

1 **Original article for Food Chemistry**

2

3 **Application of sequential and orthogonalised-partial least squares (SO-**
4 **PLS) regression to predict sensory properties of Cabernet Sauvignon wines**
5 **from grape chemical composition**

6

7 Jun Niimi^{†‡*}, Oliver Tomic[¶], Tormod Næs[§], David W. Jeffery[†], Susan E. P. Bastian[†], Paul K.
8 Boss[‡]

9 [†]*School of Agriculture, Food and Wine, The University of Adelaide, PMB 1, Glen Osmond,*

10 *SA 5064, Australia*

11 [‡]*CSIRO - Agriculture and Food, PMB 2, Glen Osmond, SA 5064, Australia*

12 [¶]*Faculty of Science and Technology, Norwegian University of Life Sciences, Ås, 1432*
13 *Norway*

14 [§]*Nofima – Norwegian Institute of Food, Fishery and Aquaculture, Ås, 1431, Norway*

15

16 **Corresponding author*

17 **CONTACT DETAILS:** jun.niimi@adelaide.edu.au, +61 8 8313 0284

18

19 **Abstract**

20 The current study determined the applicability of sequential and orthogonalised-
21 partial least squares (SO-PLS) regression to relate Cabernet Sauvignon grape chemical
22 composition to the sensory perception of the corresponding wines. Grape samples (n = 25)
23 were harvested at a similar maturity and vinified identically in 2013. Twelve measures using
24 various (bio)chemical methods were made on grapes. Wines were evaluated using descriptive
25 analysis with a trained panel (n = 10) for sensory profiling. Data was analysed globally using
26 SO-PLS for the entire sensory profiles (SO-PLS2), as well as for single sensory attributes
27 (SO-PLS1). SO-PLS1 models were superior in validated explained variances than SO-PLS2.
28 SO-PLS provided a structured approach in the selection of predictor chemical data sets that
29 best contributed to the correlation of important sensory attributes. This new approach

30 presents great potential for application in other explorative metabolomics studies of food and
31 beverages to address factors such as quality and regional influences.

32 **Keywords:** Multi-block data analysis; data orthogonalisation; grape; wine; sensory

33

34 **1.0 Introduction**

35 The field of metabolomics analysis is rapidly expanding in the quest to improve the
36 holistic understanding of food and beverage composition in relation to nutrition, quality,
37 safety, and authenticity (Wishart, 2008). Given that consumers are an important beneficiary
38 of any improved knowledge of processes and practices, researchers continue to search for the
39 compositional factors that contribute to flavour perception of products such as wine, which in
40 turn influence consumer behaviour. Detailed studies of this nature generate substantial
41 volumes of multiple data sets, which require suitable methods for data analysis to draw
42 conclusions about complex natural phenomenon.

43 Fortunately, the field of chemometrics provides a range of multivariate statistical
44 methods available for explorative analysis, interpretation and prediction. Chemometric
45 methods also need to keep pace with the emerging trend involving the collection of multiple
46 data sets obtained from advanced instrumental technologies with enhanced measurement
47 resolution (and are therefore very information-rich). However, the large amounts of data that
48 are generated pose a major challenge in the subsequent analysis to be able to interpret their
49 meaning (Johnson, Ivanisevic, Benton, & Siuzdak, 2015).

50 Recent developments in partial least squares (PLS)-based analyses have involved
51 extensions to multiple input data by using the PLS algorithm, including sequential and
52 orthogonalised-PLS (SO-PLS) and parallel orthogonalised-PLS (PO-PLS) (Næs, Tomic,
53 Mevik, & Martens, 2011). These techniques orthogonalise multiple data blocks, which
54 maintain the integrity of each block and can account for their respective variation to the
55 overall model. This may represent an advantage over joining multiple data sets together into
56 one large table (concatenation) for analysis with PLS (multi-block-PLS) (Westerhuis, Kourti,
57 & MacGregor, 1998). Specifically for the SO-PLS, data blocks of independent variables (X_1 ,
58 X_2 , etc.) are sequentially added to the analysis one at a time in succession to determine the
59 progressive change in explained variance of the global response data (Y) (Næs, Tomic,
60 Mevik, & Martens, 2011). Such new data analysis methods give a structured approach to the

61 analysis of highly complex data, thereby providing the best chance of properly modeling the
62 phenomenon being studied.

63 Despite being prime tools for investigating the results of something as complex as
64 human perception of foodstuffs, multi-block data analyses have been applied in only a few
65 studies to describe sensory perception and consumer preferences. Perhaps this is mainly due
66 to absence of such methods in standard commercial data analysis software that provide their
67 access through graphical user interfaces. Multi-block-PLS analysis was applied to the study
68 of aroma perception and release in cheeses in an attempt to determine the reasons behind
69 large inter-individual differences in aroma release (Feron, Ayed, Qannari, Courcoux, Laboure,
70 & Guichard, 2014). In this case, multiple predictor data sets were pre-processed prior to
71 concatenation, the PLS algorithm applied, followed by multi-block redundancy analysis.
72 Without taking the multi-block-PLS approach, mastication parameters as well as bolus
73 characteristics of cheese would not have been identified as influencing aroma compound
74 release in the mouth. PO-PLS was compared against conventional PLS for consumer
75 preference mapping of flavoured water and jams (Måge, Menichelli, & Næs, 2012). Although
76 the two analyses did not change in the output of the model such as the overall explained
77 variance, the PO-PLS approach added further information in the contribution of predictor
78 data blocks to the overall model and the number of components required for each block.
79 Lastly, SO-PLS was applied to path modeling to determine how consumer demographics,
80 purchase behaviour, and neophobia influenced each other (Menichelli, Almoy, Tomic, Olsen,
81 & Naes, 2014). Many other applications of SO-PLS or PO-PLS to food and beverage
82 research can be envisaged when diverse data sets need to be modelled (e.g., chemical, human
83 sensory, biochemical), such as the case of working with grape and wine data.

84 Given the size and economic importance of the global wine industry, understanding
85 the drivers of grape quality and how this translates into a finished wine that consumers
86 appreciate is a fundamental requirement for winemakers. Although much research has
87 elucidated grape chemical compounds that dictate some unique flavour characters in wines,
88 there is still much to be understood concerning why wine taste the way they do. For instance,
89 it is unknown why different wines from the same cultivar (i.e., cultivated variety, such as
90 Cabernet Sauvignon) possess different sensory characteristics, as a result of the complicated
91 interdependencies that occur between the chemical constituents within the grape berry,
92 overlaid with the effects of microbial metabolism during winemaking. Chemometric
93 treatment of data, in particular using multi-block data analysis methods, has the potential to

94 provide improved understanding of the grape chemical measures that best contribute to the
95 variation in wine style as determined from instrumental and/or sensory profiles of the wines.

96 The objective of this study was to explore the applicability of SO-PLS to model the
97 sensory characteristics of Cabernet Sauvignon wines as determined by human assessors using
98 a suite of chemical measurements made on the grapes, and hence to determine the data blocks
99 that most contribute to the models. In particular, focus will be on strategies for incorporating
100 blocks of data when the number by far exceeds the standard size of 2-3 input blocks modelled
101 by the approach. Focus will be on interpretation as well as prediction ability and how to
102 assess reliability of the interpretations.

103

104 **2.0 Materials and methods**

105 *2.1 Sampling*

106 To encompass a range of compositional differences, Cabernet Sauvignon grapes were
107 sampled from different vineyards within the following eight viticultural regions across South
108 Australia during the 2013 vintage (number of vineyards sampled given in parentheses):
109 Barossa Valley (2) (BV), Clare Valley (2) (CV), Coonawarra (4) (CWA), Eden Valley (2)
110 (EV), Langhorne Creek (2) (LC), McLaren Vale (2) (McL), Riverland (9) (RVL), and
111 Wrattenbully (2) (WBY). This provided a sample set of 25, with the Riverland having more
112 vineyards sampled due to its comparatively larger sampling area.

113 In each vineyard site, grape bunches were randomly sampled throughout the vineyard
114 block and were picked randomly from all areas of the canopy (Calderon-Orellana, Mercenaro,
115 Shackel, Willits, & Matthews, 2014) to give a total parcel size of 60 kg. Two subsamples of
116 grape berries (500 g each) were randomly taken from each sample parcel; one subsample was
117 used for wet chemistry (see supplementary experimental section) and the other was snap
118 frozen with liquid nitrogen and stored at -80 °C for compositional analysis. The pulp and skin
119 of frozen grape samples were later separated from the seeds, and the frozen pulp and skin
120 were homogenized, all the while under liquid nitrogen, in preparation for the various analyses.
121 The remaining fresh grape parcels were processed into single wines, corresponding to each
122 grape sample, and vinified identically using a small scale procedure (Niimi, Boss, Jeffery, &
123 Bastian, 2017).

124 2.2 Chemical and sensory analyses

125 Twelve physical and chemical analyses were made on the grape samples (Table 1).
126 For brevity, details of each analytical method are described in the Supplemental Information
127 sections. Many of the methods have been described previously in literature and the reader is
128 referred to those cited for the following measures of grapes; harvest measures, amino acids,
129 non-targeted volatile compounds, total phenolics and total tannins, anthocyanins, detailed
130 tannins, flavonols, and lipoxygenase enzyme activity pathway. Details for the remaining
131 blocks of data including certain aspects of harvest measures, targeted and bound volatile
132 compounds, colour, and fatty acid analyses are described in S-1 to S-5. The sensory analysis
133 procedure was based on previous literature, however specific details pertaining to the current
134 study are described in S-6. Sensory evaluations were conducted under the approval of The
135 University of Adelaide's Human Research Ethics Committee (H-2014-057).

136 2.3 Data Analysis

137 2.3.1 Data pre-processing

138 The following pre-processing steps were used in the current study: (I) computing
139 descriptive statistics for the variables to inspect the distribution of the data and identify
140 possible severe outlier data points; and (II) analysing the data of each variable (both X- and
141 Y-blocks) using one-way analysis of variance (ANOVA) to identify variables that
142 significantly discriminated the wine samples ($p < 0.05$). Non-significant variables were
143 removed from further analysis in an attempt to remove noise and ease the burden of
144 computing the SO-PLS models (described further below); (III) analysis of each block with
145 principal component analysis (PCA) to get an overview of the systematic variance in each
146 block and how the variables contributed to the variance as well as visualize any outliers; (IV)
147 computation of PLS2 models with Y- and one X-block at a time (totalling 12 PLS2 models)
148 for an understanding of the predictive power of each X-block; (V) X-blocks that yielded
149 PLS2 models with less than 10% validated explained variance were left out to remove noise
150 from data. Further, single Y-variables with less than 10% validated explained variance when
151 modelled using PLS1 with any of the X-blocks were also left out in another attempt to
152 remove noise from the data.

153 Note that initially, all blocks except X_{01} consisted of 75 rows (for complete
154 information on the block dimensions see Table 1) consisting of 3 replicates for each of the 25

155 measured wines. Prior to analysis as in (III), (IV), (V), and with SO-PLS, averages were
156 computed across the 3 replicates, such that the number of rows of each block was decreased
157 to 25. The order of the 25 wine samples was identical for all blocks. Moreover, all blocks
158 were mean centred prior to analysis in (III), (IV), (V) and SO-PLS.

159 *2.3.2 Analysis of multiple blocks using SO-PLS*

160 The SO-PLS approach used has been generically described previously; for further
161 information the reader is encouraged to read the following reference (Næs, Tomic, Mevik, &
162 Martens, 2011). However, as the current paper will present two variants of SO-PLS, the
163 distinction between single Y-variable SO-PLS and global SO-PLS (which analyses the entire
164 Y-variables table) will be referred hereinafter as SO-PLS1 and SO-PLS2, respectively,
165 following the conventional naming of the two PLS methods (i.e., PLS1 and PLS2).

166 A progressive model optimization approach was chosen, where X-blocks were added
167 sequentially to the SO-PLS model (Menichelli, Almoy, Tomic, Olsen, & Naes, 2014) until
168 there was no further reduction in root mean square error of cross validation (RMSECV) for
169 Y-block. This is less likely to lead to over-fitting of the data, than finding the most optimal
170 combination of X-blocks by model parameters using the global modeling approach
171 (Menichelli, Almoy, Tomic, Olsen, & Naes, 2014). The reason for this is that the more
172 possibilities there are, the higher the likelihood of obtaining a good result by chance. The
173 progressive model optimization approach is described in detail below.

174 As a first step, the Y-block was fitted to each of the available X-blocks by computing
175 one PLS2 model for each X-block. The models then were compared by means of lowest
176 RMSECV for block Y and the appropriate number of components. The X-block with the
177 most optimal parameters was selected to be the first X-block for subsequent SO-PLS2 models
178 that will be computed in the following steps. The appropriate number of components found
179 for the most optimal X-block at this point was set and held constant for the SO-PLS2 models.

180 As a second step, SO-PLS2 models with two X-blocks were calibrated. In each of
181 these SO-PLS2 models the first X block and its number of components were fixed (as found
182 in step 1 above) and the second X-block was one of the remaining X-blocks after step 1. This
183 step determined the second X-block that most decreased RMSECV in Y-block, where the 2
184 block SO-PLS2 model with the lowest RMSECV was considered the most optimal model.
185 The optimal number of components associated with the best second X-block as well as the

186 improvement in validated explained variance were also determined. Thus, the combination of
187 decreased RMSECV, number of components, and the improved validated explained variance
188 from the initial PLS2 were used as indicators that the 2 block SO-PLS2 model was better at
189 describing the variation in block Y.

190 Upon determining an improved model with 2 block SO-PLS2 model, the number of
191 components for the second block was set and held constant, as with the first X block. The
192 process was repeated by adding a third X-block from the remaining X-blocks, making a 3
193 block SO-PLS2 model (still the number of components is fixed in the first two). For the data
194 set below, the third block did not improve predictions and therefore only two blocks were
195 considered throughout the study. It should be noted that although the present approach is
196 more conservative than other selection strategies, prediction ability measures should be
197 validated further with data from harvests of following years.

198 Once the final model was identified, the explained variance of Y was computed after
199 the sequential addition of each X-block to show their progressive contribution to the total
200 explained variance in Y. Moreover, cross validation-ANOVA (CV-ANOVA) (Indahl & Naes,
201 1998) was used to test whether incorporating the X-blocks is statistically significant. This
202 method is based on comparing squared cross-validated residuals for different models using
203 paired t-tests. More specifically, the cross-validated residuals are calculated for zero, one and
204 two blocks in the model. One-block residual is compared with zero-block (using only the
205 mean of Y as an estimate) residual and two-block residual with one-block residual. When
206 more than two output variables are used, the sum of the cross-validated residuals is used in
207 the t-test.

208 As the last step of the SO-PLS2 procedure, Y was predicted from the optimal model
209 and subsequently analysed using principal components of prediction (PCP) (Langsrud & Naes,
210 2003). PCP implies that a PCA is first run for the predicted Y-values before these predicted
211 Y-values are related to all the X-variables using regression analysis. The X-blocks are then
212 standardized, put together in one block and then regressed onto the principal components of
213 predicted Y. Scores and loading plots for both Y and combinations of X-blocks were then
214 plotted for interpretation of the results.

215 The above progressive SO-PLS2 modeling procedure was further performed with
216 SO-PLS1, where one Y-variable was modelled at a time with the purpose to investigate
217 which variables in the X-blocks were responsible for high or low intensities of that specific

218 Y-variable, in this case a particular sensory attribute. As above with SO-PLS2 this approach
219 was used to reduce the heavy computational burden as much as possible, as well as reducing
220 chances of serious overfitting, as described above.

221 All analyses were performed using the Python programming language (Python
222 version 3.5) utilizing the Python packages *numpy* (Pérez & Granger, 2007), *IPython*,
223 (Oliphant, 2007), *pandas* (McKinney, 2010), and *statsmodels* (Seabold & Perktold, 2010).
224 The Python implementation of SO-PLS was coded in-house.

225

226 **3.0 RESULTS AND DISCUSSION.**

227 *3.1 Data pre-processing checks*

228 Data were initially checked with descriptive statistics to determine the distribution of
229 data in all blocks and from it, unusual distributions were not detected. Each data block was
230 analysed using one-way ANOVA. Data blocks X_{02} , X_{04} , X_{05} , X_{09} , X_{11} , and Y contained non-
231 significant variables and were further reduced to 24, 26, 51, 9, 15, and 16, respectively (see
232 Table 1 for initial number of variables). The number of variables in the remaining data blocks
233 were unchanged, as differences across samples within each block were significant for all the
234 variables (one-way ANOVA, $p < 0.05$). After removal of non-significant variables all data
235 blocks were analysed using PCA and systematic variation was investigated using scores and
236 loadings plots (not shown). From this approach outliers were not detected in any of the blocks
237 analysed.

238 *3.2 Global model of multiple Y-variables using SO-PLS2*

239 Data was initially analysed with PLS2 (see point IV in the data check and
240 preparation described above) prior to SO-PLS2 modelling. Individual PLS2 models were
241 determined for every X-Y block combination and found that two out of twelve individual X-
242 blocks gave low predictive power below 10% validated explained variance (4.1% for X_{03} and
243 a collapsed model with negative validated explained variance for X_{12}). The remaining ten X-
244 blocks were therefore used for subsequent SO-PLS2 modelling.

245 To determine how the chemical measures from the series of ten X-blocks related to
246 the Y-block (descriptive sensory analysis), the data was modelled using SO-PLS2. With so
247 many X-blocks at hand and no intuitive ordering thereof, one could have computed a vast

248 number of SO-PLS models, considering that models could be based on: (I) different
249 combinations of only two X-blocks up to as many as ten X-blocks and (II) different order of
250 X-blocks. The simplest SO-PLS model would consist of only the Y-block and any two X-
251 blocks in any order, i.e. models with both X_{01} and X_{02} , X_{01} and X_{04} , X_{01} and X_{05} and so on
252 (note that block X_{03} was left out due to low predictive power). This alone results in $\frac{n!}{(n-r)!} =$
253 $\frac{10!}{(10-2)!} = 90$ different SO-PLS models with $r = 2$ X-blocks out of $n = 10$ X-blocks to choose
254 from. Addition of blocks, so $r = 3$ X-blocks or 4 X-blocks, rapidly increases the number of
255 models to additional 720 and 5040, respectively. This clearly illustrates that the vast number
256 of models to be computed would go beyond the practicalities of time and computational
257 power, as well as in addition an enormous chance of overfitting. Given this situation, the
258 strategy of progressive modeling approach was used for selection and ordering of the X-
259 blocks to find a more robust SO-PLS model to describe the variation in the Y-block, in an
260 efficient manner.

261 The most optimal 2 block SO-PLS2 model consisted of blocks X_{06} (CIELab colour
262 measures using 1 component) followed by orthogonalisation of X_{02} (amino acids using 2
263 components) with respect to the one component from X_{06} (Fig 1.). This resulted in an
264 RMSECV of 1.103, and calibrated and validated explained variances of 59.1 % and 43.2 %,
265 respectively. Additional orthogonalisation of a third X block to compute 3 block SO-PLS2
266 models did not improve the model further based on RMSECV, number of components, and
267 increased validated explained variance criteria, thus data will only be interpreted up to 2
268 block SO-PLS2.

269 The CV-ANOVA gave p -values equal to 0.06 and 0.17 for the inclusion of X_{06} and
270 X_{02} , respectively. This may indicate that although there is a clear improvement in RMSECV
271 in both cases they are not strictly significant. Not knowing the power of the CV-ANOVA test,
272 this result only tells us that one should be careful and not over-interpret the findings as
273 reported below.

274 The projection of scores from the most optimal 2 block SO-PLS2 model can be
275 visualized on the principal components of prediction (PCP) plots (Figure 2). The explained
276 variance in the PCP plots show that the first and second PCs accounted for 75.2 % and
277 18.8 % of the variation, respectively, in Y-block that was predicted from the 2 block SO-
278 PLS2 model with X_{06} and X_{02} using 1 and 2 components, respectively.

279 The projection of scores showed relatively close groupings of samples by harvest
280 origin; RVL samples were mostly grouped in the area of both negative PC 1 and 2.
281 Specifically, the eight RVL samples, the EV samples and BV2 projected negatively along PC
282 1 had little association with the majority of the Y loadings (Fig 2). The CWA, LC, McL, BV1,
283 and CV1 samples were projected on positive PC 1, and were characterized by higher values
284 for many of the sensory attribute loadings, as well as three variables from X₀₆ from colour
285 measures; hue angle, Chroma, A* (Chroma and A* overlapped with each other), and amino
286 acids from X₀₂; methionine, isoleucine and γ -aminobutyric acid (GABA) (Fig 2). Although it
287 seems that the RVL samples were negatively projected on PC 1 of the scores plot with many
288 of the Y and X loadings, the samples had high values of two variables from the CIELab
289 measures (L* and b* that denote for lightness and yellowness, respectively). The RVL region
290 is known for their hot weather and high growing degree days, which can hinder production of
291 anthocyanin and phenolic compounds. The consequence is an influence on wine colour
292 (Ojeda, Andary, Kraeva, Carbonneau, & Deloire, 2002), which was also reflected in the 2
293 block SO-PLS2 model. The amino acids are known to contribute to wine flavour indirectly
294 by their metabolism by microorganisms and hence the resultant secondary metabolites
295 (Styger, Prior, & Bauer, 2011a, 2011b). These results show that the SO-PLS method provides
296 a strategy to select the X blocks used for modeling and some important ways of visualizing
297 the results.

298 *3.3 Modeling of single Y-variables using SO-PLS1*

299 The SO-PLS2 was taken to further detailed models for single Y-variables using SO-
300 PLS1 to determine the chemical data blocks that explain the differences in specific wine
301 sensory attributes, particularly those that are anecdotally known to be important for wine
302 quality. The SO-PLS1 procedure in block selection was performed in the same manner as
303 described above for SO-PLS2 models. Optimal combinations of data blocks and number of
304 components for each block were first determined for each single Y-variable with PLS1 prior
305 to SO-PLS1. Comparative model parameters were determined for the number of components,
306 RMSECV, and explained variances for calibration and validation (Table 2). Overall, all
307 models were determined with up to four components per data block when modelled with SO-
308 PLS1, the same limit that was set as SO-PLS2 modelling to prevent overfitting of the data. To
309 place the focus on SO-PLS1 applicability, only the sensory attributes (Y-variables) modelled
310 with two or more X-blocks are reported. Four Y-variables were modelled with only one X
311 block (PLS1), which were attributes hue, dark fruit and savoury aromas, and alcohol

312 mouthfeel. Orthogonal addition of a second X-block did not improve the SO-PLS1 models
313 for these three attributes, therefore these models are not reported. As such, a total of 12 Y-
314 variables were modelled using SO-PLS1: these included sensory attributes (number in
315 brackets) associated with colour (1), aroma (2), taste (1), flavour (4), mouthfeel (3), and
316 aftertaste (1) (Table 2).

317 The Y-variable modelled with the highest validated explained variance was savoury
318 flavour (F_Savoury, Table 2) (a negative contributor to red wine quality) (Johnson, Hasted,
319 Ristic, & Bastian, 2013). This attribute was modelled with both bound (X_{05}) and non-targeted
320 volatile (X_{04}) compound measures, suggesting that the perception of this sensory attribute
321 was driven by grape-derived volatile compounds. The RMSECV values obtained with SO-
322 PLS1 were consistently lower than the SO-PLS2 model (with the exception of depth of
323 colour), indicating that SO-PLS1 was more effective in describing the systematic variance in
324 the single Y-variables. This is because SO-PLS2, must compromise to fit X blocks to a suite
325 of Y-variables, whereas SO-PLS1 finds the optimal fit of X blocks to only one Y-variable.
326 There are no official cut-off limits with RMSECV values, meaning that the researcher must
327 choose what is acceptable based on the context of the data. The general rule, however, is that
328 lower RMSECV values are more desirable because they denote lower error margins related to
329 the means of the original input data, leading to better prediction accuracy. This must however
330 be balanced with the relative increases in explained variances from the addition of more
331 components and should there be little reduction in RMSECV, it is advisable to use simpler
332 and robust models with lower components.

333 The model for savoury flavour was interpreted by plotting separately the two PLS
334 models from the two steps in SO-PLS1 (Fig 3). Plotting of the first X block showed regional
335 separation of samples, a useful piece of information to demonstrate the chemical differences
336 by provenance that impact their sensory perception in wines. In particular, the samples from
337 LC and CWA were distinguished by higher relative intensities of both savoury flavours and
338 bound volatile compounds (2,6-dimethoxyphenol, an actinidole, 4-vinylphenol, guaiacol,
339 methyl vanillate, and benzyl alcohol, Fig 3A). The orthogonalised second X block
340 comprising non-targeted volatiles captured additional explained variance. In this instance,
341 McL2, CV2, and CWA2 were projected in the same direction as savoury flavour, along with
342 β -damascenone (a potent grape-derived volatile that enhances fruity aromas and suppresses
343 herbaceous ones), (Pineau, Barbe, Van Leeuwen, & Dubourdieu, 2007) benzaldehyde and to
344 a certain extent, 1-butanol. Meanwhile, RVL3 and RVL6–8, and LC2 had lower levels of

345 savoury flavour, along with low concentrations of benzaldehyde and β -damascenone (Fig 3B).
346 These samples had higher concentrations of (Z)-2-penten-1-ol, 2-methylbutanal, 3-
347 methylbutanal, and benzeneacetaldehyde. The advantages of SO-PLS1, specific Y-variables
348 can be predicted with multiple X-blocks that are orthogonalised with each other.

349 The 2 block SO-PLS1 models were extended to 3 block SO-PLS1 modeling. Further
350 variance could not be captured beyond two X blocks with meaningful improvement in
351 validated explained variance. This was seen as either improvement being minor in validated
352 explained variance or requirement of a high number of components, leading to over-fitting.
353 Therefore, at this point the modeling procedure was ended.

354 *3.4 Prospective of SO-PLS method in metabolomics research and considerations*

355 The SO-PLS2 modeling was able to incorporate two blocks from a total of 12 blocks.
356 The remaining ten X-blocks did not contain further additional systematic variation that would
357 lead to higher explained variances, and thus more comprehensive models. It may be that
358 modeling was limited because of several factors; the large biological transformation in the
359 sample matrix between grapes (X data) and wine (Y data), the nature of the measurements
360 (with X data being chemical and Y data being perceptual, elaborated below), and perhaps that
361 there were unknown relevant metabolites that were not measured. In spite of this, the results
362 reveal the great potential of the multi-block data analysis approach, in this case by using
363 diverse grape compositional data sets (determined instrumentally) to predict wine sensory
364 properties (as perceived by humans), thereby having profound implications for pre-
365 determining wine sensory characteristics (thus quality and style attributes) “in the vineyard”.
366 The SO-PLS approach could conceivably apply to other research fields, whether food and
367 beverage or biomedical. The current study has shown that it is possible to use up to two X-
368 blocks in SO-PLS1 models to describe a large part of the variation in single Y variables using
369 progressive model optimization, a conservative modelling approach that reduces the chances
370 of over-fitting.

371 The modeling approach taken with SO-PLS in the current study used a maximum of
372 four components for one X data block. More components can be used, thereby potentially
373 increase the validated explained variance in Y while further reducing the RMSECV, but with
374 a caveat of being wary of over-fitting the model. It is thus important for the data analyst to
375 choose the appropriate number of components suitable for the data type used in the model.
376 Unlike spectroscopic data, where the number of components can be high (Næs, Tomic,

377 Afseth, Segtnan, & Måge, 2013), the inherent noise associated with sensory measurements
378 (unavoidable inter-panellist variation) calls for a conservative approach that uses fewer
379 components in the models. The exact number may depend on the data type, however the
380 maximum number of components used in this study is most likely suitable to model other
381 sensory data, depending on the degree of increase in explained variance of the model with
382 each component.

383 Overall, the foremost advantage of SO-PLS is the ability to systematically select X
384 data blocks for analysis of the RMSECV decrease as a function of component combinations,
385 particularly when using many X data blocks. This allows for the selection of the model with
386 the largest increases in explained variances. SO-PLS1 provided substantially increased
387 explained variances for some of the attributes compared to SO-PLS2, which most likely
388 stems from the enhanced ability to match the best fitting X data blocks to each single Y-
389 variables rather than to a whole Y data set. In our case, future steps involve applying the SO-
390 PLS1 method to multiple grape data sets to explore their correlations with sensory perception
391 of wines across multiple vintages, and to determine the consistency in the contribution of
392 grape measures to the modeling of wine sensory attributes. The SO-PLS method has great
393 potential for application in any field that requires prediction of Y data from multiple X blocks,
394 irrespective of research field.

395 Several limitations are worth considering from the current study. The first is in
396 measurements used for modeling and may include collection of the chemical data, although a
397 good degree of accuracy can be expected from modern analytical instrumentation. It is
398 mainly the perceptual data, having inherent variation due to the nature of using human
399 assessors, such that a certain margin of error in the models is unavoidable thus making
400 predictions extremely challenging. Secondly the choice of the data blocks used in the
401 modeling requires scrutiny from the data analyst to decide whether the optimal X-blocks in
402 the models make sense in the context of the research field. For example, the attribute
403 astringency was best predicted by X₀₆ and X₀₈, which is sensible given that these were data
404 sets for colour and anthocyanins, respectively. Pigments contributing to red wine colour
405 (including anthocyanins) are among a range of polyphenolic compounds extracted from red
406 grapes during winemaking that are known to contribute directly to astringency (Brossaud,
407 Cheynier, & Noble, 2001). Should unrelated data blocks give optimal models for astringency,
408 say X₀₄ and X₀₅ (non-targeted and bound volatile compounds, respectively), a direct
409 relationship would be difficult to explain and likely be correlative than causative. Making

410 these judgements for models of attributes that have not previously been related to chemical
411 predictors on the other hand will be challenging. Lastly, despite the samples being vinified
412 identically, there will always be unavoidable variations arising from the vinification
413 procedure. This will inevitably be captured as unexplained variance and reflected in the SO-
414 PLS models.

415 It must be emphasized that the methodology (and results) presented here can be prone
416 to overfitting or over-optimism due to a number of reasons, the most important being the
417 relatively few samples available for building the calibration model as compared to the
418 number of variables and choices/selections that are made. Another reason could be that there
419 is a certain tendency of grouping according to wine region and therefore, all sub-models in
420 the full cross-validation are tested on samples, which are similar to at least one sample in the
421 training set. The fact that the least predictive Y-variables are eliminated, could also have a
422 slight impact. All this means that the prediction results reported should be validated by new
423 data. This also holds for the interpretation. The conclusion of this is that the present study
424 should be considered a feasibility study with some clear indication of how estimation and
425 model fitting can be done and what types of results that can be obtained.

426

427 **4.0 Conclusions**

428 In the age of big data and using the power of metabolomics, improved methods for
429 modeling diverse datasets and complex phenomenon are still required to reveal underlying
430 relationships that can be overlooked with typical modeling approaches. Thus SO-PLS
431 methodologies were investigated to link grape compositional measures with wine sensory
432 traits determined by human sensory assessment. Modeling of the data with SO-PLS2 showed
433 overall that two X-blocks could be modelled to fit the entire sensory profile of the wines.
434 Further modeling of single Y-variables using SO-PLS1 resulted in lower cross validation
435 error and higher explained variances. Conducting SO-PLS1 with X-data blocks
436 orthogonalised to maintain their data integrity enhanced the modeling of sensory data for
437 single Y-variables. SO-PLS1 was able to determine components that were optimal for each
438 X-data block, which together led to models that better represented the data with higher
439 explained variances than SO-PLS2. The use of SO-PLS provides a strategy for researchers to
440 tackle the issue of analysing and screening multiple data sets to achieve optimal modeling
441 with only important data blocks. The present work has demonstrated the value of the SO-PLS

442 analysis method in wine analysis and it is expected that this data analysis approach would
443 greatly assist in the advancement of metabolomics research more generally.

444

445 **Acknowledgements**

446 The authors would like to thank the industry partners CCW Co-operative Ltd, Yalumba
447 Wine Company and Treasury Wine Estates for generously allowing access to vineyards and
448 grape samples. The study was funded by Australia's grape growers and wine makers through
449 their investment body Wine Australia with matching funding from the Australian Federal
450 Government (CSP1201). Sandra Olarte-Mantilla and Trent Johnson are acknowledged for
451 assistance in data collection, and Sue Maffei and Emily Nicholson are acknowledged for their
452 assistance in collecting the grape samples and chemical analyses and the WIC Winemaking
453 service is thanked for producing the small scale wines.

454

455 **References**

- 456 Bonnet, J.-L., & Croljzet, J. (1977). Lipoxygenase from tomato fruit: partial purification and
457 study of some properties. *Journal of Food Science*, 42(3), 625-628.
- 458 Boss, P., Pearce, A., Zhao, Y., Nicholson, E., Dennis, E., & Jeffery, D. (2015). Potential
459 Grape-Derived Contributions to Volatile Ester Concentrations in Wine. *Molecules*,
460 20(5), 7845.
- 461 Böttcher, C., Boss, P. K., & Davies, C. (2012). Delaying Riesling grape berry ripening with a
462 synthetic auxin affects malic acid metabolism and sugar accumulation, and alters wine
463 sensory characters. *Functional Plant Biology*, 39(9), 745-753.
- 464 Brossaud, F., Cheynier, V., & Noble, A. C. (2001). Bitterness and astringency of grape and
465 wine polyphenols. *Australian Journal of Grape and Wine Research*, 7(1), 33-39.
- 466 Calderon-Orellana, A., Mercenaro, L., Shackel, K. A., Willits, N., & Matthews, M. A. (2014).
467 Responses of Fruit Uniformity to Deficit Irrigation and Cluster Thinning in
468 Commercial Winegrape Production. *American Journal of Enology and Viticulture*.
- 469 Downey, M. O., & Rochfort, S. (2008). Simultaneous separation by reversed-phase high-
470 performance liquid chromatography and mass spectral identification of anthocyanins
471 and flavonols in Shiraz grape skin. *Journal of Chromatography A*, 1201(1), 43-47.

472 Feron, G., Ayed, C., Qannari, E. M., Courcoux, P., Laboure, H., & Guichard, E. (2014).
473 Understanding aroma release from model cheeses by a statistical multiblock approach
474 on oral processing. *PloS One*, 9(4).

475 Hanlin, R. L., & Downey, M. O. (2009). Condensed Tannin Accumulation and Composition
476 in Skin of Shiraz and Cabernet Sauvignon Grapes during Berry Development.
477 *American Journal of Enology and Viticulture*, 60(1), 13-23.

478 Iland, P. G., Bruner, N., Edwards, G., Caloghiris, S., & Willkes, E. (2013). *Chemical analysis*
479 *of grapes and wine: Techniques and concepts*. Campbelltown, SA: Patrick Iland Wine
480 Promotions.

481 Indahl, U. G., & Naes, T. (1998). Evaluation of alternative spectral feature extraction
482 methods of textural images for multivariate modelling. *Journal of Chemometrics*,
483 12(4), 261-278.

484 Johnson, C. H., Ivanisevic, J., Benton, H. P., & Siuzdak, G. (2015). Bioinformatics: The Next
485 Frontier of Metabolomics. *Analytical Chemistry*, 87(1), 147-156.

486 Johnson, T. E., Hasted, A., Ristic, R., & Bastian, S. E. P. (2013). Multidimensional scaling
487 (MDS), cluster and descriptive analyses provide preliminary insights into australian
488 Shiraz wine regional characteristics. *Food Quality and Preference*, 29(2), 174-185.

489 Kalua, C. M., & Boss, P. K. (2009). Evolution of Volatile Compounds during the
490 Development of Cabernet Sauvignon Grapes (*Vitis vinifera* L.). *Journal of*
491 *Agricultural and Food Chemistry*, 57(9), 3818-3830.

492 Langsrud, Ø., & Næs, T. (2003). Optimised score plot by principal components of predictions.
493 *Chemometrics and Intelligent Laboratory Systems*, 68(1-2), 61-74.

494 Måge, I., Menichelli, E., & Næs, T. (2012). Preference mapping by PO-PLS: Separating
495 common and unique information in several data blocks. *Food Quality and Preference*,
496 24(1), 8-16.

497 McKinney, W. (2010). Data structures for statistical computing in python. In *Proceedings of*
498 *the 9th Python in Science Conference*, vol. 445 (pp. 51-56).

499 Menichelli, E., Almoy, T., Tomic, O., Olsen, N. V., & Naes, T. (2014). SO-PLS as an
500 exploratory tool for path modelling. *Food Quality and Preference*, 36, 122-134.

501 Mercurio, M. D., Damberg, R. G., Herderich, M. J., & Smith, P. A. (2007). High
502 Throughput Analysis of Red Wine and Grape Phenolics - Adaptation and Validation
503 of Methyl Cellulose Precipitable Tannin Assay and Modified Somers Color Assay to
504 a Rapid 96 Well Plate Format. *Journal of Agricultural and Food Chemistry*, 55(12),
505 4651-4657.

506 Næs, T., Tomic, O., Afseth, N. K., Segtnan, V., & Måge, I. (2013). Multi-block regression
507 based on combinations of orthogonalisation, PLS-regression and canonical correlation
508 analysis. *Chemometrics and Intelligent Laboratory Systems*, 124, 32-42.

509 Næs, T., Tomic, O., Mevik, B. H., & Martens, H. (2011). Path modelling by sequential PLS
510 regression. *Journal of Chemometrics*, 25(1), 28-40.

511 Niimi, J., Boss, P. K., Jeffery, D., & Bastian, S. E. P. (2017). Linking the sensory properties
512 and chemical composition of *Vitis vinifera* cv. Cabernet Sauvignon grape berries to
513 wine. *American Journal of Enology and Viticulture*, 68(3), 357-368.

514 Ojeda, H., Andary, C., Kraeva, E., Carbonneau, A., & Deloire, A. (2002). Influence of pre-
515 and postveraison water deficit on synthesis and concentration of skin phenolic
516 compounds during berry growth of *Vitis vinifera* cv. Shiraz. *American Journal of*
517 *Enology and Viticulture*, 53(4), 261-267.

518 Oliphant, T. E. (2007). Python for Scientific Computing. *Computing in Science &*
519 *Engineering*, 9(3), 10-20.

520 Pérez, F., & Granger, B. E. (2007). IPython: A System for Interactive Scientific Computing.
521 *Computing in Science & Engineering*, 9(3), 21-29.

522 Pineau, B., Barbe, J.-C., Van Leeuwen, C., & Dubourdieu, D. (2007). Which Impact for β -
523 Damascenone on Red Wines Aroma? *Journal of Agricultural and Food Chemistry*,
524 55(10), 4103-4108.

525 Sarneckis, C. J., Damberg, R. G., Jones, P., Mercurio, M., Herderich, M. J., & Smith, P. A.
526 (2006). Quantification of condensed tannins by precipitation with methyl cellulose:
527 development and validation of an optimised tool for grape and wine analysis.
528 *Australian Journal of Grape and Wine Research*, 12(1), 39-49.

529 Seabold, S., & Perktold, J. (2010). Statsmodels: Econometric and statistical modeling with
530 python. In *Proceedings of the 9th Python in Science Conference*, (pp. 57-61).

531 Styger, G., Prior, B., & Bauer, F. F. (2011a). Wine flavor and aroma. *Journal of Industrial*
532 *Microbiology and Biotechnology*, 38(9), 1145-1159.

533 Styger, G., Prior, B., & Bauer, F. F. (2011b). Wine flavor and aroma. *Journal of Industrial*
534 *Microbiology & Biotechnology*, 38(9), 1145.

535 Tesnière, C., & Verriès, C. (2000). Molecular cloning and expression of cDNAs encoding
536 alcohol dehydrogenases from *Vitis vinifera* L. during berry development. *Plant*
537 *Science*, 157(1), 77-88.

538 Vick, B. A. (1991). A spectrophotometric assay for hydroperoxide lyase. *Lipids*, 26(4), 315-
539 320.

- 540 Westerhuis, J. A., Kourti, T., & MacGregor, J. F. (1998). Analysis of multiblock and
541 hierarchical PCA and PLS models. *Journal of Chemometrics*, 12(5), 301-321.
- 542 Wishart, D. S. (2008). Metabolomics: applications to food science and nutrition research.
543 *Trends in Food Science & Technology*, 19(9), 482-493.

Table 1. Data blocks assigned for data analysis based on the measurements performed and analysis method used for each Cabernet Sauvignon grape sample ($n = 25$) arising from different vineyards and regions (with references and cross-reference to Supporting Information for additional details of the methods, where applicable).

Data block*	Measurement	Data Dimensions [†]	Analysis method	Method Reference and Supplementary Experimental Section
X ₀₁	Harvest measures	25 × 6	Weight, TSS [‡] , pH, TA [#]	(Böttcher, Boss, & Davies, 2012), S-1
X ₀₂	Amino acids	75 × 25	HPLC	(Boss, Pearce, Zhao, Nicholson, Dennis, & Jeffery, 2015)
X ₀₃	Targeted volatile compounds	75 × 12	GC-MS	S-2
X ₀₄	Non-targeted volatile compounds	75 × 27	GC-MS	(Kalua & Boss, 2009)
X ₀₅	Bound volatile compounds	75 × 62	GC-MS	S-3
X ₀₆	Color	75 × 5	CIELab tristimulus	S-4
X ₀₇	Total phenolics and tannins	75 × 3	UV spectrophotometry	(Iland, Bruner, Edwards, Caloghiris, & Willkes, 2013; Mercurio, Dambergs, Herderich, & Smith, 2007; Sarneckis, Dambergs, Jones, Mercurio, Herderich, & Smith, 2006)
X ₀₈	Anthocyanins	75 × 11	HPLC	(Downey & Rochfort, 2008)
X ₀₉	Tannins	75 × 11	HPLC	(Hanlin & Downey, 2009)
X ₁₀	Flavonols	75 × 7	HPLC	(Downey & Rochfort, 2008)
X ₁₁	Fatty acids	75 × 31	GC-MS	S-5
X ₁₂	Lipoxygenase pathway enzyme activities	75 × 3	Spectrophotometric	(Bonnet & Croljzet, 1977; Tesnière & Verriès, 2000; Vick, 1991)
Y	Sensory profiles	75 × 28	Descriptive analysis	(Niimi, Boss, Jeffery, &

* X block measurements were made on grapes and the Y block measurement was made on wines.

† X01 consists of 25 samples as the inputs were averaged prior to analysis, whereas the remaining X blocks included triplicates of 25 samples, giving a total of 75.

‡ TSS is total soluble solids.

Titratable acidity.

Table 2. Most optimal 2X block SO-PLS1 models determined for Y-variables (sensory attributes pertaining to wine quality).

Y-Variable ^a	1 st X	2 nd X	Comp ^b	RMSECV	Cal (%) ^c	Val (%) ^d
F_Savory	X ₀₅	X ₀₄	3_2	0.492	95.0	69.1
C_Depth	X ₀₇	X ₁₀	2_2	1.920	78.8	68.6
A_Overall	X ₁₀	X ₀₆	2_4	0.433	81.5	66.3
AT_Phenolic Length	X ₀₆	X ₀₂	1_3	0.551	81.6	64.6
MF_Tannin quality	X ₀₆	X ₀₉	1_3	0.647	77.5	64.1
T_Acid	X ₀₄	X ₁₁	3_4	0.475	94.5	63.4
MF_Astringency	X ₀₆	X ₀₈	1_1	0.938	64.9	56.5
A_Dried fruit	X ₀₆	X ₀₂	1_3	0.714	74.2	51.1
F_Pepper	X ₀₁	X ₀₂	2_2	0.562	76.3	50.1
F_Dried fruit	X ₁₀	X ₀₆	2_3	0.707	69.1	50.1
F_Dark fruit	X ₀₆	X ₁₁	1_1	0.948	58.6	44.2
MF_Body	X ₀₇	X ₀₂	2_2	0.661	65.9	42.8

^a C_ - color, A_ - aroma, T_ - taste, F_ - flavor, MF_ - mouthfeel, AT_ - aftertaste.

^b Component values in SO-PLS1 (e.g. 3_4) denotes for number of components in first and second block, respectively

^c Cal – calibrated explained variance

^d Val – validated explained variance

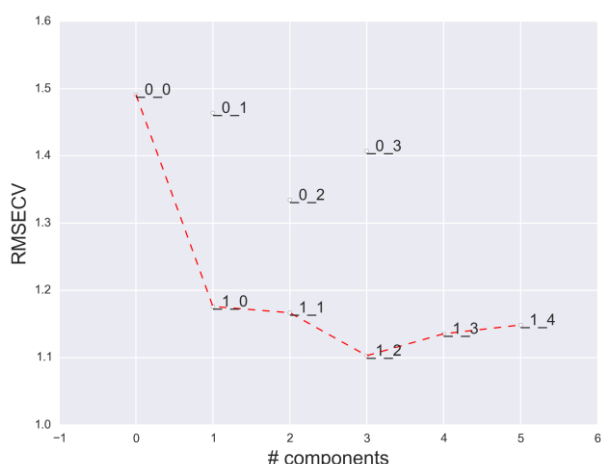


Fig 1. Måge plot showing reduction in RMSECV as a function of total number of model components for 2 block SO-PLS2 with X_{06} and X_{02} . The numbers above points with underscores denote for the number of components for each data set, i.e., components for 1st and 2nd data blocks. In this particular case, using one component for X_{06} and two components for X_{02} (hence $_1_2$) led to the lowest RMSECV.

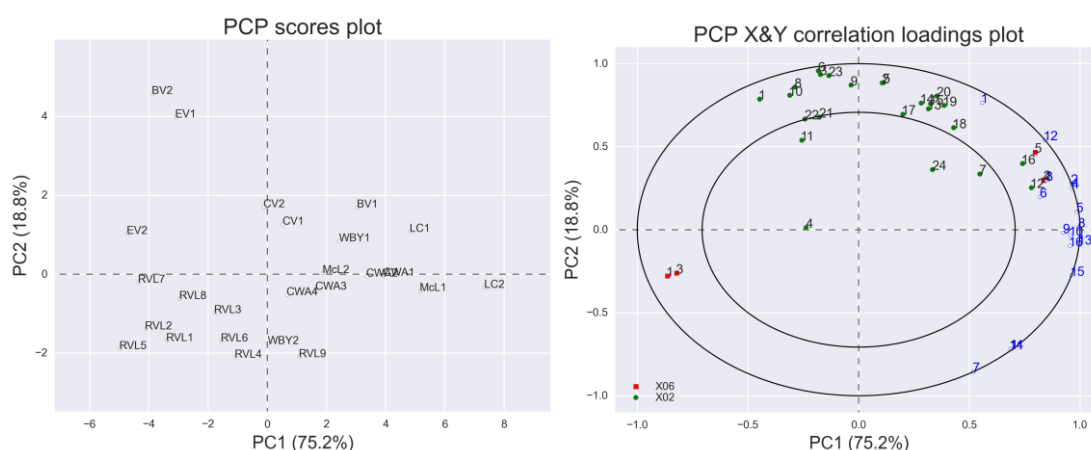


Fig 2. The PCP scores plot (left) and the X and Y correlation loadings plot (right) on the first two principal components, obtained from 2 block SO-PLS2 modeling of Y-block with X_{06} (CIELab) followed by X_{02} (amino acids) chemical data blocks. Variables in blue on the correlation loadings plot denote for Y loadings (sensory attributes) and those in green and red denote for X loadings belonging to X_{06} and X_{02} blocks, respectively. Numbers corresponding to each X and Y loadings are provided in supplementary information (Table S-4). The outer and inner ellipse on the correlations loadings plot indicate 100 % and 50 % of explained variance, respectively. Sample symbols denote for the following: BV-Barossa Valley, CV-Clare Valley, CWA-Coonawarra, EV-Eden Valley, LC-Langhorne Creek, McL-McLaren Vale, RVL-Riverland, and WBY-Wrattonbully.

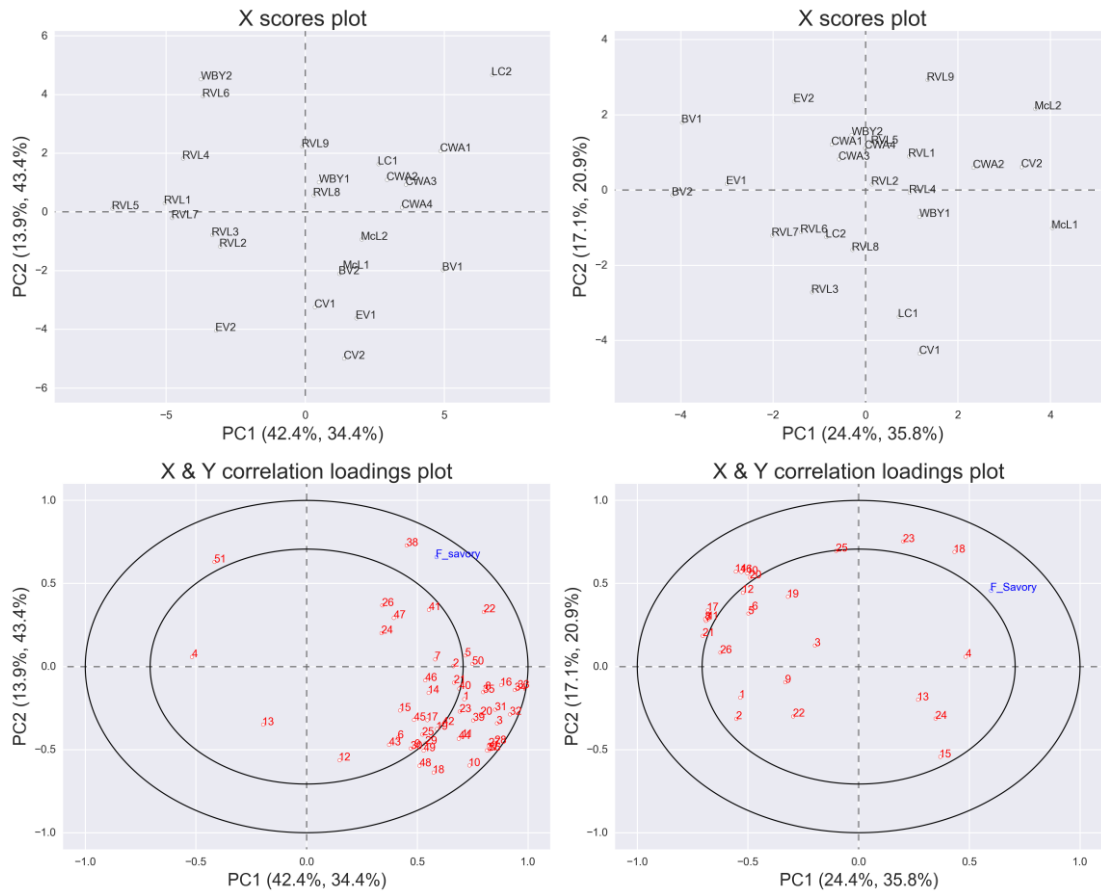


Fig 3. PLS scores and loadings plots for savory flavor attribute using 2 block SO-PLS1, showing (A) scores (top) and loadings (bottom) plots model with the first block X_{05} (left plots) (bound volatiles, S-7), and the second block X_{04} (non-targeted volatiles, Table S-5) after orthogonalization with respect to the 3 components from X_{05} (right plots). The outer and inner ellipse on the correlations loadings plot indicate 100% and 50% of explained variance, respectively. Sample symbols denote for the following: BV-Barossa Valley, CV-Clare Valley, CWA-Coonawarra, EV-Eden Valley, LC-Langhorne Creek, McL-McLaren Vale, RVL-Riverland, and WBV-Wrattonbully.

Supplementary Material

[Click here to download Supplementary Material: Supplementary Material.docx](#)