA was
179 repeated a number of times and in each run a randomly selected model set was used for
180 model construction and feature selection.

181 The UVE-PLSDA procedure followed in this work can be described in four steps:

182 *i. Matrix augmentation:* The original data matrix \mathbf{X} was augmented column-wise by an
183 artificial variable matrix \mathbf{R} ($N \times 250$) with random values drawn from a standard distribution
184 with mean zero and standard deviation 10^{-9} [21].

185 *ii. Calculation of PLSDA submodels:* A series of N submodels for the augmented data matrix
186 was calculated by leave-one-out CV. Accordingly, for each submodel, $(N-1)$ samples were
187 used to calculate an inner PLSDA model of complexities $a=\{1,\dots,A\}$, which was subsequently
188 used for prediction of the y_i value of the remaining sample. The procedure was repeated
189 until all samples were predicted once as validation sample. The PLSDA model complexity
190 was selected from dQ^2 values calculated using the predicted \mathbf{y} values.

191 *iii. Calculation of the reliability coefficient for each variable:* The reliability coefficient was
192 calculated from the set of N regression vectors according to Equation 2.

193 *iv. Cut-off selection and UVE:* In this work, the α -quantile value of 99% was selected as cut-off
194 value for a pre-classification of real variables as informative. To reduce the variability due to
195 the use of random artificial variables during variable selection, the UVE-PLSDA process was
196 repeated a total of 1000 times. By doing this, a frequency (γ) of pre-classification as
197 informative was obtained for each variable. As informative variables are expected to be
198 retained more frequently than uninformative variables, $\gamma=99\%$ was used as a second
199 threshold value for variable discrimination.

200 The predictive performance of the PLSDA models after variable selection by both procedures
201 was estimated by 4-fold 2CV, as described elsewhere [9]. The random selection of 2CV
202 training + validation and test sample subsets was repeated M times ($M=20$ in this work) to
203 reduce the influence of the split on the results. Again, dQ^2 values calculated by leave-one-out
204 CV within each training set were used to optimize the number of LVs of each inner model.
205 Finally, average NMC and dQ^2 statistics were calculated from the obtained 2CV results.

206

207 2.4 Permutation test

208 In order to estimate the statistical significance of figures of merit obtained after variable
209 selection, a permutation test was carried out to create a null distribution. Accordingly, the
210 evaluation of the PLSDA performance using the selected variables was repeated 2000 times
211 using randomly permuted class labels. P-values for the figures of merit were calculated
212 either empirically [3] or by tail approximation to a generalized pareto distribution (GPD)
213 [13].

214 In the empirical approach, the p-value was computed as the fraction of permuted statistics
215 that are at least as extreme as the test statistic obtained from the original data, as described
216 elsewhere [3]. As the minimum type-I risk after z iterations calculated this way is $1/z$, this
217 empirical approximation becomes impractical when the calculation of each random statistic
218 is computing intensive. In the second approach the (right) tail of the distribution of
219 permutation values (i.e. x in Equation 3) using a maximum of 250 points is fitted to a
220 generalized Pareto distribution with the following cumulative distribution function:

$$221 \quad F(x) = 1 - (1 - kx\alpha^{-1})^{1/k}, \text{ for } k \neq 0 \text{ (Equation 3)}$$

222 [This method is based on tail approximation and reduces the number of required](#)
223 [permutations to accurately provide small p-values \[13\].](#) Detailed descriptions of [the method,](#)
224 [the](#) data pretreatment and fitting procedure can be found in literature [13]. After estimation
225 of both k and α parameters in Equation 3, the Anderson-Darling statistic (A^2) was calculated
226 for the estimation of the goodness of fit of the data to a GPD [25]. If the test failed, the
227 smallest exceedance was eliminated and the GPD fit was tested again. It has been
228 demonstrated that this method provides accurate p-values with a reduced number of
229 permutations as compared to the standard empirical approach. This advantage is of special
230 importance when the number of permuted values exceeding the test statistic (f) is very low

231 ($f \leq 10$, in this work) and the permutation approach is computing intensive. Nonetheless, in
232 situations where the GPD fitting failed, the empirical p-value was used.

233

234 **3. Results and discussion**

235 *3.1 Simulated data sets*

236 First a study assessing the potential of the proposed approach to estimate the statistical
237 significance of chance correlations on the improvement of PLSDA models after variable
238 selection was performed using simulated data. Then the same approach was applied to the
239 **Gaucher** data set.

240 *Null data sets*

241 Since predictors and responses were randomly generated in the null data sets, only non-
242 statistically significant PLSDA models were expected before variable selection [2,3,9-11,15].
243 Accordingly, 2CV figures of merit (NMC and dQ^2) obtained before variable selection showed
244 no class difference (e.g. NMC around 50% of the samples) for all four null data sets
245 independently of their size (see Table 1). Nonetheless, the higher the variables-to-samples
246 ratio, the higher the probability of finding a subset of variables with different distributions
247 between classes because of sheer coincidence [7, 15-19]. Consequently, after variable
248 selection the number of variables selected as informative increased and figures of merit of
249 the submodels improved with the numbers of variables in the original data set (see Table 1).

250 For example, whereas the NMC for the Null dataset with 250 variables is reduced from 22 to
251 6 or 13, for the Null dataset with 2000 variables, the NMC can be artificially reduced from 24
252 to 0 after variable selection. This effect could be clearly seen using the b_{cv} -PLSDA approach
253 where a correlation among the number of variables in the original dataset, the number of
254 retained variables and overoptimistic CV results were found. In spite of that, the
255 permutation test showed the lack of statistical significance of the PLSDA submodels,

256 expressed by the calculated p-values for both NMC and dQ^2 obtained using real class labels
257 in comparison to those obtained from permutation testing, all giving p-values > 0.05 .

258 This overoptimistic effect was further confirmed by results obtained for the simulated
259 external test sets: whereas using variables selected by both b_{cv} -PLSDA and UVE-PLSDA
260 approaches the number of misclassified samples employing cross-validation decreased
261 rapidly (see Table 1), the number of misclassified samples in the external validation sets
262 remained constant as shown in Table 2. [For example, for the Null dataset \(20 x 1000\) the](#)
263 [NMC in the external validation set before and after variable selection remains constant](#)
264 [\(equal to 11\)](#). Additionally, the number of selected variables for each null data set was close
265 to the mean value of the number of retained variables using randomly permuted class labels.
266 This can be appreciated from Figures 2a and 3a for a **Null** (40 x 540) data set. Likewise,
267 Figures 2b-c and 3b-c confirm that also NMC and dQ^2 values obtained for the same null data
268 set are close to the mean value obtained for randomly permuted class labels, using both
269 variable selection approaches.

270 In summary, results obtained from null data sets demonstrate that the evaluation of the
271 statistical significance of figures of merit obtained after variable selection can be used to
272 conclude whether there is a statistically significant difference between classes in the original
273 data set. Nonetheless, this procedure is computing intensive and alternative approaches can
274 also provide the same information faster with the same accuracy level [9].

275 *SIMUIN data sets*

276 Table 1 summarizes results obtained for the **SIMUIN5**, **SIMUIN15** and **SIMUIN25** data sets
277 in which, while the variables-to-samples ratio was kept constant ([540/40](#)), the number of *a*
278 *priori* informative variables increased [from 5/540 up to 25/540](#). Results showed that the
279 b_{cv} -PLSDA method retained percentages of *a priori* informative variables in the 53-60%
280 range, and NMC as well as dQ^2 values were substantially improved after variable selection.

281 Moreover, as depicted in Figure 2 for the SIMUIN data sets, the number of variables retained
282 was higher than those kept using randomly permuted class labels. Besides, statistically
283 significant p-values were obtained for the figures of merit of the submodels thus indicating
284 that the hypothesis that figures of merit using real and random class labels were equal could
285 be rejected ($\alpha=0.05$) and so, improvements were not exclusively due to existing chance
286 correlations in the original data set. The suitability of both variable selection methods and
287 the significance tests was also supported by lower NMC in the external validation sets after
288 variable selection, as summarized in Table 2.

289 Likewise, results obtained from UVE-PLSDA for data sets **SIMUIN15** and **SIMUIN25**
290 provided improved PLSDA figures of merit as shown in Tables 1 and 2 concerning the NMC
291 and dQ^2 values obtained from 2CV as well as for the external validation set. Also results
292 depicted in Figure 3 are in good agreement with those obtained by b_{cv} -PLSDA showing the
293 same trends in the number of selected variables, NMC and dQ^2 for SIMUIN data sets.
294 Although the NMC and dQ^2 values for **SIMUIN5** employing UVE-PLSDA had slightly
295 improved after variable selection, indicating an improvement of submodel performance, p-
296 values obtained for both, NMC and dQ^2 (see Table 1) indicated that the results obtained after
297 variable selection were comparable to those due to chance correlations. This could also be
298 confirmed by the NMC in the external validation set, increasing from 6 to 10 after variable
299 selection. Indeed, whereas 3 out of 5 informative variables were selected using the b_{cv} -
300 PLSDA approach, none of the those variables was retained by UVE-PLSDA. The observed
301 differences between results found after b_{cv} -PLSDA and UVE-PLSDA selection in case of
302 **SIMUIN5** were likely due to the effect of α and γ values on the set of retrieved UVE variables:
303 whereas low values increase the number of both informative and uninformative variables
304 retained, high thresholds may lead to a loss of useful information thus reducing the
305 predictive capabilities of PLS models calculated after variable elimination.

306 *Gaucher data set*

307 The **Gaucher** (40 x 590) data set was obtained from a study focusing on the measurement of
308 the protein profiles of serum of symptomatic type I Gaucher patients (n=20) and controls
309 (n=20) [22]. A total of 52 and 47 variables, 11 in common, were retained in the final models
310 using b_{cv} -PLSDA and UVE-PLSDA, respectively. Results of this study using the two
311 considered variable selection approaches provided p-values <0.05 for both dQ^2 and NMC as
312 it can be seen in Table 1. In Figures 2 and 3 it can be appreciated that the numbers of
313 selected variables, NMC and dQ^2 values are different from the mean values obtained from
314 permutation testing lying at the side of the random distributions as confirmed by the p-
315 values shown in Table 1. It is interesting that all 10 variables identified in a previous work
316 [22] as those with the largest contribution to the discrimination were included in both
317 variable subsets.

318 When comparing figures of merit before and after UVE-PLSDA variable selection, results
319 obtained were worse than it could have been expected. For example, the NMC after variable
320 selection for the **SIMUIN_5** external test set were clearly worse than those found by using
321 b_{cv} -PLSDA. The same effect was observed for the **Gaucher** data set where variable selection
322 did not reduce the NMC.

323 Whereas the effect of chance correlations could not be eliminated (i.e. CV after variable
324 selection provided overoptimistic figures of merit as it can be seen comparing the NMCs
325 included in Tables 1 and 2), permutation testing provided a straightforward way to assess
326 up to which extent the observed improvements in the predictive properties of PLSDA
327 models after variable selection could be attributed to chance, and to compare different
328 variable selection methods or conditions.

329

330

331 **4. Conclusions**

332 The elimination of variables irrelevant for classification is an important task that improves
333 the predictive capabilities of multivariate models and facilitates their interpretation. Still, if
334 the effect of chance correlations is unknown, variable selection must be performed in
335 combination with an assessment of the obtained PLSDA models. Using simulated data sets as
336 well as a real data set it could be shown that the inclusion of variable selection in the
337 statistical validation process provides an estimation of its statistical significance, being
338 useful when no external validation set is available. This procedure increases confidence in
339 the variable selection process, which might be relevant for biological interpretation and
340 development of further analysis methods (i.e. development of target methods) based on the
341 obtained results. Furthermore, in spite of being computing intensive, this approach can also
342 be useful to compare variable selection methods or conditions.

343

344 **Acknowledgements**

345 JE and JK acknowledge the “Sara Borrell” grant (CD11/00154 and CD12/00667) from the
346 Instituto Carlos III (Ministry of Economy and Competitiveness). DPG acknowledge the “V
347 Segles” grant provided by the University of Valencia to carry out this study. MV
348 acknowledges the FISPI11/0313 grant from the Instituto Carlos III (Ministry of Economy
349 and Competitiveness). AF acknowledges the DPI2011-28112-C04-02 grant from Spanish
350 Ministry of Science and Innovation (MICINN). GQ acknowledges financial support from the
351 Spanish Ministry of Economy and Competitiveness (SAF2012-39948).

352

354 **References**

355 [1] J.C. Lindon, J.K. Nicholson, E. Holmes, *The handbook of metabonomics and metabolomics*,
356 first ed., Elsevier, Amsterdam, 2007.

357 [2] K. Wongravee, G.R. Lloyd, J. Hall, M.E. Holmboe, M.L. Schaefer, R.R. Reed, J. Trevejo,
358 R.G. Brereton, *Metabolomics* 5 (2009) 387-406.

359 [3] S. Smit, M.J. van Breemen, H.C.J. Hoefsloot, A.K. Smilde, J. Aerts, C.G. de Koster, *Anal.*
360 *Chim. Acta* 592 (2007) 210-217.

361 [4] R. Brereton, *Chemometric pattern recognition*, first ed., Elsevier, Amsterdam, 2009.

362 [5] B. Simonetti, A. Lucadamo, M.R. González Rodríguez, *Curr Anal Chem* 8(2) 2012 266-272.

363 [6] S. Wold, M. Sjostrom, L. Eriksson, *Chemometr. Intell. Lab.* 58 (2001) 109-130.

364 [7] S. Smit, *Statistical data processing in clinical proteomics*, PhD Dissertation, University of
365 Amsterdam, Amsterdam, 2009.

366 [8] K.H. Esbensen, P. Geladi, *J. Chemometrics* 24 (2010) 168-187.

367 [9] J.A. Westerhuis, H.C.J. Hoefsloot, S. Smit, D.J. Vis, A.K. Smilde, E.J.J. van Velzen, J.P.M.
368 van Duijnhoven, F.A. van Dorsten, *Metabolomics* 4 (2008) 81-89.

369 [10] C.M. Rubingh, S. Bijlsma, E. Derks, I. Bobeldijk, E.R. Verheij, S. Kochhar, A.K. Smilde,
370 *Metabolomics* 2 (2006) 53-61.

371 [11] P. Filzmoser, B. Liebmann, K. Varmuza, *J. Chemometrics* 23 (2009) 160-171.

372 [12] L. Gidskeaug, E. Anderssen, B.K. Alsberg, *Chemometr. Intell. Lab.* 93 (2008) 1-10.

373 [13] T.A. Knijnenburg, L.F.A Wessels, M.J.T. Reinders, I. Shmulevich, *Bioinformatics* 25 (2009)
374 I161-I168.

375 [14] G. Quintás, N. Portillo, J.C. García, J.V. Castell, A. Ferrer, A. Lahoz, *Metabolomics* 8(1)
376 (2012) 86-98.

377 [15] J.G. Topliss, R.J.J. Costello, *J. Med. Chem.* 15 (1972) 1066-1076.

378 [16] K. Baumann, N.J. Stiefl, *J. Comput. Aid. Mol. Des.* 18 (2004) 549-562.

- 379 [17] K. Baumann, *Qsar Comb. Sci.* 24 (2005) 1033-1046.
- 380 [18] K. Baumann, *Abstr. Paper Am. Chem. Soc.* 227 (2004) U1026-U1027.
- 381 [19] J.A. Westerhuis, E.J.J. van Velzen, H.C.J. Hoefsloot, A.K. Smilde, *Metabolomics* 4 (2008)
382 293-296.
- 383 [20] Matlab R2011a, Mathworks, Natick, MA, 2011.
- 384 [21] V. Centner, D.L. Massart, O.E. deNoord, S. de Jong, B.M. Vandeginste, M.C. Sterna, *Anal.*
385 *Chem.* 68 (1996) 3851-3858.
- 386 [22] M.M. Hendriks, S. Smit, W.L. Akkermans, T.H. Reijmers, P.H. Eilers, H.C. Hoefsloot, C.M.
387 Rubingh, C.G. de Koster, J.M. Aerts, A.K. Smilde, *Proteomics* 7(20) 3673-3680.
- 388 [23] J. Moros, J. Kuligowski, G. Quintás, S. Garrigues, M. de la Guardia, *Anal. Chim. Acta* 630
389 (2008) 150-160.
- 390 [24] M. Daszykowski, I. Stanimirova, B. Walczak, F. Daeyaert, M.R. de Jonge, J. Heeres, L.M.H.
391 Koymans, P.J. Lewi, H.M. Vinkers, P.A. Janssen, D.L. Massart, *Talanta* 68 (2005) 54-60.
- 392 [25] V. Choulakian, M.A. Stephens, *Technometrics* 43 (2001) 478-484.
- 393
- 394

395 **Legends of Figures**

396

397 **Figure 1.** Sample classification according to calculated y values of simulated data matrices.

398 Note: blue circles: class A samples; red circles: class B samples; dotted line: class threshold.

399

400 **Figure 2.** Histograms of the number of selected variables (a), misclassified (NMC) samples

401 (b) and discriminant Q^2 (c) in the simulated **Null** (40 x 540), **SIMUIN** and **Gaucher** data sets

402 after variable selection using permuted class labels and the b_{cv} -PLSDA approach. Colored

403 dots indicate values obtained using the original class labels.

404

405 **Figure 3.** Histograms of the number of selected variables (a), misclassified (NMC) samples

406 (b) and discriminant Q^2 (c) in the simulated **Null** (40 x 540), **SIMUIN** and **Gaucher** data sets

407 after variable selection using permuted class labels and the UVE-PLSDA approach. Colored

408 dots indicate values obtained using the original class labels.

409

410 **Legends of Tables**

411

412 **Table 1.** Figures of merit of PLSDA models established by 2CV and calculated for different
413 data sets before and after variable selection. Standard deviations were obtained from 4-fold
414 2CV results (see text for details).

415

416 **Table 2.** Number of misclassified (NMC) samples included in the test sets before and after
417 variable selection.