



ORIGINAL ARTICLE

Principal components analysis of descriptive sensory data: Reflections, challenges, and suggestions

Tormod Næs^{1,2}  | Oliver Tomic³ | Isabella Endrizzi⁴ | Paula Varela¹ 

¹Nofima, Ås, Norway

²Department of Food Science, Faculty of Sciences, University of Copenhagen, Copenhagen, Denmark

³Faculty of Science and Technology, Norwegian University of Life Sciences, Ås, Norway

⁴Department of Food Quality and Nutrition, Research and Innovation Centre, Fondazione Edmund Mach (FEM), San Michele all'Adige, Italy

Correspondence

Tormod Næs, Nofima, Oslovegen 1, Ås 1433, Norway.

Email: tormod.naes@nofima.no

Funding information

Research Council of Norway; The Norwegian Levy on Agricultural Products

Abstract

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

Practical Application

The article is about using principal components analysis (PCA) in sensory science. The applicability of the methods and ideas presented in this article are relevant for all types of descriptive sensory data. The ideas are general and comprise areas such as standardization, validation, and many correlated variables. The target group of readers for the article is the sensory scientist who uses PCA on a daily basis and who may have questions regarding how to use the method the best possible way.

1 | INTRODUCTION

When analyzing data from quantitative descriptive analysis (QDA, see for example, Stone, Bleibaum, & Thomas, 2021), a number of choices are made more or less consciously based on tradition or habits. Some of these choices, however, can have an impact on the solution, and for proper interpretation of results it is important to be aware of their consequences. Special emphasis here will be on the use and interpretation of results from principal components analysis (PCA). Five selected aspects are described briefly below and will be discussed in more detail later in the article using examples with real data. We emphasize that, this is not an exhaustive list covering all possible aspects of PCA.

1.1 | Aspect 1: using all individual data or aggregated data

For sensory panels, data contain one intensity score value for each assessor, sample, attribute, and replicate. These can be analyzed either

simultaneously in this initial form, or one can average across assessors and replicates, which is often done in practice. This results in a data matrix with samples as rows and attributes as columns. In this article, we will discuss pros and cons of the two approaches and point at different analysis methods that are suitable in the two cases.

1.2 | Aspect 2: standardization

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

A common way of making variables comparable is to standardize them to the same variance (obtained by dividing the observations for each variable by its standard deviation), but in many applications this

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2021 The Authors. *Journal of Sensory Studies* published by Wiley Periodicals LLC.

is not done. It is important to stress that standardization is not primarily a statistical and technical issue, but goes to the core of how to interpret the sensory attributes and to how the assessors are trained and calibrated. In other words, the variability of a sensory attribute is a consequence not only of the difference of the products but also of how the panel is calibrated. If the panel training is properly done, the first two principal components used for visualization—with or without standardization—will, however, usually coincide quite well if nonsignificant variables are eliminated. In some cases other types of standardization than the standard deviation scaling, like for instance, Pareto scaling (Eriksson, Johansson, Kettaneh-Wold, & Wold, 1999) may be appropriate.

1.3 | Aspect 3: many highly correlated variables

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

We emphasize that there is nothing wrong with using PCA on the full data set, it will always reflect the internal correlation structure in the whole data set. The potential problem is that the assessment of the relative importance of underlying sensory dimensions may be biased and sometimes sensory dimensions may appear more/less important than they deserve.

1.4 | Aspect 4: validation

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

1.5 | Aspect 5: QDA used in relation to consumer data

In some cases, not all sensory attributes are important for the purpose they are used for. An example is preference mapping, where for

instance a certain spice or salt level may be important for consumer preference, but its effect is blurred by the presence of a large number of attributes that are irrelevant for this problem. If for instance, only two principal components are considered in external preference mapping, the effect of a single important variable appearing in the third component may pass unnoticed. Another example is studies of satiety, where in most cases only the texture attributes will be relevant (Nguyen, Næs, Almøy, & Varela, 2019), not the whole sensory profile.

The present article is a discussion of these five aspects with focus on interpretation and what type of effects they may have on the results. Both personal experience, concrete results from sensory data and basic principles will be important in the discussion. The main purpose is to provide guidelines for the sensory analyst in industry and science and suggestions of how to use PCA in a safe and reliable way. The article is not intended for the specialist statistician, but for the more typical users of these methods in their daily activities and practice. Some possible pitfalls are underlined and some new suggestions and tools will be presented and discussed. A short introduction to PCA is provided here, but for a thorough description of several more aspects of PCA we refer to Jolliffe (2010). At the end of the article (Section 10), a number of conclusions and recommendations are given for each of the issues discussed. The phenomena discussed will be illustrated by examples using real sensory data sets.

2 | STRUCTURE OF DESCRIPTIVE SENSORY DATA

The focus of the present article is the use of PCA for descriptive sensory data (QDA data). In most cases, the entries in such data sets will lie between a lower and an upper limit on some sort of intensity scale. The different attributes are calibrated to be positioned within this interval. It should be mentioned that although PCA is a very important tool in this context, a proper analysis and interpretation of each of the attributes separately is always recommended.

For the purpose of interpretation and also for some of the tools proposed, the sensory data will be thought of as generated according to an experimental design with assessors and products as the two factors in the design. In more technical terms, each sensory variable can be considered a sum of contributions from the two factors, product and assessor, that is,

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

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

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

$$Y = XB + E \tag{2}$$

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

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

Another way of organizing QDA data is by using a three-way array structure with the rows corresponding to samples*replicates, columns to attributes and slices to the different assessors (Figure 1b). This type of data structure can be analyzed by so-called multi-way methods such as PARAFAC (Bro, Qanari, Kiers, Næs, & Frost, 2008), or one of the Tucker methods (Tucker, 1964), which are extensions of standard PCA. The data set organized as in Equation (2) is referred to as a three-way data set, which has been unfolded (see Figure 1b) vertically. The data structure to the right in Figure 1b corresponds to Y in Figure 1a and Equation (2). The three-way structure and analysis will not be pursued further here.

3 | SHORT DESCRIPTION OF PCA

Principal component analysis is a so-called component method. This means that it is based on the idea that a large number of variables in Y can be approximated by a small number of so-called components T (sometimes called axes or latent variables) calculated as linear combination YW , where W is the matrix of so-called loading weights

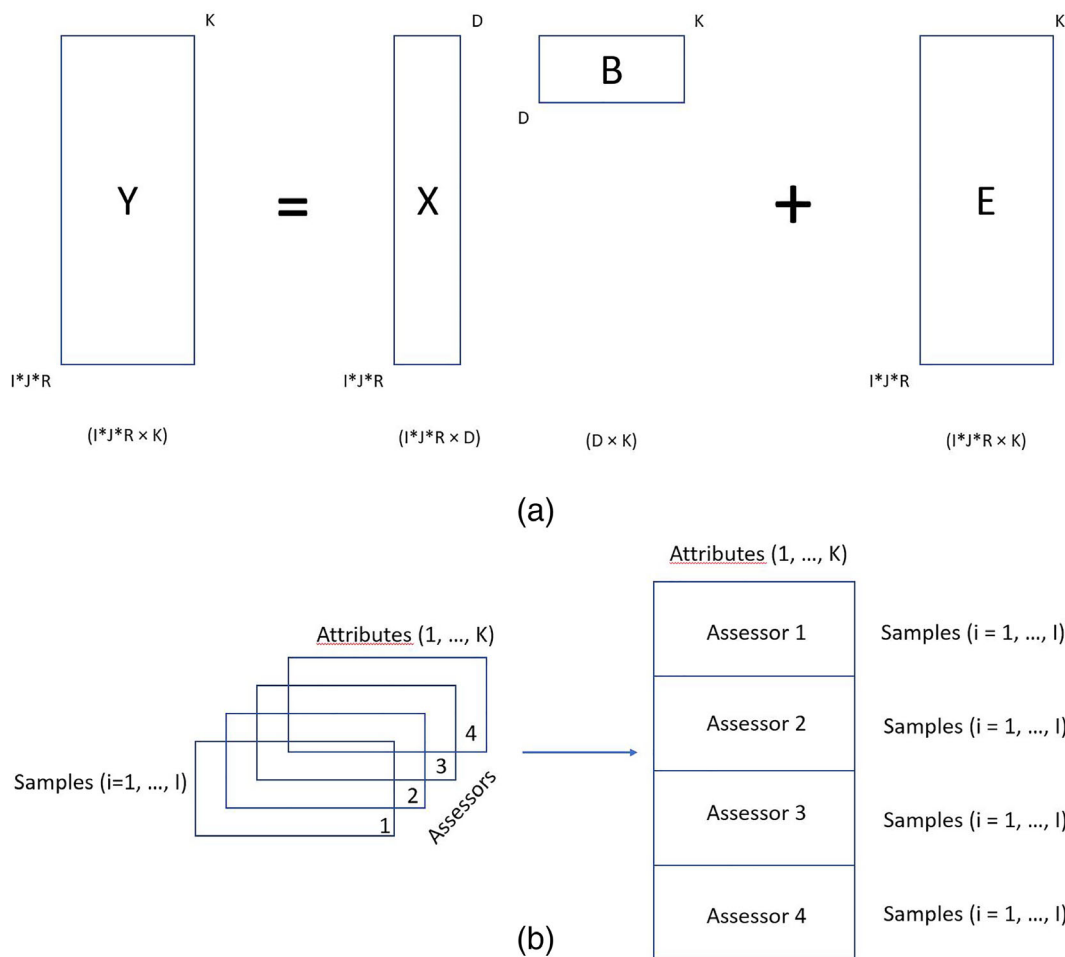


FIGURE 1 (a) Illustration of the setup in Equation (2). The D now represents the number of design variables (including product and assessor factors plus interactions). (b) Data structure for quantitative descriptive analysis presented as a three-way data set and an unfolded data set. The illustration is for simplicity only for four assessors. If replicates are present, the vertical dimension will be samples*replicates ($I \times R$)

(columns of \mathbf{W} have length = 1). The components are found by maximizing their variance and such that each new component extracted is orthogonal/uncorrelated with previous ones. The first component describes the most of the variability, the second is the next in the order etc. A consequence of the criterion used is that variables or variable groups with large variance will have a stronger impact on the solution than the rest. Usually one extracts only a few components treating the rest of the variability as noise. After calculation of the components, they can be related to \mathbf{Y} by regression in order to find the loadings \mathbf{P} . The model for PCA can be written as

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

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

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

4 | PCA FOR ORIGINAL OR AVERAGED DATA?

4.1 | Averaged data for studying product differences

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

4.2 | Different types of panel averages

It should be mentioned that there are different ways of obtaining panel averages (or a panel consensus). One of them is to use

straightforward averaging as will be focused here. Other possibilities are Generalized Procrustes analysis (Gower, 1975), STATIS (see e.g., Schlich, 1996), multiple factors analysis (MFA, Escofier and Pages (1995)), and various scaling techniques (Romano, Brochoff, Hersleth, Tomic, and Næs (2008)). Generalized Procrustes analysis rotates, reflects, and scales (isotropic scaling) the individual assessor data matrices to make them as similar as possible and then afterward calculates the consensus as the average. The STATIS method calculates a weighted average of the individual (cross-product) matrices, where the weights depend on the RV coefficients between them. MFA concatenates the individual data matrices horizontally and essentially runs a PCA on the combined matrix after a specific individual scaling of each of them. The resulting scores matrix of this PCA is then used as a consensus for the individual assessors. An alternative to MFA, with a similar underlying idea is the Tucker-2 method used in Dahl and Næs (2009). The scaling methods in Romano et al. (2008) are used to eliminate additive and multiplicative differences among assessors before averaging. Note that, all these methods are also suitable for investigating individual differences among assessors (see e.g., Næs et al., 2018).

4.3 | PCA for original data

If focus is also on individual differences between assessors, one can use the original \mathbf{Y} data in Equation (2) directly without averaging. There will be several more points in the score plot, one score for each replicate, assessor and sample combination. For improved interpretation one can include colors and sample averages as will be illustrated here. This plot can be useful for visualizing differences/disagreement among assessors.

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

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

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

5 | STANDARDIZATION

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

For PCA in general, many different types of standardization are used, but here we confine ourselves to the most used namely division

by standard deviation. It should be mentioned that using PCA on standardized data is what some authors phrase as using the correlation matrix as the basis for the calculation of components.

5.1 | Standardization is not primarily a statistical issue

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

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

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

5.2 | Interpretation of PCA with and without standardization

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

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

It is important to mention that when using standardization, the variance of all variables will be the same. This implies that only the number of variables related to a sensory dimension will be the driver for order of the components. If for instance, one phenomenon is described using four highly correlated sensory attributes and another phenomenon is represented by one attribute only, the first principal component will represent the phenomenon with the four attributes and the second component will represent the other variable. Therefore, in such cases, importance of dimensions (in terms of explained variance) is driven by the number of correlated attributes representing the same phenomenon rather than by the most dominating sensory dimension. This shows that it is not obvious how to define the concept of common concept of “most important sensory dimensions” using QDA and PCA.

5.3 | Eliminate nonsignificant attributes

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

5.4 | Using correlation loadings plot

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

It is tempting to think of correlation loadings as a way of eliminating the problem of standardization. However, this is not always the case since correlation loadings only represent a post processing procedure after the principal components have been estimated. The method may be better at highlighting the relations between variables with a small initial variance (and which therefore have little influence

on the solution) and the components, but this does not change the data for which PCA is calculated. For standardized data, the two are the same except for a scaling factor. We here use the unit circle scaling for the correlation loadings.

6 | CORRELATIONS BETWEEN VARIABLES

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

6.1 | Avoiding highly correlated variables

It is generally recommended that too much repetition of information should be avoided in order to reduce unnecessary bias and focus for the PCA. Some of these repetitions may be quite obvious such as using confounding attributes as for example, dark/light and hard/tender (see introduction), while others may be more subtle and difficult to identify directly without data analysis. Assessors may for instance have problems discriminating between two or more cognitively similar attributes and will automatically score them similarly. This is known as halo dumping effect. It comes from the human desire of consistent cognitive structures and has been widely described in the sensory literature (see e.g., Clark & Lawless, 1994). Correlation between unrelated attributes may also happen when one salient negative attribute causes another to be rated in the same direction, Such correlations are known as horn effects, common when describing defective samples (Lawless and Heyman (2010)). This is an unfortunate situation and having tools to detect such cognitive coincidence is important for more relevant analysis and interpretation of PCA and for improved training of the panel. One of the objectives of panel training is to achieve de-correlation of the attributes, and avoid redundancy leading to particular issues in multi-product panels, as some attributes can be correlated for one product but not for another.

6.2 | Correlations at different levels

Correlation between attributes/columns in \mathbf{Y} can be due to correlation induced by the design (\mathbf{X} in Equation (2), representing sample, assessor and interaction) and by the random error \mathbf{E} in the model. The correlations between variables in \mathbf{XB} are the most important since these are functions of the design of the study. Correlations among the variables in \mathbf{E} are, however, conceptually more problematic. This calls for investigating the correlation structure for \mathbf{XB} and \mathbf{E} separately and

sometimes also for the products and assessors separately. We will next discuss a possible tool to use for detecting correlations among the variables in the before we describe briefly a few methods for studying \mathbf{XB} by PCA.

6.3 | Partial correlation for detecting correlations among random errors in Equation (2)

The concept of partial correlation between variables was developed for the purpose of correlating two variables with each other after they have been conditioned upon a third variable (or set of variables). This is equivalent to correlating the residuals \mathbf{E} for the two variables with each other after they have been regressed onto the same variables. If the partial correlation among two variables is high, one should consider eliminating one of them from the PCA to avoid the problem discussed above. This type of information may also be important for retraining the panel and to improve its performance. Since this type of correlation will most typically be present at the individual level, correlation between residuals at an individual level will be given the strongest focus here.

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

6.4 | PCA for the systematic part \mathbf{XB} of Equation (2)

An important PCA based methods for analyzing the systematic part \mathbf{XB} is ASCA (Jansen et al., 2005). PCA plots for this method can be used to reveal cases with highly overlapping attributes as discussed above. The effects of the assessor and product (and their interactions) are first estimated using the model (1) and standard ANOVA methods. Then the effects for the different factors are further analyzed by PCA using all the response variables. This is equivalent to estimating \mathbf{B} in Equation (2), then splitting the \mathbf{XB} contribution into three parts, the assessor part, product part and the interaction part. Analyzing each of them by PCA results in three separate PCA models. In mathematical terms this means that \mathbf{XB} is essentially written 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 the \mathbf{B} 's. In this way, information is obtained about the variability structure of the sensory attributes for the assessors, products, and interactions separately (see Liland, Smilde, & Næs, 2018). This means that this method can reveal correlation structure at the sample level and assessor level separately. The PC-ANOVA (Luciano and Næs (2009)) is related, but reverses the order of ANOVA and PCA. First a PCA is run for \mathbf{Y} and then the scores for the first few components are related separately to the design using the model (1).

7 | VALIDATION OF PCA MODELS

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

7.1 | External validity

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

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

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

7.2 | Internal validity

Internal validity of a component means that a component is more meaningful or describes a larger percentage of variance than the variance that can be obtained by chance, that is, in data sets without an underlying structure. Therefore, comparing true explained variance with what is obtained by chance is a possibility. This type of validity is only referring to the data set under study and will not tell anything about how well the model represents a population of other samples. The CV as defined by Wold (1978), which is based on successively

creating subsets for validation by eliminating entries according to a diagonal pattern of the data set, can be considered an internal validation method. Here we will, however, concentrate on a method based on permutations as proposed in Endrizzi, Gasperi, Rødbotten, and Næs (2014) and later studied and modified by Vitale et al. (2017). We will here use the original version.

7.2.1 | Permutation testing

The idea behind the method is that for each new component to be tested, the residuals from the model based on all previous components are permuted (for each column separately) and then orthogonalized with respect to both columns and rows (since this is the case for the true residuals in a PCA). Then, one calculates the explained variance of the permuted residuals data set and compares it with the true explained variance. This is done by comparing the explained variances for the component considered relative to the variance left in their respective data sets (permuted residuals and true residuals). The procedure is repeated for a large number of permutations (e.g., 1,000, as used here). The results are then presented in a plot with component number on the X-axis and the explained variances as described above on the Y-axis. For the real data, there is only one point for each component, but for the permuted data, we will here present three values, the median, the lower 5% percentile and the upper 5% percentile, obtained from a large number of permutations. The lower and upper values are there for assessing the uncertainty of the estimates. If the true value falls clearly above the confidence band obtained by the two percentiles, the component can be judged significantly different from that generated by chance and therefore worth looking at. Although assessing the number of components is essentially a one-sided test, we here prefer the setup used to indicate the uncertainty in both directions. For details we refer to Endrizzi et al. (2014).

7.2.2 | Assessor based CV

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

7.3 | Validation using external information

In some cases, there may be other data available about the samples, for instance, chemistry data, spectroscopy data, or simply the

experimental design. In such cases it is possible to regress the (for instance) average sensory attribute scores (across assessor and replicates) onto the external data and then evaluate how much of the sensory data that can be accounted for by the external variables/measurements. Such a method was used in Dahl and Næs (2004) for relating the average sensory profile to external near infrared (NIR) spectra. Explained variance of the sensory profile obtained from the NIR data was then used as criterion of validity. In the article, the same was also done for each individual assessor separately in order to identify outliers.

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

7.4 | Validation using confidence intervals

In addition to focusing directly on the significance of a component, confidence intervals or ellipsoids for each sample is a good option. They are primarily meant for assessing stability of solutions, but can also be useful for indicating how many components that are worth considering. Bootstrap procedures as illustrated for instance in Cadoret and Husson (2013) are the most important to use in this case. The method is based on resampling assessors at random (the same number as in the original panel) and calculating the scores for each selection (after averaging over assessors). These are then projected onto the scores plot of the original averaged PCA and confidence ellipses are drawn based on this for each sample.

8 | IMPLICATIONS FOR RELATIONS TO CONSUMER DATA

As mentioned in the introduction, very often a sensory data set is not only used for understanding the variability in the sensory properties of samples. A typical example is preference mapping where the main focus is on relating consumer liking to sensory data. One can do this by analyzing one sensory attribute at a time, but a more typical way is to use PCA of the sensory data (or PLS regression) and regress the

liking for different consumers onto the first couple of components (often only 2). If then a specific attribute with minor relation to the main variability of the sensory data set, has an important influence on the liking, it will not be visible in standard external preference mapping analysis with two components. Typical examples are salt level and spices which may influence liking strongly, but do not account for much variability in the sensory data. One should therefore inspect more than two components or supplement (or replace) the analysis with an internal preference mapping, where PCA is applied to the liking data and sensory data are regressed onto these principal components. PLS regression could be another alternative for such data (see e.g., Næs et al., 2018).

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

9 | CASE STUDIES

9.1 | Data sets used

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

9.2 | Case 1. Should one average or not before computing PCA on sensory data? Exemplified using yogurt data

The data used for visualizing the differences between using the PCA for average data and for the individual data before averaging is a yogurt dataset with 8 samples and 21 attributes, (Nguyen et al. (2019)). An experimental design with three factors at two levels is used for producing the samples. In this case, we focus on standardized data for visualization (after elimination of the single nonsignificant attribute at 5% level).

The results are presented for panel averages and raw data in Figure 2 and Figure 3. In Figure 3, the average component scores across assessors for each sample are superimposed using diamond shapes. As can be seen, the loadings are quite similar for the two PCA models, but the explained variances are larger for the averaged data due to the averaging process, as explained above. The main difference in loadings is that dryness in mouth and astringent form an own group of attributes for the individual data while for standardized data they are grouped together with sandy, stale odor, and so on. There are quite large individual differences around each sample average in Figure 3 (scores with same color). Still, the average scores for each sample are quite similar to the scores in Figure 2. This means that the essential information is similar for the two analyses. The former provides a simpler plot, while the second gives an opportunity for

TABLE 1 Overview of quantitative descriptive analysis data sets used

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

Note: For the olive oil data set also the small subset is tested. For the bread data also consumer liking data for a number of consumers were available.

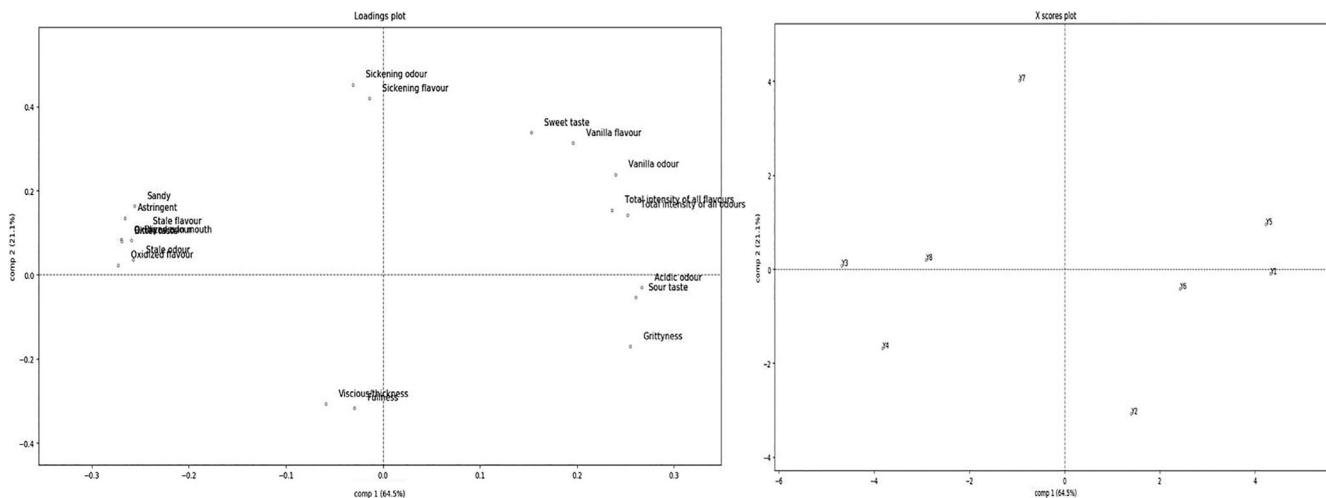


FIGURE 2 Yogurt data. Standardized PCA on consensus data, 20 significant attributes

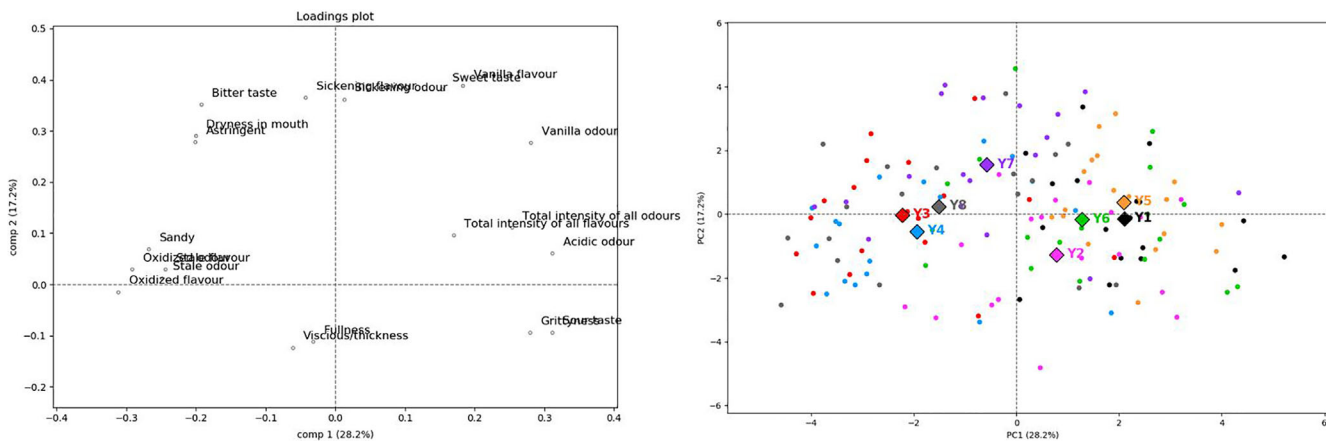


FIGURE 3 Yogurt data. Scores and loadings for the standardized PCA based on individual data, 20 significant attributes

studying individual differences. As will be seen below, the latter also allows for an ANOVA test for the components. In practice choosing between the two is often a matter of scope of the study and need for simplicity. Most of the discussion below will be focused on average data.

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

An illustration of the effect of standardization will be given using data from sensory analysis of olive oil (based on averages over

assessors). The results are presented in Figure 4–d. Figure 4a gives results from PCA on the full set of variables without standardization, while in Figure 4b, PCA is based on the full set of standardized variables, Figure 4c shows results of PCA for only significant variables, not standardized, while Figure 4d shows PCA results for significant standardized variables. In all cases the explained variances were high, about 90% after three components. The three components look significant using leave-one-out CV, and this is also confirmed by the other permutation based method to be shown below.

The Figure 4a shows that loadings and correlation loadings plot are quite different without standardization. The Figure 4b shows that

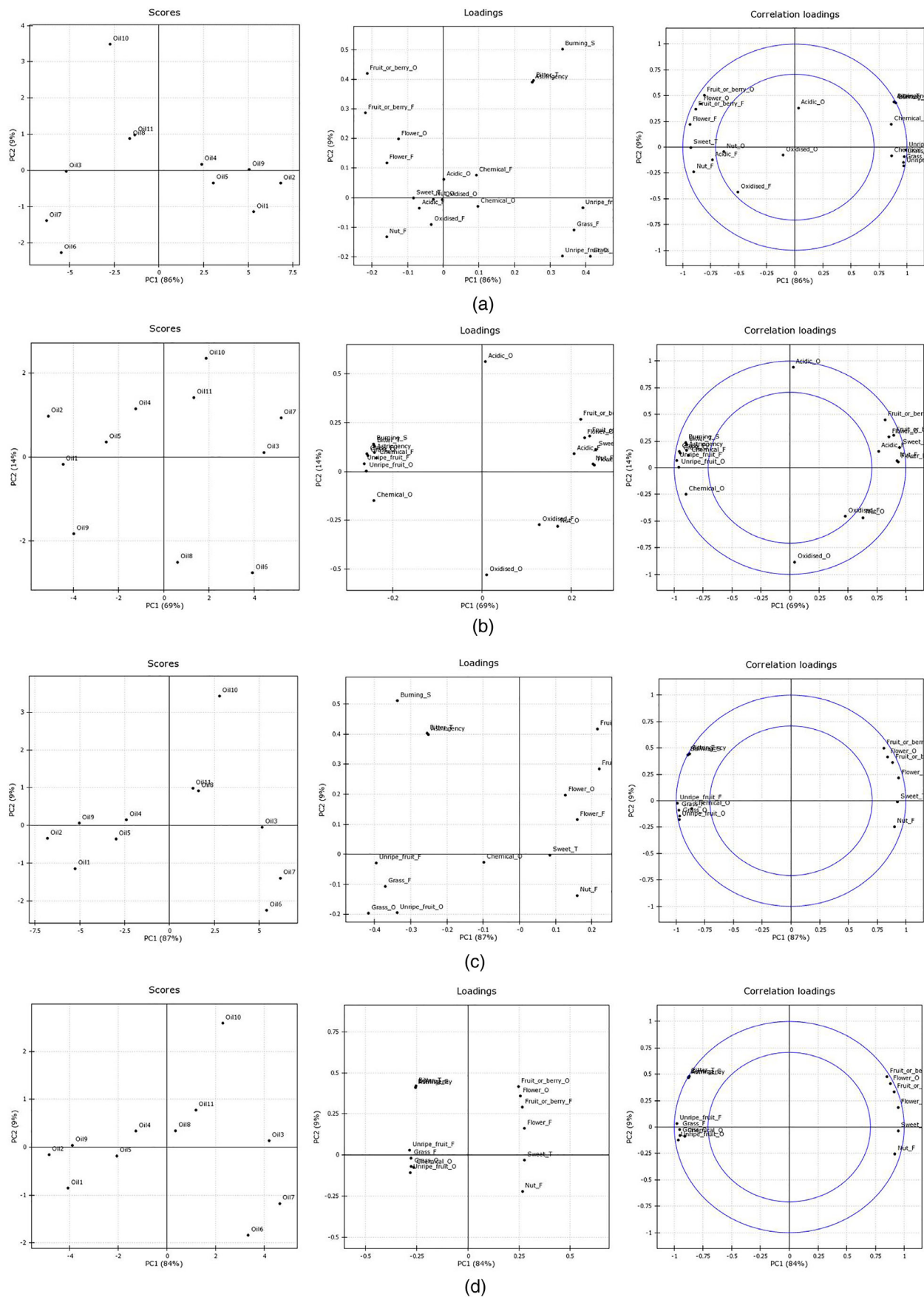


FIGURE 4 (a) Olive oil data. Full data set nonstandardized. (b) Olive oil data. Full data set standardized. (c) Olive oil data. Reduced data set nonstandardized. (d) Olive oil data. Reduced data set standardized

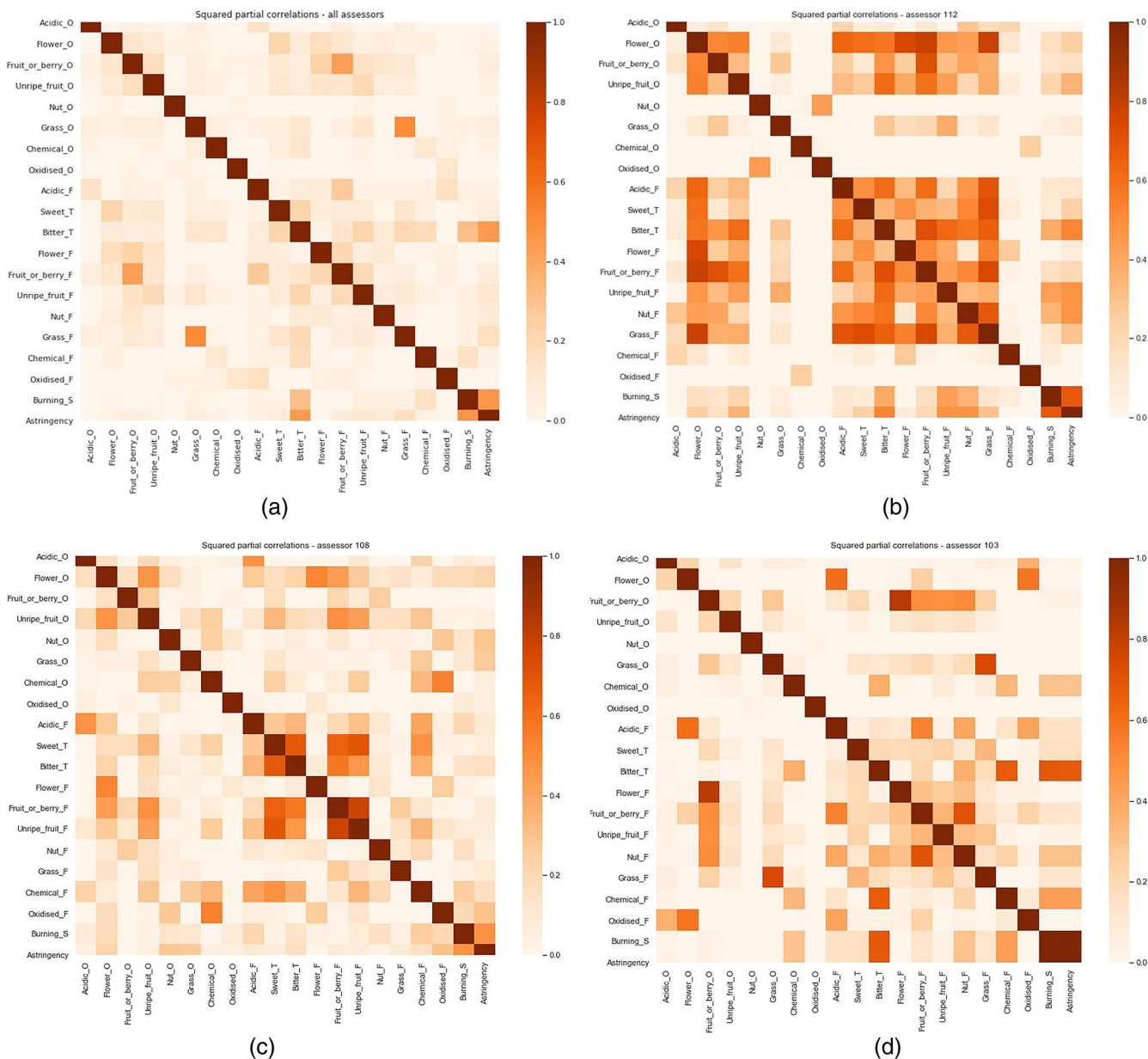


FIGURE 5 Olive oil data. Heat map of correlations between residuals for different attributes. Over all assessor in (a). The other three, (b–d), represent three individual assessors

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

After eliminating nonsignificant variables (Figure 4c. Six attributes eliminated), we see that the scores are back again to the ones obtained without standardization for the full set of variables (Figure 4a). Correlation loadings and loadings are still different, but less so if we compare with the full data set. Standardization (Figure 4d) now has little effect (for reduced data) on the loadings except for one variable close to the middle. Scores are almost the same for Figure 4c,d. After standardization, loadings and correlation loadings in Figure 4d are identical except for the scaling.

In conclusion, after elimination of nonsignificant variables, the results are similar regardless of whether one standardized or not. This is true for both scores and loadings.

Comparing full and reduced data sets, we see that scores are almost the same except for the standardized full data set (Figure 4b). Two of the attributes (acidic-O and oxidized-O) that show up in the full data set along the second component are not present in Figure 4c,d since they are nonsignificant. They are also less visible in Figure 4a. These two are examples of variables that are “inflated” when standardized. This phenomenon is quite frequent with off-flavors or other attributes that may appear in low intensities (i.e., spicy). After standardization low scoring attributes will get a larger importance in the outcome.

Our advice is to eliminate nonsignificant variables since it then matters less what is done regarding standardization. The standardized results

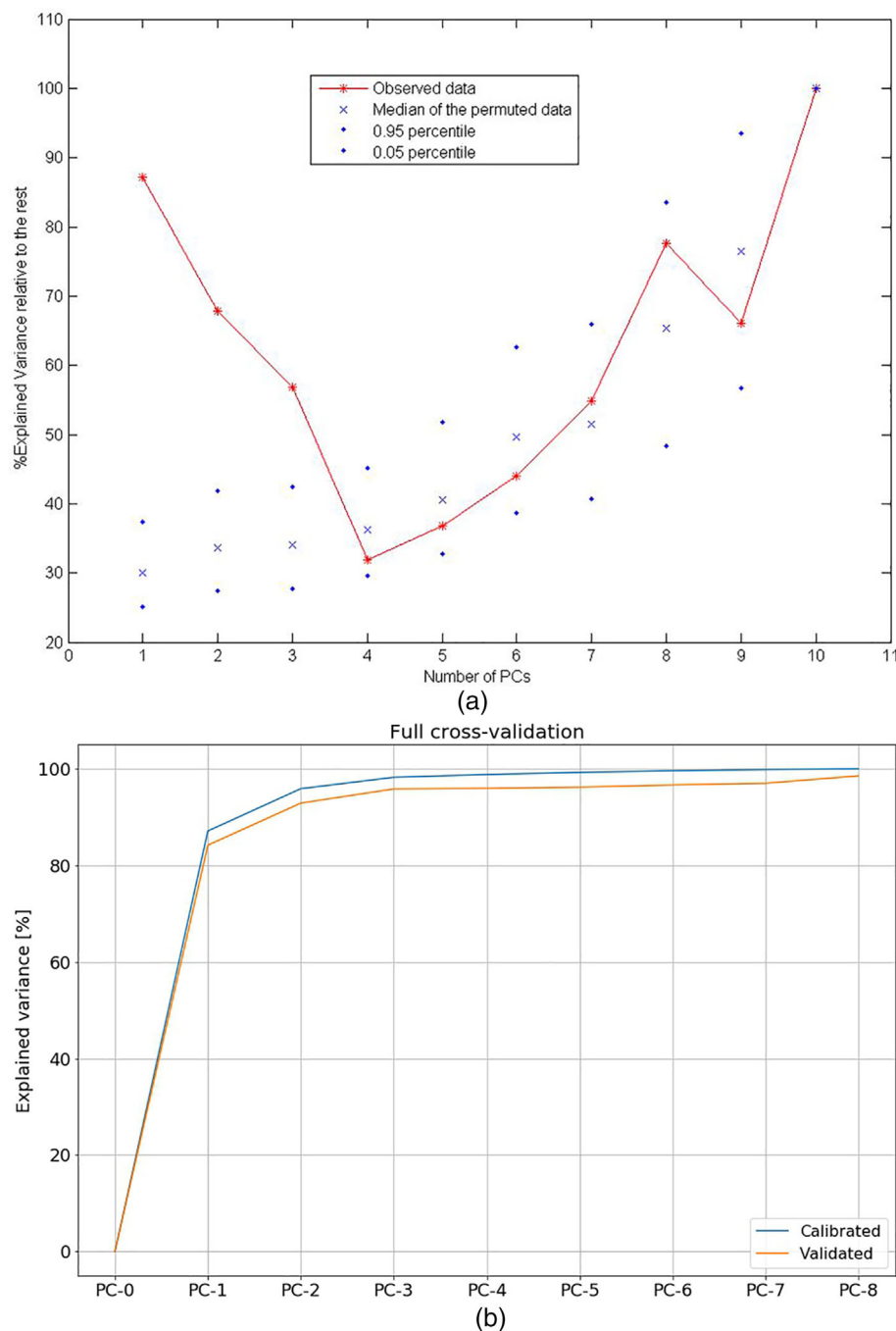


FIGURE 6 Olive oil data. Nonstandardized PCA, 14 significant attributes. The illustration in (a) shows the curve obtained by the permutation method. The points represent the quantiles for each of the number of components. In (b) is presented explained variance for fitting/calibration and leave-one-out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

with all variables, including nonsignificant ones, are the most different from the rest. One should focus on a good training for the low scoring attributes when relevant for the products or objective of the study.

9.4 | Case 3. Many correlated sensory variables. Exemplified using yogurt and olive oil data

Figure 2 shows PCA results from the yogurt experiment in Nguyen et al. (2019) (based on a 2^3 design). Most of the variables contrast each other along the first axis. This means that the large variability accounted for along this axis to a large extent is due to the many

variables measuring more or less the same phenomenon. This is important information per se, but it clearly gives a biased impression of the relative importance of the two components or underlying dimensions (62 and 20%). Eliminating several of the highly correlated variables along the first component, leads to a different relative weighting of the two axes. In other words, the relative importance of the components is dependent on how many strongly correlated variables that are in the data set.

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

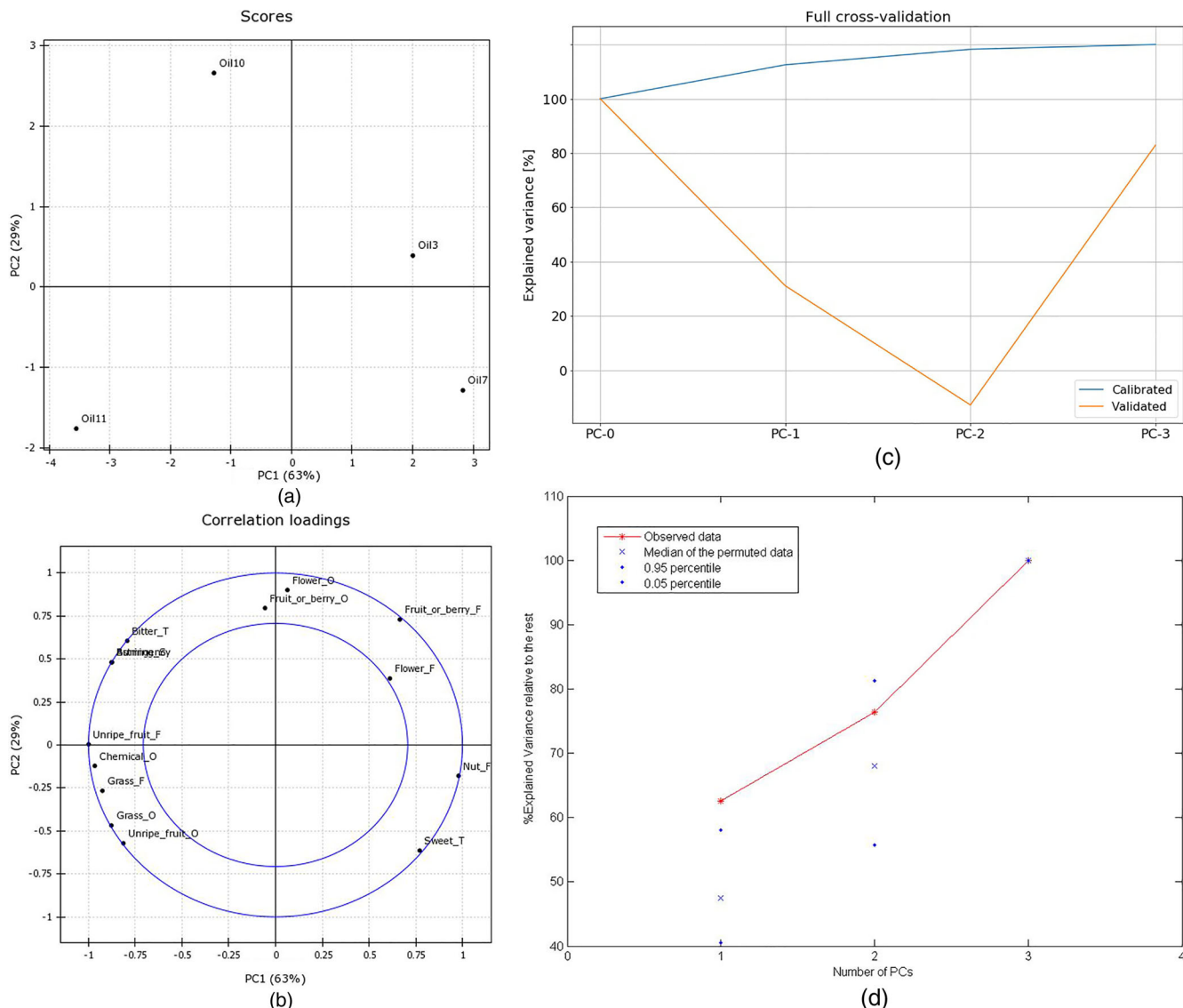


FIGURE 7 Olive oil data. Four samples, standardized PCA, 14 significant attributes. (a) scores, (b) correlation loadings, (c) cross-validation, and (d) permutation testing

9.4.1 | Partial correlation results

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

9.5 | Case 4. Validation based on CV and permutation testing. Exemplified using olive oil data

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

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

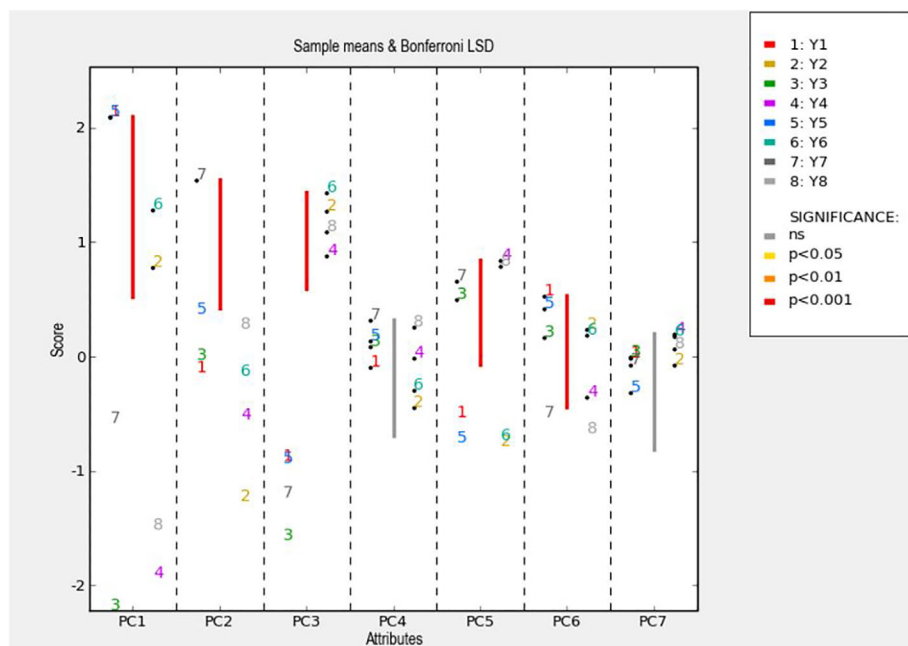
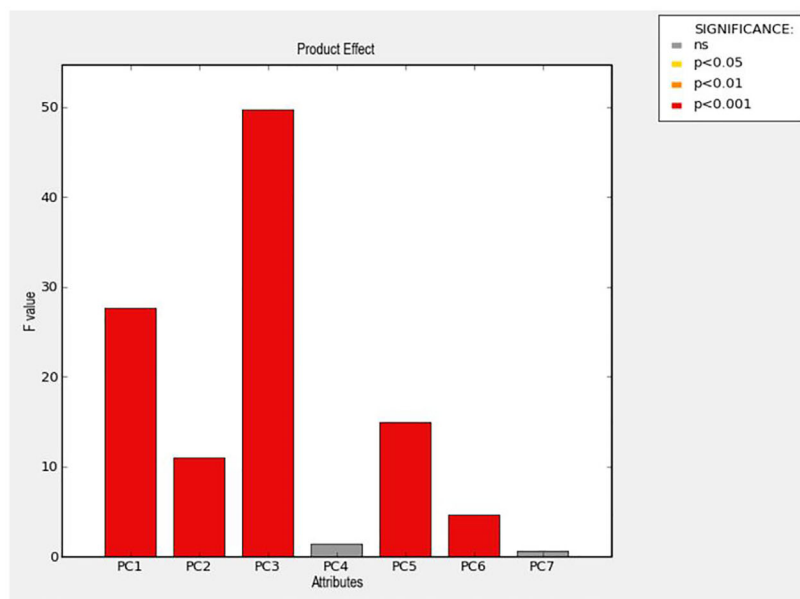


FIGURE 8 Yogurt data. PCA-ANOVA results, standardized PCA, 20 significant attributes. (a) Multiple comparisons for products. Line indicates range of no significant differences. (b) F-values for the product effect factor. The significance is indicated with color as given in panel in the upper right corner

(a)



(b)

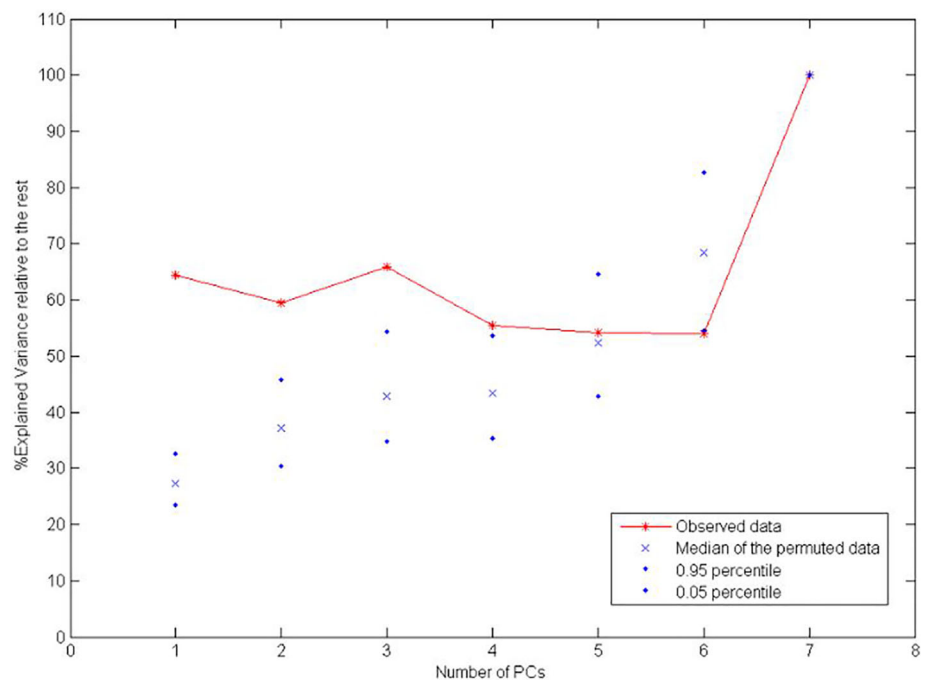
9.5.1 | An illustration based on reduced data

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

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

An interesting observation is that the loadings plot change when a subset (oils 3, 7, 10, and 11) of the full set of samples (oil 1–11) is

FIGURE 9 Yogurt data. Standardized PCA. Twenty significant attributes. Permutation test for PCA based on averages over assessors



used (see Figure 4d). This underlines that interpretation of a subset of samples only relates to this specific subset at hand and cannot be generalized to the sensory space of the full set of samples. Conclusions will then always be local and of limited value for saying something about a larger set of “similar” samples.

9.5.2 | The use of PC-ANOVA for validation

PC-ANOVA (Luciano & Næs, 2009) was applied to the standardized yogurt data and compared to the use of the permutation test for the consensus/average data set. The results are presented in Figures 8a,b and 9. As can be seen, the results correspond reasonably well, the first three components are obviously significant, while number 4 is more questionable. It seems that the PC-ANOVA finds significance further out (components 5 and 6), but these components represent so small variance that they are not very interesting in practice. Also, the fact that component number 4 is nonsignificant is an indication that one should not consider further components after Component 3. The explained variances for the 5 first consensus components are 64.4, 21.1, 9.5, 2.7, and 1.2. For the PCA done on raw data the corresponding 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 the first to the second component.

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

For this example based on external preference mapping, a bread data set with 8 samples (based on a 2^3 design) and 13 attributes is

used. The data set consists of both QDA data and consumer liking of the same samples. Only the averages will be considered for QDA.

In Figure 10a,b, correlation loadings plots of Component 1 versus Component 2 and for Component 1 versus Component 3 are shown. As can be seen, there is a major tendency in liking toward Component 3 dominated by salt taste. This tendency is not visible in the plot of Component 1 versus Component 2 where salt is lying well within the 50% explained variance circle.

This shows that relying only on a two-dimensional external preference mapping plot can leave important drivers of liking undetected.

10 | CONCLUSIONS AND SUGGESTIONS

10.1 | Using averages over assessors or raw data

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

10.2 | Standardization

The calibration and training procedure should be considered and evaluated for making a decision on whether to standardize or not. The focus should be on the meaningfulness of relying on actual

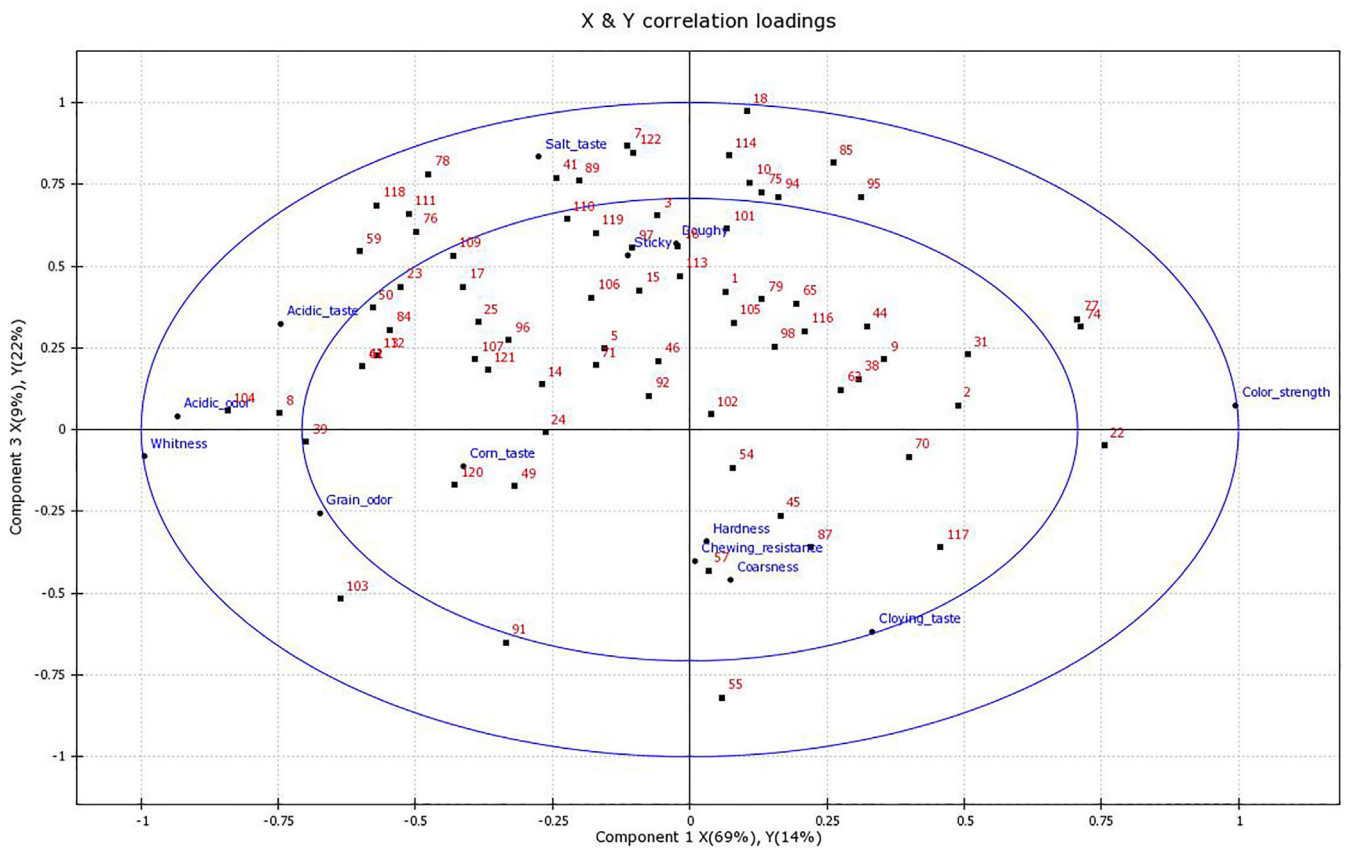
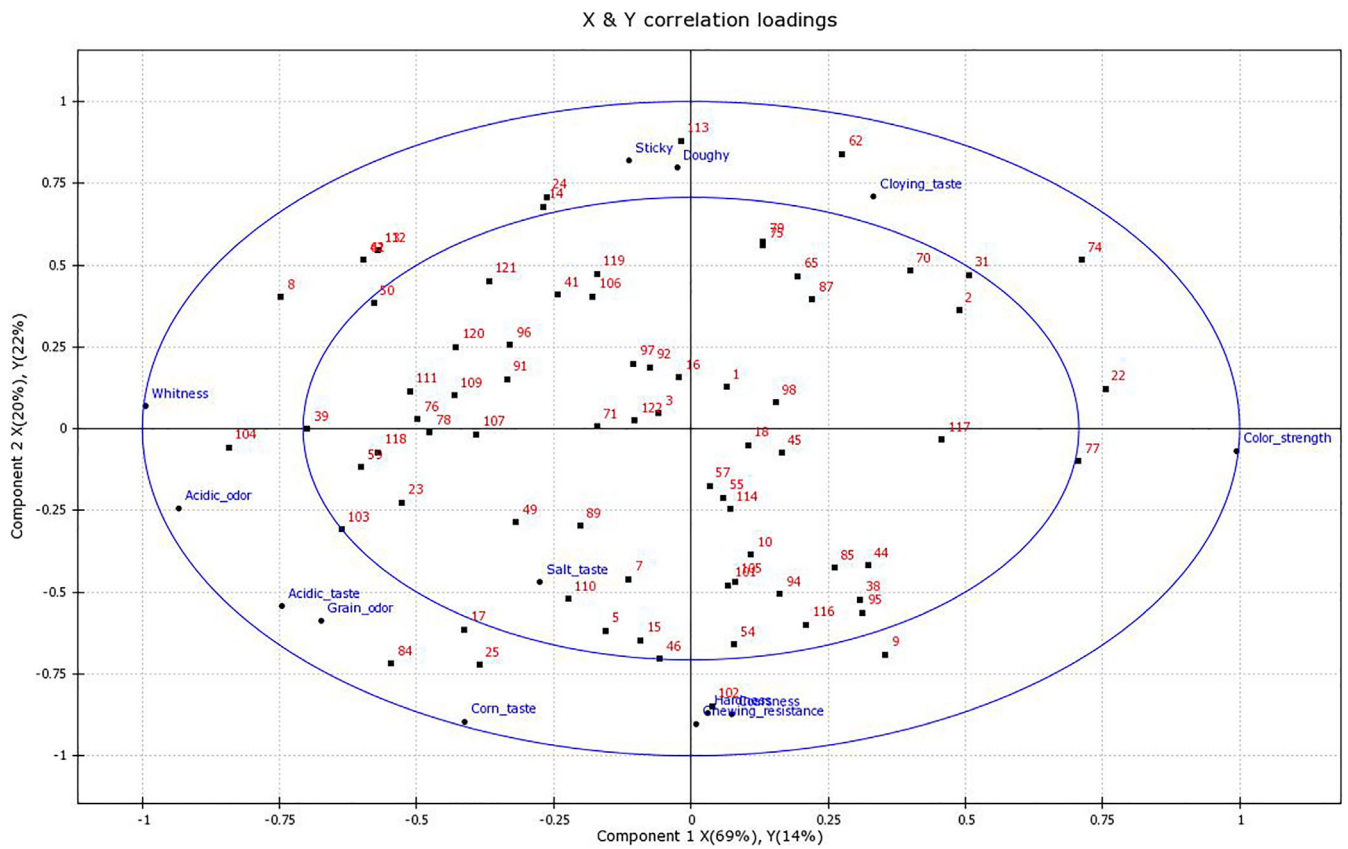


FIGURE 10 Bread data. Nonstandardized PCA. Correlation loadings for external preference mapping. (a) Component 1 versus Component 2. (b) Component 1 versus Component 3

differences in variability of different attributes (possibly belonging to different sensory modalities) in the analysis. If these are not meaningful, one should standardize. This is an interesting aspect when comparing results from different panels. In such cases, the need for standardization is stronger unless the training procedure is harmonized between the labs. If clearly nonsignificant variables are present, one should be careful about incorporating them in a standardized analysis.

10.3 | Using all attributes or eliminating obvious overlap

Eliminating highly correlated variables will in most cases have only a moderate effect on the interpretation. One should be careful about strong statements about what are the most important sensory dimensions since this will depend on the number of attributes that represent it. A tool based on partial correlations is presented that can enhance insight into nontrivial overlap among attributes.

10.4 | Validation of components

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

10.5 | Relating sensory QDA data to consumer liking data

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

ACKNOWLEDGMENTS

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

ORCID

Tormod Næs  <https://orcid.org/0000-0001-5610-3955>

Paula Varela  <https://orcid.org/0000-0003-2473-8678>

REFERENCES

Bro, R., Qanari, E. M., Kiers, H. A. L., Næs, T., & Frost, M. B. (2008). Multi-way models for sensory profiling data. *Journal of Chemometrics*, 22, 36–45.

- Brockhoff, P. B., Schlich, P., & Skovgaard, I. (2015). Taking individual scaling differences into account by analyzing profile data with the mixed assessor model. *Food Quality and Preference*, 39, 156–166.
- Cadoret, M., & Husson, F. (2013). Construction and evaluation of confidence ellipses applied at sensory data. *Food Quality and Preference*, 28, 106–115.
- Clark, C. C., & Lawless, H. T. (1994). Limiting response alternatives in time-intensity scaling: An examination of the halo-dumping effect. *Chemical Senses*, 19, 583–594.
- Dahl, T., & Næs, T. (2004). Outlier and group detection in sensory analysis using hierarchical clustering and the Procrustes distance. *Food Quality and Preference*, 15(3), 195–208.
- Dahl, T., & Næs, T. (2009). Identifying outlying assessors in sensory profiling using fuzzy clustering and multi-block methodology. *Food Quality and Preference*, 20, 287–294.
- Dahl, T., Tomic, O., Wold, J. P., & Næs, T. (2008). Some new tools for visualising multi-way sensory data. *Food Quality and Preference*, 19, 103–113.
- Dijksterhuis, G. (1995). Assessing panel consonance. *Food Quality and Preference*, 6(1), 7–14.
- Endrizzi, I., Gasperi, F., Rødbotten, M., & Næs, T. (2014). Interpretation, validation and segmentation of preference mapping models. *Food Quality and Preference*, 32, 198–209.
- Eriksson, L., Johansson, E., Kettaneh-Wold, N., & Wold, S. (1999). Scaling. In *Introduction to multi- and megavariate data analysis using projection methods (PCA & PLS) Umetrics* (pp. 213–225). Umeå: Umetrics Academy.
- Escofier, B., & Pages, J. (1995). Multiple factor analysis. *Computational Statistics and Data Analysis*, 18, 121–150.
- Gower, J. C. (1975). Generalized Procrustes analysis. *Psychometrika*, 40, 33–51.
- Jansen, J., van der Hoefsloot, J., Greef, M., Timmerman, E., Westerhuis, J., & Smilde, A. K. (2005). ASCA: Analysis of multivariate data obtained from an experimental design. *Journal of Chemometrics*, 19(9), 469–481.
- Jolliffe, I. T. (2010). *Principal component analysis*. New York: Springer.
- Lawless, H. T., & Heyman, H. (2010). *Sensory evaluation of food: Principles and practices*. New York, NY: Springer Science and Business Media.
- Liland, K. H., Smilde, A., & Næs, T. (2018). Confidence ellipsoids for ASCA models based on multivariate regression theory. *Journal of Chemometrics*, 32, 1–13. <https://doi.org/10.1002/cem.2990>
- Luciano, G., & Næs, T. (2009). Interpreting sensory data by combining principal component analysis and analysis of variance. *Food Quality and Preference*, 20(3), 167–175.
- Martens, H., & Martens, M. (2001). *Multivariate analysis of quality: An introduction*. Chichester, UK: John Wiley and Sons.
- Næs, T., Varela, P., & Berget, I. (2018). *Analysing individual differences in sensory and consumer science*. UK: Elsevier.
- Nguyen, Q. C., Næs, T., Almøy, T., & Varela, P. (2019). Portion size selection as related to product and consumer characteristics studied by PLS path Modelling. *Food Quality and Preference*, 79, 103613.
- PanelCheck software. (n.d.) Available from <https://sourceforge.net/projects/sensorytool/>, <https://doi.org/10.5281/zenodo.10768>
- Romano, R., Brochhoff, P. B., Hersleth, M., Tomic, O., & Næs, T. (2008). Correcting for different use of the scale and the need for further analysis of individual differences in sensory analysis. *Food Quality and Preference*, 19, 197–209.
- Schlich, P. (1996). Defining and validating assessor compromises about product distance and attribute correlations. In T. Næs & E. Risvik (Eds.), *Multivariate analysis of data in sensory science*. Amsterdam: Elsevier.
- Stone, H., Bleibaum, R., & Thomas, H. (2021). *Sensory evaluation in practice* (4th ed.). UK: Elsevier.
- Stone, M. (1974). Cross-validators choice and assessment of statistical prediction. *Journal of the Royal Statistical Society, Series B*, 36, 111–133.

- Tomic, O., Forde, C., Delahunty, C., & Næs, T. (2013). Performance indices in descriptive sensory analysis – A complimentary screening tool for assessor and panel performance. *Food Quality and Preference*, 28, 122–133.
- Tomic, O., Luciano, G., Nilsen, A., Hyldig, G., Lorensen, K., & Næs, T. (2010). Analysing sensory panel performance in a proficiency test using the PanelCheck software. *European Food Research and Technology*, 230(3), 497–511.
- Tomic, T., Nilsen, A., Martens, M., & Næs, T. (2007). Visualization of sensory profiling data for performance monitoring. *LTW*, 40, 262–269.
- Tucker, L. R. (1964). The extension of factor analysis to three-dimensional matrices. In N. Fredriksen & H. Gulliksen (Eds.), *Contributions to mathematical psychology*. New York, NY: Holt, Rinehart and Winston.
- Vitale, R., Westerhuis, J. A., Næs, T., Smilde, A. K., de Noord, O. E., & Ferrer, A. (2017). Selecting the number of factors in principal component analysis by permutation testing - Numerical and practical aspects. *Journal of Chemistry*, 31, 12.
- Wold, S. (1978). Cross-validators estimation of the number of components in factors analysis and principal component models. *Technometrics*, 20, 397–406.

How to cite this article: Næs, T., Tomic, O., Endrizzi, I., & Varela, P. (2021). Principal components analysis of descriptive sensory data: Reflections, challenges, and suggestions. *Journal of Sensory Studies*, 36(5), e12692. <https://doi.org/10.1111/joss.12692>