

1 **Principal components analysis of descriptive sensory data;**
2 **reflections, challenges and suggestions.**

3

4 Tormod Næs*#, Oliver Tomic+, Isabella Endrizzi** and Paula Varela*

5

6 * Nofima, Oslovegen 1, 1433 Ås, Norway

7 # Dept of Food science, Faculty of Sciences, University of Copenhagen, Rolighetsvej 30,

8 1958 Fredriksberg, Copenhagen, Denmark.

9 + Faculty of Science and Technology, Norwegian University of Life Sciences, 1433, Ås,

10 Norway

11 ** Department of Food Quality and Nutrition, Research and Innovation Centre, Fondazione

12 Edmund Mach (FEM), Via E. Mach 1, 38010 S. Michele all'Adige, Italy

13

14

15 **Abstract**

16 This paper presents a discussion of principal components analysis of descriptive sensory data.
17 Focus is on standardisation, many correlated variables, validation and the use of descriptive
18 data in preference mapping. Different ways of performing the analysis are presented and
19 discussed with focus on how to obtain informative and reliable results. The results will be
20 commented on in light of experience. All methods will be illustrated by calculations based on
21 real data. The paper ends with a list of suggestions for all the topics covered.

22 **Practical application**

23 The paper is about using PCA in sensory science. The applicability of the methods and ideas
24 presented in this paper are relevant for all types of descriptive sensory data. The ideas are
25 general and comprise areas such as standardisation, validation and many correlated variables.
26 The target group of readers for the paper is the sensory scientist who uses PCA on a daily
27 basis and who may have questions regarding how to use the method the best possible way.

28 **Key words:** QDA, PCA, validation, standardisation, partial correlation

29

30

31 **1. Introduction**

32 When analyzing data from quantitative descriptive analysis (QDA, see e.g. Stone et al.
33 (2012)), a number of choices are made more or less consciously based on tradition or habits.
34 Some of these choices, however, can have an impact on the solution, and for proper
35 interpretation of results it is important to be aware of their consequences. Special emphasis
36 here will be on the use and interpretation of results from principal components analysis
37 (PCA). Five selected aspects are described briefly below and will be discussed in more detail
38 later in the paper using examples with real data. We emphasise that this is not a exhaustive list
39 covering all possible aspects of PCA.

40 *Aspect 1: Using all individual data or aggregated data*

41 For sensory panels, data contain one intensity score value for each assessor, sample, attribute
42 and replicate. These can be analysed either simultaneously in this initial form, or one can
43 average across assessors and replicates, which is often done in practice. This results in a data
44 matrix with samples as rows and attributes as columns. In this paper we will discuss pros and
45 cons of the two approaches and point at different analysis methods that are suitable in the two
46 cases.

47 *Aspect 2: Standardisation*

48 An important first choice that has to be made when using PCA is whether the variables should
49 be used as they are in their original units or to weight/standardise them in some way. Centring
50 of variables is always done in PCA since interpretation for interval scale data is always easier
51 with a basis at the data centre than in the origin. But how to weigh the relative influence of
52 variables is less obvious.

53 A common way of making variables comparable is to standardise them to the same variance
54 (obtained by dividing the observations for each variable by its standard deviation), but in
55 many applications this is not done. It is important to stress that standardisation is not primarily
56 a statistical and technical issue, but goes to the core of how to interpret the sensory attributes
57 and to how the assessors are trained and calibrated. In other words, the variability of a sensory
58 attribute is a consequence not only of the difference of the products but also of how the panel
59 is calibrated. If the panel training is properly done, the first two principal components used for
60 visualization - with or without standardisation – will, however, usually coincide quite well if
61 non-significant variables are eliminated. In some cases other types of standardisation than the
62 standard deviation scaling, like for instance Pareto scaling (Eriksson et al. (1999)) may be
63 appropriate.

64 *Aspect 3: Many highly correlated variables*

65 Another choice that has to be made when using PCA is which variables to incorporate into the
66 analysis. Should one use all variables or only a subset reflecting the most important
67 dimensions? If for instance the same phenomenon is described by several variables, the PCA
68 plots may give a biased impression of the relative importance of the underlying sensory
69 dimensions. Obvious examples of this are variables describing the odour and flavour of the
70 same phenomenon and contrasting attributes such as dark/light and soft/hard, but other less
71 obvious examples related to the cognitive or sensing process may also be envisioned. In this
72 paper we will discuss this phenomenon in some detail and give advice regarding what to do in
73 practice. Partial correlation analysis will be proposed as a useful tool in this context. This
74 method may be useful both for making PCA results more relevant to the user and also for
75 obtaining a deeper insight that can lead to improved panel training.

76 We emphasise that there is nothing wrong with using PCA on the full data set, it will always
77 reflect the internal correlation structure in the whole data set. The potential problem is that the

78 assessment of the relative importance of underlying sensory dimensions may be biased and
79 sometimes sensory dimensions may appear more/less important than they deserve.

80 *Aspect 4: Validation*

81 Validation is another important issue when using PCA (Næs et al. (2018)). In most
82 applications of PCA one will be interested in knowing to which degree one can rely on the
83 different components extracted. One can of course always consider PCA as only an empirical
84 way of looking at the data, but some assessment of confidence in the components is also often
85 wanted. In this paper we discuss a number of ways of how this can be done. Different types of
86 validity will also be discussed.

87 *Aspect 5: QDA used in relation to consumer data*

88 In some cases, not all sensory attributes are important for the purpose they are used for. An
89 example is preference mapping, where for instance a certain spice or salt level may be
90 important for consumer preference, but its effect is blurred by the presence of a large number
91 of attributes that are irrelevant for this problem. If for instance only two principal components
92 are considered in external preference mapping, the effect of a single important variable
93 appearing in the third component may pass unnoticed. Another example is studies of satiety,
94 where in most cases only the texture attributes will be relevant (Nguyen et al. (2019)), not the
95 whole sensory profile.

96 The present paper is a discussion of these five aspects with focus on interpretation and what
97 type of effects they may have on the results. Both personal experience, concrete results from
98 sensory data and basic principles will be important in the discussion. The main purpose is to
99 provide guidelines for the sensory analyst in industry and science and suggestions of how to
100 use PCA in a safe and reliable way. The paper is not intended for the specialist statistician, but
101 for the more typical users of these methods in their daily activities and practice. Some

102 possible pitfalls are underlined and some new suggestions and tools will be presented and
103 discussed. A short introduction to PCA is provided here, but for a a thorough description of
104 several more aspects of PCA we refer to Jolliffe (2010). At the end of the paper (Section 10)
105 a number of conclusions and recommendations are given for each of the issues discussed. The
106 phenomena discussed will be illustrated by examples using real sensory data sets.

107 **2. Structure of descriptive sensory data**

108 The focus of the present paper is the use of PCA for descriptive sensory data (QDA data). In
109 most cases the entries in such data sets will lie between a lower and an upper limit on some
110 sort of intensity scale. The different attributes are calibrated to be positioned within this
111 interval. It should be mentioned that although PCA is a very important tool in this context, a
112 proper analysis and interpretation of each of the attributes separately is always recommended.
113 For the purpose of interpretation and also for some of the tools proposed, the sensory data will
114 be thought of as generated according to an experimental design with assessors and products as
115 the two factors in the design. In more technical terms, each sensory variable can be considered
116 a sum of contributions from the two factors, product and assessor, i.e.

$$117 \quad y_{ijr} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + \varepsilon_{ijr} \quad (1)$$

118 where y_{ijr} is the measurement for product i ($i=1,\dots,I$), assessor j ($j=1,\dots,J$) and replicate r
119 ($r=1,\dots,R$). The α represents the product effect, β the assessor effect, $\alpha\beta$ the interaction
120 between the two and ε represents the random error. Note that when the samples are obtained
121 according to an experimental design, one can replace the samples effect α by separate effects
122 for the design factors (see e.g. Næs et al. (2018)). It should be mentioned that for ANOVA
123 purposes, more sophisticated models than (1) have also been proposed (Brockhoff et al.
124 (2015)).

125 If we combine the models in (1) for the all sensory attributes (K), the joint model can be
126 written as

$$127 \quad \mathbf{Y} = \mathbf{XB} + \mathbf{E} \quad (2)$$

128 where \mathbf{Y} is the matrix of sensory data (each column of \mathbf{Y} represents an attribute), the \mathbf{X} is a
129 dummy matrix (containing zeros and ones) representing the design, \mathbf{B} is the matrix of
130 unknown regression coefficients and \mathbf{E} is the random error, i.e. the variation in \mathbf{Y} not
131 accounted for by the design. The different columns of \mathbf{B} represent the coefficients for the
132 different sensory variables, i.e. they correspond to the Greek letters in Equation (1). The
133 number of columns/attributes in the data matrix \mathbf{Y} is K and the number of rows will be equal
134 to $I*J*R$ (products*assessors*replicates). We refer to Figure 1a for an illustration of the data
135 structure in Equation (2). Some places below, the data set \mathbf{Y} without any prior modifications
136 or transforms will be called the raw data.

137 The data can be analysed by PCA directly using \mathbf{Y} in Equation (2) or using the data matrix
138 obtained after averaging across assessors and replicates. In this case \mathbf{Y} is sometimes referred
139 to as a consensus matrix and consists of I rows and K columns.

140 Another way of organising QDA data is by using a three-way array structure with the rows
141 corresponding to samples*replicates, columns to attributes and slices to the different assessors
142 (Figure 1b). This type of data structure can be analysed by so-called multi-way methods such
143 as PARAFAC (Bro et al. (2008)), or one of the Tucker methods (Tucker (1964)), which are
144 extensions of standard PCA. The data set organised as in Equation (2) is referred to as a three-
145 way data set which has been unfolded (See Figure 1b) vertically. The data structure to the
146 right in Figure 1b corresponds to \mathbf{Y} in Figure 1a and Equation (2). The three-way structure
147 and analysis will not be pursued further here.

148 **3. Short description of PCA.**

149 Principal component analysis is a so-called component method. This means that it is based on
 150 the idea that a large number of variables in \mathbf{Y} can be approximated by a small number of so-
 151 called components \mathbf{T} (sometimes called axes or latent variables) calculated as linear
 152 combination \mathbf{YW} , where \mathbf{W} is the matrix of so-called loading weights (columns of \mathbf{W} have
 153 length= 1). The components are found by maximising their variance and such that each new
 154 component extracted is orthogonal/uncorrelated with previous ones. The first component
 155 describes the most of the variability, the second is the next in the order etc. A consequence of
 156 the criterion used is that variables or variable groups with large variance will have a stronger
 157 impact on the solution than the rest. Usually one extracts only a few components treating the
 158 rest of the variability as noise. After calculation of the components, they can be related to \mathbf{Y}
 159 by regression in order to find the loadings \mathbf{P} . The model for PCA can be written as

$$160 \quad \mathbf{Y} = \mathbf{TP}^T + \mathbf{E} \quad (3)$$

161 Here \mathbf{T} represents the few components extracted to approximate \mathbf{Y} and the \mathbf{E} is usually
 162 thought of as noise. The \mathbf{T} 's are called scores and the \mathbf{P} 's loadings and are usually plotted in
 163 scatter plots for interpretation of results.

164 Although there is an arbitrary choice related to the scaling of \mathbf{T} relative to \mathbf{P} , one usually
 165 organises the solution such that the length of the loading vectors, columns in \mathbf{P} , is equal to 1.
 166 Then the variance of the columns of \mathbf{T} represent variability along the unit axes defined by the
 167 loadings. The components and loadings can be found using the singular value decomposition
 168 (SVD), which is a standard mathematical tool for decomposing a general matrix. For a
 169 thorough introduction to PCA we refer to Jolliffe (2010). In this paper we will consider the
 170 components in the order they appear according to explained variance and no focus will be on
 171 rotations.

172 4. PCA for original or averaged data?

173 *Averaged data for studying product differences.*

174 In most cases in the literature, panel averages are used both for interpretation and for
175 estimating relations with other data, for instance chemical data. This is a sensible strategy if
176 focus is on product differences, but should always be accompanied with proper checking of
177 the panellist quality. If an assessor is clearly outlying/different, it is questionable to keep
178 him/her as a part of the analysis. This is in particular true if the number of assessors is low
179 since in such cases outliers may have a larger impact on the analysis. A number of methods
180 have been developed for the purpose of checking panel performance (see e.g. PanelCheck
181 software, Dijksterhuis (1995), Tomic et al. (2007), Tomic et al. (2010), Dahl and Næs (2004,
182 2009)) and Dahl et al. (2008), Tomic et al. (2013)).

183 *Different types of panel averages*

184 It should be mentioned that there are different ways of obtaining panel averages (or a panel
185 consensus). One of them is to use straightforward averaging as will be focused here. Other
186 possibilities are Generalised Procrustes analysis (Gower (1975)), STATIS (see e.g. Schlich
187 (1996)), multiple factors analysis (MFA, Escofier and Pages (1995)) and various scaling
188 techniques (Romano et al (2008)). Generalised Procrustes analysis rotates, reflects and scales
189 (isotropic scaling) the individual assessor data matrices to make them as similar as possible
190 and then afterwards calculates the consensus as the average. The STATIS method calculates a
191 weighted average of the individual (cross-product) matrices, where the weights depend on the
192 RV coefficients between them. MFA concatenates the individual data matrices horizontally
193 and essentially runs a PCA on the combined matrix after a specific individual scaling of each
194 of them. The resulting scores matrix of this PCA is then used as a consensus for the individual
195 assessors. An alternative to MFA, with a similar underlying idea is the Tucker-2 method used
196 in Dahl and Næs (2009). The scaling methods in Romano et al. (2008) are used to eliminate
197 additive and multiplicative differences among assessors before averaging. Note that all these

198 methods are also suitable for investigating individual differences among assessors (See e.g.
199 Næs et al (2018)).

200 *PCA for original data*

201 If focus is also on individual differences between assessors, one can use the original **Y** data in
202 (2) directly without averaging. There will be several more points in the score plot, one score
203 for each replicate, assessor and sample combination. For improved interpretation one can
204 include colours and sample averages as will be illustrated here. This plot can be useful for
205 visualising differences/disagreement among assessors.

206 If the assessor points for each sample deviate strongly from each other, it provides evidence
207 that the assessors disagree to a larger extent. But in general, the differences will always look
208 quite large in this case due to noise and different use of the scale. For this reason, it is also
209 possible, to centre (and also standardise) each of the assessor data matrices before PCA. By
210 doing this one eliminates differences in intensity level on the scale between assessors before
211 analysis (see also Romano et al. (2008)).

212 Note that the explained variances when using the original data will normally be smaller for
213 the original data than for the averages since averaging reduces noise (see also example
214 below).

215 If focus is only on product differences, we recommend to use averaged data because of
216 simpler plots.

217 **5. Standardisation**

218 Different practices for standardisation in PCA exist, but whether to do it or not may
219 sometimes seem to be more a matter of habit than of serious reflection and consideration. The
220 issue of standardisation is important both for panel averages and for individual data.

221 For PCA in general, many different types of standardisation are used, but here we confine
222 ourselves to the most used namely division by standard deviation. It should be mentioned that
223 using PCA on standardised data is what some authors phrase as using the correlation matrix as
224 the basis for the calculation of components.

225 *Standardisation is not primarily a statistical issue*

226 It is important to emphasize that standardisation is not primarily a statistical issue. Whether to
227 do it or not is strongly related to how the sensory attributes are calibrated and interpreted. This
228 is clearly a decision with a subjective element, made by the panel leader or agreed upon by
229 the panel during the training session. One could easily envision that two panels with the same
230 sensitivity to product differences could be calibrated in a different way leading to another
231 ratio between the variability of for instance sweetness and hardness and then possibly
232 different PCA results. Culture and context will also have an influence on this matter, which
233 can lead to different plots and varying interpretation of results.

234 The complexity of the attributes will play a role (i.e. training and calibration on complex
235 attributes as for example creaminess is not straightforward), as well as the variability of
236 references. Taste and flavor attributes are usually easier to anchor with reference solutions or
237 products as compared to texture attributes.

238 A crucial question is whether one can justify that two attributes, possibly representing
239 different modalities, can be compared directly or not. Let us for instance consider two non-
240 standardised variables hardness and sweetness, the former with standard deviation equal to 1
241 and the other with standard deviation equal to 3. From this it seems that the variability of
242 hardness is 3 times larger than the variability of sweetness. The question is how to interpret
243 this in an appropriate manner. Can variability in hardness and in sweetness really be
244 compared this simply?

245

246 *Interpretation of PCA with and without standardisation*

247 If no standardisation is done, the rationale is that the ratio of the standard deviations of the
248 attributes is considered meaningful. In other words, without standardisation, one relies on the
249 meaningfulness of the subjective decisions made in the calibration phase. A consequence of
250 this is that the variables with the larger variance will have the strongest influence on the PCA
251 solution.

252 If on the other hand the variables are standardised by their standard deviation (or span or other
253 multiplicative constants), the relative differences in standard deviation are disregarded. This
254 corresponds conceptually to saying that for each of the attributes, the anchors (defining the
255 span) used for calibration of the different attributes are placed approximately at the same
256 place on the scale. This implies that differences between two samples are always interpreted
257 relative to the same variability or span. This means that variables with for instance initial
258 standard deviations equal to 1 and 3, will end up being compared as though they have the
259 same standard deviation.

260 It is important to mention that when using standardisation, the variance of all variables will be
261 the same. This implies that only the number of variables related to a sensory dimension will
262 be the driver for order of the components. If for instance one phenomenon is described using
263 four highly correlated sensory attributes and another phenomenon is represented by one
264 attribute only, the first principal component will represent the phenomenon with the four
265 attributes and the second component will represent the other variable. Therefore, in such
266 cases, importance of dimensions (in terms of explained variance) is driven by the number of
267 correlated attributes representing the same phenomenon rather than by the most dominating

268 sensory dimension. This shows that it is not obvious how to define the concept of common
269 concept of ‘most important sensory dimensions’ using QDA and PCA

270 *Eliminate non-significant attributes*

271 If one decides to standardise the data, it is important to recognise that variables with very
272 small variability will then be comparable (i.e. have the same influence) to the rest. A possible
273 problem with this is that variables containing mainly noise may become important in the
274 analysis and results. A pragmatic approach to avoid this problem is to test all attributes for
275 significant product effect, using ANOVA based on the model (1) above, or a more
276 sophisticated model as proposed in Brockhoff et al. (2015). If an attribute is non-significant,
277 the variable should be disregarded, thus reducing the amount of noise in the data. It is
278 important to emphasise that this approach should be used with care since significance of a
279 variable is not an objective concept and that significance of an attribute can be deflated due to
280 a few of the assessors only. Another aspect of eliminating non-significance variables is that
281 variables with low significance are eliminated and one is left only with variables which have
282 already proved their significance in the data. Generally, it is our view that, it is most often
283 better, from a pragmatic point of view, to remove non-significant variables in order to avoid
284 further problems with noisy attributes.

285 *Using correlation loadings plot*

286 Correlations loadings (Martens and Martens (2001)) are defined as the correlations between
287 the original variables and the components. This provides a plot similar to the standard
288 loadings plot with two axes, but is in addition most often equipped with circles indicating
289 100% and 50% explained variance. The correlations loadings have the advantage that they
290 highlight variables with low variance that may have a strong correlation with the components.

291 It is tempting to think of correlation loadings as a way of eliminating the problem of
292 standardisation. However, this is not always the case since correlation loadings only represent
293 a post processing procedure after the principal components have been estimated. The method
294 may be better at highlighting the relations between variables with a small initial variance (and
295 which therefore have little influence on the solution) and the components, but this does not
296 change the data for which PCA is calculated. For standardized data, the two are the same
297 except for a scaling factor. We here use the unit circle scaling for the correlation loadings.

298 **6. Correlations between variables**

299 A PCA solution is determined by the variance-covariance structure among all the variables in
300 **Y**. More precisely, PCA tries to explain as much as possible of the variance in **Y**. This means
301 for instance that if several variables describe the same phenomenon, this phenomenon may
302 represent more variability than the underlying phenomenon deserves, possibly only because a
303 panel leader may have chosen to have the panel evaluate these variables. To PCA it will then
304 look more important than other dimensions which may be represented only by one single
305 attribute.

306 *Avoiding highly correlated variables*

307 It is generally recommended that too much repetition of information should be avoided in
308 order to reduce unnecessary bias and focus for the PCA. Some of these repetitions may be
309 quite obvious such as using confounding attributes as for example dark/light and hard/tender
310 (see introduction), while others may be more subtle and difficult to identify directly without
311 data analysis. Assessors may for instance have problems discriminating between two or more
312 cognitively similar attributes and will automatically score them similarly. This is known as
313 halo dumping effect. It comes from the human desire of consistent cognitive structures and
314 has been widely described in the sensory literature (see for example Clark and Lawless

315 (1994)). Correlation between unrelated attributes may also happen when one salient negative
316 attribute causes another to be rated in the same direction, Such correlations are known as horn
317 effects, common when describing defective samples (Lawless and Heyman (2010)). This is an
318 unfortunate situation and having tools to detect such cognitive coincidence is important for
319 more relevant analysis and interpretation of PCA and for improved training of the panel. One
320 of the objectives of panel training is to achieve de-correlation of the attributes, and avoid
321 redundancy leading to particular issues in multi-product panels, as some attributes can be
322 correlated for one product but not for another.

323 *Correlations at different levels*

324 Correlation between attributes/columns in \mathbf{Y} can be due to correlation induced by the design
325 (\mathbf{X} in Equation (2), representing sample, assessor and interaction) and by the random error \mathbf{E}
326 in the model. The correlations between variables in \mathbf{XB} are the most important since these are
327 functions of the design of the study. Correlations among the variables in \mathbf{E} are, however,
328 conceptually more problematic. This calls for investigating the correlation structure for \mathbf{XB}
329 and \mathbf{E} separately and sometimes also for the products and assessors separately. We will next
330 discuss a possible tool to use for detecting correlations among the variables in the before we
331 describe briefly a few methods for studying \mathbf{XB} by PCA.

332 *Partial correlation for detecting correlations among random errors in equation (2)*

333 The concept of partial correlation between variables was developed for the purpose of
334 correlating two variables with each other after they have been conditioned upon a third
335 variable (or set of variables). This is equivalent to correlating the residuals \mathbf{E} for the two
336 variables with each other after they have been regressed onto the same variables. If the partial
337 correlation among two variables is high, one should consider eliminating one of them from
338 the PCA to avoid the problem discussed above. This type of information may also be

339 important for retraining the panel and to improve its performance. Since this type of
340 correlation will most typically be present at the individual level, correlation between residuals
341 at an individual level will be given the strongest focus here.

342 There are different ways of implementing this idea, but here we will confine ourselves to
343 results obtained from the residuals for all variables after a full two-way ANOVA of the data
344 (Equation (1)). The true partial correlations will be presented, but for the individual assessors
345 we will only consider correlations between the residuals from the full ANOVA of all
346 assessors.

347 *PCA for the systematic part \mathbf{XB} of equation (2)*

348 An important PCA based methods for analysing the systematic part \mathbf{XB} is ASCA (Jansen et al
349 (2005). PCA plots for this method can be used to reveal cases with highly overlapping
350 attributes as discussed above. The effects of the assessor and product (and their interactions)
351 are first estimated using the model (1) and standard ANOVA methods. Then the effects for
352 the different factors are further analysed by PCA using all the response variables. This is
353 equivalent to estimating \mathbf{B} in Equation (2), then splitting the \mathbf{XB} contribution into three parts,
354 the assessor part, product part and the interaction part. Analysing each of them by PCA results
355 in three separate PCA models. In mathematical terms this means that \mathbf{XB} is essentially written
356 as $\mathbf{X}_1\mathbf{B}_1+\mathbf{X}_2\mathbf{B}_2+\mathbf{X}_3\mathbf{B}_3$ and each of the terms is treated separately by PCA after estimation of
357 the \mathbf{B} 's. In this way information is obtained about the variability structure of the sensory
358 attributes for the assessors, products and interactions separately (see Liland et al (2018)). This
359 means that this method can reveal correlation structure at the sample level and assessor level
360 separately. The PC-ANOVA (Luciano and Næs (2009)) is related, but reverses the order of
361 ANOVA and PCA. First a PCA is run for \mathbf{Y} and then the scores for the first few components
362 are related separately to the design using the model (1).

363

364 **7. Validation of PCA models**

365 When using PCA, there is always a question of how many dimensions/components that can be
366 interpreted safely, regardless of whether it is applied to individual assessor data or panel
367 averages. PCA will always provide a model or solution, but the question is whether it is valid
368 in the sense that it is reproducible. Before considering methods for assessing validity, we will
369 discuss different types of validity.

370 **7.1. External validity.**

371 This validity looks into whether the model can tell something about a larger population of
372 samples or not. In sensory science this case is often not of highest interest since the samples
373 considered are the samples at hand and very often these are not selected to represent a larger
374 population. Typically, the samples are from product development, quality control or another
375 more specific situation and as such, the samples do not represent something else than
376 themselves and the perceptual space they span. The fact that the number of samples is often
377 also very small and sometimes based on an experimental design, makes it even more difficult
378 to interpret them as representing something bigger.

379 Leave one-out cross-validation (CV) of samples is a method which was originally developed
380 for external validation of regression models (Stone (1974)). It can also in principle be applied
381 for PCA if the explained variance of \mathbf{Y} is used as a criterion. As argued among others in Næs
382 et al. (2018), this method is for the above reasons not always suitable in PCA studies of
383 sensory data. It may give reasonable indications of number of components to rely on in
384 medium size data sets, but one should, always be careful with small data sets (for instance 4-5
385 samples) , especially if the samples were designed to be very different from each other . In the
386 results section we will give an example for a very small data set and a normally sized set.

387 For standardised data, the leave-one-out CV can be done in slightly different ways. Here we
388 have used the following procedure: every time an object is left out, the remaining data are
389 standardised prior to PCA. Then the sample which is left out is corrected for the mean and the
390 standard deviations from the samples used for model building, before calculating how well it
391 fits.

392

393 **7.2 Internal validity.**

394 Internal validity of a component means that a component is more meaningful or describes a
395 larger percentage of variance than the variance that can be obtained by chance, i.e. in data sets
396 without an underlying structure. Therefore, comparing true explained variance with what is
397 obtained by chance is a possibility. This type of validity is only referring to the data set under
398 study and will not tell anything about how well the model represents a population of other
399 samples. The cross-validation as defined by Wold (1978), which is based on successively
400 creating subsets for validation by eliminating entries according to a diagonal pattern of the
401 data set, can be considered an internal validation method. Here we will, however, concentrate
402 on a method based on permutations as proposed in Endrizzi et al. (2014) and later studied and
403 modified by Vitale et al. (2017). We will here use the original version.

404 *Permutation testing*

405 The idea behind the method is that for each new component to be tested, the residuals from
406 the model based on all previous components are permuted (for each column separately) and
407 then orthogonalised with respect to both columns and rows (since this is the case for the true
408 residuals in a PCA). Then, one calculates the explained variance of the permuted residuals
409 data set and compares it with the true explained variance. This is done by comparing the
410 explained variances for the component considered relative to the variance left in their

411 respective data sets (permuted residuals and true residuals). The procedure is repeated for a
412 large number of permutations (for instance 1000, as used here). The results are then presented
413 in a plot with component number on the X-axis and the explained variances as described
414 above on the Y-axis. For the real data, there is only one point for each component, but for the
415 permuted data, we will here present three values, the median, the lower 5% percentile and the
416 upper 5% percentile, obtained from a large number of permutations. The lower and upper
417 values are there for assessing the uncertainty of the estimates. If the true value falls clearly
418 above the confidence band obtained by the two percentiles, the component can be judged
419 significantly different from that generated by chance and therefore worth looking at. Although
420 assessing the number of components is essentially a one-sided test, we here prefer the setup
421 used to indicate the uncertainty in both directions. For details we refer to Endrizzi et al.
422 (2014).

423 *Assessor based cross-validation*

424 If original data are available at individual assessor level, another possible internal validation
425 method is to compare results for the different assessors, i.e. to cross-validate the assessors
426 instead of the samples. We here refer to the block splitting according to assessor illustrated to
427 the right in Figure 1b. A possible way of doing this is to project each assessor, i.e. each
428 segment removed, onto the space spanned by the rest of the assessors and compute the
429 average explained variance over the segments. This method can also be used to identify
430 outlying assessors by looking at the individual contributions to the explained variance.

431 **7.3 Validation using external information.**

432 In some cases, there may be other data available about the samples, for instance chemistry
433 data, spectroscopy data or simply the experimental design. In such cases it is possible to
434 regress the (for instance) average sensory attribute scores (across assessor and replicates) onto

435 the external data and then evaluate how much of the sensory data that can be accounted for by
436 the external variables/measurements. Such a method was used in Dahl and Næs (2004) for
437 relating the average sensory profile to external near infrared (NIR) spectra. Explained
438 variance of the sensory profile obtained from the NIR data was then used as criterion of
439 validity. In the paper the same was also done for each individual assessor separately in order
440 to identify outliers.

441 If PCA is run on the raw data \mathbf{Y} (equation 2), the PC-ANOVA method mentioned above can
442 also be used for validation. Each principal component for the full data set is now regressed
443 onto the design variables (product, assessor and interactions) using the model (1). Note that
444 this can be done in all possible cases with more than one replicate since the sample factor here
445 only refers to the samples tested and not necessarily to a particular experimental design for the
446 samples. It must be stressed, however, that the significance tests in such a model may be quite
447 strong tests due to the large number of observations. One should therefore in addition to
448 looking at degree of significance also look at the explained variances of the components in
449 order to evaluate relevance. A component with very small explained variance and only
450 borderline significant product factors is usually not worth focusing on too much. Significance
451 testing in this case may therefore in general be more useful for assessing the significance of
452 the first 2-3 components rather than evaluating how many components further out that are
453 significant.

454 **7.4 Validation using confidence intervals.**

455 In addition to focusing directly on the significance of a component, confidence intervals or
456 ellipsoids for each sample is a good option. They are primarily meant for assessing stability of
457 solutions, but can also be useful for indicating how many components that are worth
458 considering. Bootstrap procedures as illustrated for instance in Cadoret and Husson (2013) are
459 the most important to use in this case. The method is based on resampling assessors at random

460 (the same number as in the original panel) and calculating the scores for each selection (after
461 averaging over assessors). These are then projected onto the scores plot of the original
462 averaged PCA and confidence ellipses are drawn based on this for each sample.

463 **8. Implications for relations to consumer data**

464 As mentioned in the introduction, very often a sensory data set is not only used for
465 understanding the variability in the sensory properties of samples. A typical example is
466 preference mapping where the main focus is on relating consumer liking to sensory data. One
467 can do this by analysing one sensory attribute at a time, but a more typical way is to use PCA
468 of the sensory data (or PLS regression) and regress the liking for different consumers onto the
469 first couple of components (often only 2). If then a specific attribute with minor relation to the
470 main variability of the sensory data set, has an important influence on the liking, it will not be
471 visible in standard external preference mapping analysis with 2 components. Typical
472 examples are salt level and spices which may influence liking strongly, but don't account for
473 much variability in the sensory data. One should therefore inspect more than 2 components or
474 supplement (or replace) the analysis with an internal preference mapping, where PCA is
475 applied to the liking data and sensory data are regressed onto the these principal components.
476 PLS regression could be another alternative for such data (see e.g. Næs et al. (2018).)

477 Satiety studies is another important example where the whole sensory profile is not needed for
478 explaining consumer data. This was demonstrated in Nguyen et al (2019). In such cases, the
479 texture properties are the essential ones for relating to satiety; the rest may not add
480 information to explain the problem at hand, or can at worst blur the focus and results of the
481 study.

482 **9. Case studies**

483 **9.1 Data sets used.**

484 Table 1 shows the structure of the 3 data sets used in the different examples.

485 **9.2. Case 1. Should one average or not before computing PCA on sensory data?**

486 **Exemplified using yogurt data.**

487 The data used for visualizing the differences between using the PCA for average data and for
488 the individual data before averaging is a yoghurt dataset with 8 samples and 21 attributes,
489 (Nguyen et al. (2019)). An experimental design with 3 factors at two levels is used for
490 producing the samples. In this case we focus on standardised data for visualization (after
491 elimination of the single non-significant attribute at 5% level).

492 The results are presented for panel averages and raw data in Figure 2 and Figure 3. In Figure 3,
493 the average component scores across assessors for each sample are superimposed using
494 diamond shapes. As can be seen, the loadings are quite similar for the two PCA models, but
495 the explained variances are larger for the averaged data due to the averaging process, as
496 explained above. The main difference in loadings is that dryness in mouth and astringent form
497 an own group of attributes for the individual data while for standardised data they are grouped
498 together with sandy, stale odour, etc. There are quite large individual differences around each
499 sample average in Figure 3 (scores with same colour). Still, the average scores for each
500 sample are quite similar to the scores in Figure 2. This means that the essential information is
501 similar for the two analyses. The former provides a simpler plot, while the second gives an
502 opportunity for studying individual differences. As will be seen below, the latter also allows
503 for an ANOVA test for the components. In practice choosing between the two is often a
504 matter of scope of the study and need for simplicity. Most of the discussion below will be
505 focused on average data.

506 **9.3 Case 2. Should one standardize or not before PCA? Exemplified using olive oil data.**

507 An illustration of the effect of standardisation will be given using data from sensory analysis
508 of olive oil (based on averages over assessors). The results are presented in Figures 4a, b, c
509 and d. Figure 4a gives results from PCA on the full set of variables without standardisation,
510 while in Figure 4b, PCA is based on the full set of standardised variables, Figure 4c shows
511 results of PCA for only significant variables, not standardised, while Figure 4d shows PCA
512 results for significant standardised variables. In all cases the explained variances were high,
513 about 90% after 3 components. The three components look significant using leave-one-out
514 cross-validation, and this is also confirmed by the other permutation based method to be
515 shown below.

516 The Figure 4a shows that loadings and correlation loadings plot are quite different without
517 standardisation. The Figure 4b shows that the scores plot change significantly after
518 standardisation, but now the loadings and correlation loadings are quite similar. Correlation
519 loadings are also different in Figure 4a and Figure 4b. This means that standardisation has an
520 effect on scores and loadings if used on all variables without considering significance. Also,
521 correlation loadings may change with standardisation.

522 After eliminating non-significant variables (Figure 4c. 6 attributes eliminated), we see that the
523 scores are back again to the ones obtained without standardisation for the full set of variables
524 (Figure 4a). Correlation loadings and loadings are still different, but less so if we compare
525 with the full data set. Standardisation (Figure 4d) now has little effect (for reduced data) on
526 the loadings except for one variable close to the middle. Scores are almost the same for
527 Figure 4c and Figure 4d. After standardisation, loadings and correlation loadings in Figure 4d
528 are identical except for the scaling.

529 In conclusion. After elimination of non-significant variables, the results are similar regardless
530 of whether one standardised or not. This is true for both scores and loadings.

531 Comparing full and reduced data sets, we see that scores are almost the same except for the
532 standardised full data set (Figure 4b). Two of the attributes (acidic-O and oxidised-O) that
533 show up in the full data set along the second component are not present in Figure 4c and
534 Figure 4d since they are non-significant. They are also less visible in Figure 4a. These two are
535 examples of variables that are ‘inflated’ when standardised. This phenomenon is quite
536 frequent with off-flavours or other attributes that may appear in low intensities (i.e. spicy).
537 After standardisation low scoring attributes will get a larger importance in the outcome.
538 Our advice is to eliminate non-significant variables since it then matters less what is done
539 regarding standardisation. The standardised results with all variables, including non-
540 significant ones, are the most different from the rest. One should focus on a good training for
541 the low scoring attributes when relevant for the products or objective of the study.

542 **9.4 Case 3. Many correlated sensory variables. Exemplified using yogurt and olive oil**
543 **data.**

544 Figure 2 shows PCA results from the yogurt experiment in Nguyen et al. (2019) (based on a
545 2³ design). Most of the variables contrast each other along the first axis. This means that the
546 large variability accounted for along this axis to a large extent is due to the many variables
547 measuring more or less the same phenomenon. This is important information per se, but it
548 clearly gives a biased impression of the relative importance of the two components or
549 underlying dimensions (62% and 20%). Eliminating several of the highly correlated variables
550 along the first component, leads to a different relative weighting of the two axes. In other
551 words, the relative importance of the components is dependent on how many strongly
552 correlated variables that are in the data set.

553 In practice there is no fixed rule for how to possibly reduce the profile other than the obvious
554 ones, for instance dark/light. It is, however, important to be aware of this fact and interpret
555 results accordingly.

556 *Partial correlation results*

557 An illustration of the use of the partial correlation concept discussed above is given in Figure
558 5 for the olive oil data set, both for the whole panel (Figure 5a) and for three individual
559 assessors (presented in Figures 5b, 5c and 5d). There is some correspondence between panel
560 and individuals, but the individuals are also quite different. The panel clearly has a large
561 partial correlation between grass flavour and grass odour, between astringency and burning,
562 between astringency and bitter and between bitter and burning. The same tendency holds for
563 two of the individuals presented, but the third does not share this particular tendency. For the
564 assessor in Figure 5b, there are also many partial correlations among some of the attributes in
565 the middle of the plot, for instance between grass flavour and a number of the other attributes.
566 For this specific assessor there is good reason to question his/her interpretation of the
567 attributes involved and consider a retraining.

568 **9.5 Case 4. Validation based on cross-validation and permutation testing. Exemplified** 569 **using olive oil data**

570 Figure 6 shows results from the permutation test (a) and standard leave-one-out cross-
571 validation (b) for the olive oil data (see above for details) In the permutation test the true
572 explained variance is far outside the confidence interval for components up to 3. After that it
573 is inside, which indicates that from component 4 one cannot distinguish the component from
574 noise. Ten components is the maximum number possible and therefore no confidence interval
575 can be computed for the tenth component.

576 This data set is also quite suitable for the leave-one-out CV since there are many very similar
577 samples and no unique ones. As can be seen (based on the explained variance along the
578 vertical axis), also the CV indicates clearly that at least 3 components can be interpreted.
579 After that the improvement is negligible. The advantage of the randomisation test is that it
580 gives a statement of significance.

581 *An illustration based on reduced data*

582 For illustrating the problems with standard leave one out cross-validation for small data sets,
583 we selected a subset consisting of only 4 samples from the olive oil data and computed a new
584 PCA model based on standardized data. The scores and correlation loadings are given in
585 Figure 7a) and Figure 7b) respectively. The leave one out CV (Figure 7c) gives meaningless
586 results since each sample is unique and the model changes substantially every time one
587 sample out of four is left out during cross-validation. Note that a negative value of explained
588 variance is not possible when fitting the data by PCA, but for validation it can happen when
589 data left out (a segment or single samples) fit very poorly to the model estimated by the rest of
590 the data.

591 The permutation method (Figure 7d), on the other hand, indicates that the first component is
592 reliable, while the second is not. This means that the vertical axis has no statistical power
593 regarding interpretation. In other words, there is no general tendency (underlying common
594 component) representing common variability among samples along the second component. It
595 should be emphasized, however, that statistical properties of the permutation test for such
596 small data sets have not yet been tested out, so care must be taken not to overinterpret the
597 results. It should also be mentioned that this is a very extreme case for CV and incorporated
598 just to illustrate how problematic it can be for very small data sets.

599 An interesting observation is that the loadings plot change when a subset (oils 3, 7, 10 and 11)
600 of the full set of samples (oil 1-11) is used (see Figure 4d). This underlines that interpretation
601 of a subset of samples only relates to this specific subset at hand and cannot be generalised to
602 the sensory space of the full set of samples. Conclusions will then always be local and of
603 limited value for saying something about a larger set of ‘similar’ samples.

604 *The use of PC-ANOVA for validation*

605 PC-ANOVA (Luciano, G. and Næs, T.(2009) was applied to the standardised yogurt data and
606 compared to the use of the permutation test for the consensus/average data set. The results are
607 presented in Figure 8 and Figure 9. As can be seen, the results correspond reasonably well, the
608 first three components are obviously significant, while number 4 is more questionable. It
609 seems that the PC-ANOVA finds significance further out (components 5 and 6), but these
610 components represent so small variance that they are not very interesting in practice. Also, the
611 fact that component number 4 is non-significant is an indication that one should not consider
612 further components after component 3. The explained variances for the 5 first consensus
613 components are 64.4, 21.1, 9.5, 2.7 and 1.2. For the PCA done on raw data the corresponding
614 values are 28.2, 17.2, 10.4, 9.1 and 6.8. As can be seen, the drop in this case is smaller from
615 the first to the second component.

616 **9.6. Case 5. Relations between QDA and consumer data. Exemplified using bread data.**

617 For this example based on external preference mapping, a bread data set with 8 samples
618 (based on a 2^3 design) and 13 attributes is used. The data set consists of both QDA data and
619 consumer liking of the same samples. Only the averages will be considered for QDA.

620 In Figure 10 correlation loadings plots of component 1 vs. component 2 and for component 1
621 vs. component 3 are shown. As can be seen, there is a major tendency in liking towards

622 component 3 dominated by salt taste. This tendency is not visible in the plot of component 1
623 vs. component 2 where salt is lying well within the 50% explained variance circle.
624 This shows that relying only on a two-dimensional external preference mapping plot can leave
625 important drivers of liking undetected.

626 **10. Conclusions and suggestions**

627 *Using averages over assessors or raw data.*

628 The average data will give a simpler solution to look at, but no information about individual
629 differences across assessors in the panel. When choosing averages it is not possible to apply
630 PC-ANOVA the way presented here for deciding on the number of components. If averaging
631 is used, one should always do a proper check on the reliability of the individual assessors
632 before averaging.

633 *Standardisation*

634 The calibration and training procedure should be considered and evaluated for making a
635 decision on whether to standardise or not. The focus should be on the meaningfulness of
636 relying on actual differences in variability of different attributes (possibly belonging to
637 different sensory modalities) in the analysis. If these are not meaningful, one should
638 standardise. This is an interesting aspect when comparing results from different panels. In
639 such cases, the need for standardisation is stronger unless the training procedure is
640 harmonised between the labs. If clearly non-significant variables are present, one should be
641 careful about incorporating them in a standardised analysis.

642 *Using all attributes or eliminating obvious overlap.*

643 Eliminating highly correlated variables will in most cases have only a moderate effect on the
644 interpretation. One should be careful about strong statements about what are the most

645 important sensory dimensions since this will depend on the number of attributes that represent
646 it. A tool based on partial correlations is presented that can enhance insight into non-trivial
647 overlap among attributes.

648 *Validation of components*

649 Leave-one-out Cross-validation is often not the best choice in sensory analysis when samples
650 are unique and few.. In such cases an alternative is to use permutation testing.

651 *Relating sensory QDA data to consumer liking data*

652 In this case it is important to be aware that not all variables may be of interest. If obvious
653 candidates exist, one should consider excluding the non-informative variables. On the other
654 hand, there may be important attributes that are not so visible when considering only few
655 principal components of sensory data. It is always recommended in such cases to compute a
656 PCA model of consumer liking data to support the conclusions. Alternatively, one can take
657 the latter as point of departure and regress sensory variables individually onto the PCA
658 solution (internal preference mapping).

659

660 **Acknowledgements.**

661 We would like to thank Dr Nguyen for providing the yogurt data. The authors would like to
662 thank for financial support from Research Council of Norway.

663

664 **References**

665

666 Bro, R., Qanari, E.M., Kiers; H.A.L. , Næs, T. and Frost, M.B. (2008). Multi-way models for
667 sensory profiling data. *J. Chemometrics*, 22, 36-45..

668 Brockhoff, P.B., Schlich, P. and Skovgaard, I. (2015). Taking individual scaling differences
669 into account by analyzing profile data with the Mixed Assessor Model. *Food Quality and*
670 *Preference* 39, 156–166.

671 Cadoret, M., & Husson, F. (2013). Construction and evaluation of confidence ellipses applied
672 at sensory data. *Food Quality and Preference*, 28, 106-115.

673 Clark, C. C. and Lawless, H. T. 1994. Limiting response alternatives in time–intensity
674 scaling: An examination of the halo-dumping effect. *Chemical Senses*, 19, 583–594

675 Dahl. T. and Næs, T.(2004). Outlier and group detection in sensory analysis using hierarchical
676 clustering and the Procrustes distance. *Food Quali. Preference.*, 15 (3). 195-208.

677 Dahl, T., Tomic, O. Wold, J.P. and Næs, T. (2008). Some new tools for visualising multi-way
678 sensory data. *Food Quality and preference*, 19

679

680 Dahl, T and Næs, T. (2009). Identifying outlying assessors in sensory profiling using fuzzy
681 clustering and multi-block methodology. *Food Quali. Preference*, 20, 287-294.

682 Dijksterhuis, G. (1995). Assessing panel consonance, *Food Qual. Preference*, 6 (1), 7-14.

683 Endrizzi, I. Gasperi, F., Rødbotten, M and Næs, T. (2014). Interpretation, validation and
684 segmentation of preference mapping models. *Food Quality and Preference*. 32, 198-209.

685 Eriksson L, Johansson E, Kettaneh-Wold N, Wold S: Scaling. In *Introduction to multi- and*
686 *megavariate data analysis using projection methods (PCA & PLS)* *Umetrics*; 1999:213-225.

- 687
- 688 Escofier, B. and Pages, J. (1995). Multiple factor analysis, *Comp. Stat. Data Analysis*. 18, 121-
- 689 150.
- 690
- 691 Gower, J. C. (1975). Generalized Procrustes analysis. *Psychometrika*, 40, 33–51.
- 692
- 693 Jansen, J., van der Hoefsloot, J., Greef, M., Timmerman, E., Westerhuis, J., & Smilde, A.
- 694 K. (2005). ASCA: Analysis of multivariate data obtained from an experimental design.
- 695 *Journal of Chemometrics*, 19(9), 469–481.
- 696
- 697 Jolliffe, I.T. (2010). *Principal component analysis*. Springer.
- 698
- 699 Lawless, H.T. and Heyman, H (2010). *Sensory evaluation of food. Principles and practices*,
- 700 Springer Science and Business Media, New York.
- 701
- 702 Liland, K.H., Smilde, A. and Næs, T. (2018). Confidence ellipsoids for ASCA models based
- 703 on multivariate regression theory. *J. Chemometrics*, 32. <https://doi.org/10.1002/cem.2990>
- 704 Luciano, G. and Næs, T.(2009). Interpreting sensory data by combining principal component
- 705 analysis and analysis of variance. *Food Quality and Preference*, 20, 3, 167-175.
- 706 Martens, H. and Martens, M. (2001). *Multivariate Analysis of Quality : An Introduction*.
- 707 John Wiley and sons, Chichester.
- 708 Nguyen, Q. C., Næs, T., Almøy, T and Varela, P. (2019). Portion size selection as related to product
- 709 and consumer characteristics studied by PLS Path Modelling. *Food Quality and Preference* (in press)
- 710

- 711 Næs, T., Varela, P. and Berget, I. (2018). Analysing individual differences in sensory and
712 consumer science. Elsevier.
713
- 714 PanelCheck software: <https://sourceforge.net/projects/sensorytool/>, DOI:
715 10.5281/zenodo.10768
- 716 Romano,R., Brochoff, P.B., Hersleth, M., Tomic, O., Næs, T. (2008). Correcting for different
717 use of the scale and the need for further analysis of individual differences in sensory analysis.
718 Food Quality and Preference, 19, 197-209.
- 719 Schlich, P. (1996). Defining and validating assessor compromises about product distance and
720 attribute correlations. In Næs, T. and Risvik, E.(Eds.). Multivariate analysis of data in sensory
721 science. Elsevier, Amsterdam.
- 722 Stone, M. (1974). Cross-validatory choice and assessment of statistical prediction. J. Roy.
723 Stat. Soc., B. 111-133.
- 724 Stone, H., Bleibaum, R. and Thomas, H. (2012). Sensory evaluation in practice.4th edition.
725 Academic Press.
- 726 Tomic, T., Nilsen, A., Martens, M. and Næs, T. (2007). Visualization of sensory profiling
727 data for performance monitoring. LTW, 40, 262-269.
728
- 729 Tomic,O., Luciano, G., Nilsen, A., Hyldig, G., Lorensen, K. and Næs, T. (2010). Analysing
730 sensory panel performance in a proficiency test using the PanelCheck software. Eur. Food
731 Res. Technology, 230, 3, 497-511.
- 732 Tomic, O., Forde, C., Delahunty, C. and Næs, T. (2013). Performance indices in descriptive
733 sensory analysis – A complimentary screening tool for assessor and panel performance. Food
734 Quality and Preference, 28, 122-133

735 Tucker, L.R., (1964). The extension of factor analysis to three-dimensional matrices. In N.
736 Fredriksen and H. Gulliksen (eds). Contributions to mathematical psychology, Holt, Rinehart
737 and Winston, NY.

738 Vitale, R., Westerhuis, J.A., Næs, T., Smilde, A.K., de Noord, O.E., Ferrer, A. (2017).
739 Selecting the number of factors in Principal Component Analysis by permutation testing -
740 Numerical and practical aspects. J. Chem. 31, 12, pages ...

741

742 Wold, S. (1978). Cross-validatory estimation of the number of components in factors analysis
743 and principal component models. Technometrics, 20, 397-406.

744

<u>Data set</u>	<u>Number of samples</u>	<u>Number of attributes</u>	<u>Number of assessors</u>
<u>Yogurt</u>	8	21	9
<u>Olive oil</u>	11 and 4	20	<u>Only averages used</u>
<u>Bread</u>	8	13	<u>Only averages used</u>

745

746 **Table 1.** Overview of QDA data sets used. For the olive oil data set also the small subset is
747 tested. For the bread data also consumer liking data for a number of consumers were available

748

749 **Figure Captions**

750

751 **Figure 1a.** Illustration of the setup in Equation 2. The D now represents the number of design
752 variables (including product and assessor factors plus interactions).

753

754 **Figure 1b** Data structure for QDA presented as a three-way data set and an unfolded data set.
755 The illustration is for simplicity only for 4 assessors. If replicates are present, the vertical
756 dimension will be samples*replicates ($I \times R$)

757

758 **Figure 2,** Yogurt data. Standardised PCA on consensus data, 20 significant attributes.

759

760 **Figure 3.** Yogurt data. Scores and loadings for the standardized PCA based on individual
761 data, 20 significant attributes.

762

763 **Figure 4a,** Olive oil data. Full data set non-standardised

764

765 **Figure 4b.** Olive oil data. Full data set standardised

766

767 **Figure 4c.** Olive oil data. Reduced data set non-standardised

768

769 **Figure 4d.** Olive oil data. Reduced data set standardised

770

771 **Figure 5,** Olive oil data. Heat map of correlations between residuals for different attributes.

772 Over all assessor in a). The other three, b), c) and d), represent three individual assessors.

773

774 **Figure 6**, Olive oil data. Non-standardised PCA, 14 significant attributes. The illustration in
775 a) shows the curve obtained by the permutation method. The points represent the quantiles for
776 each of the number of components. In b) is presented explained variance for fitting/calibration
777 and leave-one-out cross-validation.

778

779 **Figure 7**. Olive oil data. Four samples, standardised PCA, 14 significant attributes. a) scores
780 and b) correlation loadings, c) cross-validation, d) permutation testing.

781

782 **Figure 8**. Yogurt data. PCA-ANOVA results, standardised PCA, 20 significant attributes. a)
783 multiple comparisons for products. Line indicates range of no significant differences. b) F-
784 values for the product effect factor. The significance is indicated with colour as given in panel
785 in the upper right corner.

786

787 **Figure 9**, Yogurt data. Standardised PCA. 20 significant attributes. Permutation test for PCA
788 based on averages over assessors.

789

790 **Figure 10**, Bread data. Non-standardised PCA. Correlation loadings for external preference
791 mapping . a) component 1 vs. Component 2. b) component 1 vs. Component 3.

792

793