



HAL
open science

Reduction of repeatability error for analysis of variance-Simultaneous Component Analysis (REP-ASCA): Application to NIR spectroscopy on coffee sample

Maxime Ryckewaert, Nathalie Gorretta, Fabienne Henriot, Federico Marini,
Jean-Michel Roger

► To cite this version:

Maxime Ryckewaert, Nathalie Gorretta, Fabienne Henriot, Federico Marini, Jean-Michel Roger. Reduction of repeatability error for analysis of variance-Simultaneous Component Analysis (REP-ASCA): Application to NIR spectroscopy on coffee sample. *Analytica Chimica Acta*, 2020, 1101, pp.23-31. 10.1016/j.aca.2019.12.024 . hal-02963035

HAL Id: hal-02963035

<https://hal.inrae.fr/hal-02963035>

Submitted on 7 Mar 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

1 Reduction of repeatability error for Analysis of
2 variance-Simultaneous Component Analysis
3 (REP-ASCA): Application to NIR spectroscopy on
4 coffee sample

5 Maxime Ryckewaert^{a,b}, Nathalie Gorretta^b, Fabienne Henriot^a, Federico
6 Marini^c, Jean-Michel Roger^b

7 ^a*Limagrain Europe, Chappes, France*

8 ^b*ITAP, Univ Montpellier, Irstea, Montpellier SupAgro, Montpellier, France*

9 ^c*Department of Chemistry, University of Rome "La Sapienza", P.le Aldo Moro 5,
10 I-00185 Rome, Italy*

11 **Abstract**

12 A method to reduce repeatability error in multivariate data for Analysis
13 of variance-Simultaneous Component Analysis (REP-ASCA) has been devel-
14 oped. This method proposes to adapt the acquisition protocol by adding a
15 set containing repeated measures for describing repeatability error. Then,
16 an orthogonal projection is performed in the row-space to reduce the re-
17 peatability error of the original dataset. Finally, ASCA is performed on the
18 orthogonalized dataset. This method was evaluated on NIR spectral data of
19 coffee beans. This study shows that the repeatability error due to physical
20 variations between measurements can alter results of the analysis of variance.
21 These effects are predominant in factors analysis and can be seen on spec-
22 tra as constant or non-constant baselines. By reducing repeatability error

Email address: maxime.ryckewaert@limagrain.com (Maxime Ryckewaert)

Preprint submitted to Analytica Chimica Acta

November 4, 2019

1 with REP-ASCA, baselines are removed and factor analysis provides more
2 information about chemical content of the factors of interest.

3 **1. Introduction**

4 When conducting measurements in experiments, measured values can dif-
5 fer from the true values of the samples. Measurement error is the difference
6 between a measured quantity and its true value. For reliable measures, the
7 measurement error must be relatively small compared to the variance of the
8 measured samples [1].

9 The reproducibility error is the difference between two measurements of
10 the same sample under two different conditions, for example, change of lab-
11 oratory or operator. When conditions of acquisitions do not change, the
12 reproducibility error is also called repeatability error [2]. In this paper, the
13 repeatability error stands for differences observed between two measurements
14 performed on the same physical sample.

15 In an experiment, the observed responses can contain repeatability error
16 and this may affect the results of processing carried out on this data. This is
17 especially true when we want to study the influence of independent variables,
18 i.e., the so-called factors.

19 The use of design of experiments (DoE) allows defining the experiments
20 to perform in an optimal way, in order to test the influence of the identified
21 factors on the observed responses [3, 4]. The statistical analysis called the
22 ANalysis Of Variance (ANOVA) of a dataset from an experimental design,
23 was historically formalized in 1925 by Fischer with the first edition of [5]. The
24 aim of this analysis is to split the total variance of the measured responses

1 into several sources of variances [6] However, the repeatability error may alter
2 the results of ANOVA, even when the factors have a relevant influence on the
3 response(s) of a studied system. Indeed, it may become more complicated to
4 differentiate which part of the variance is explained by the factors and which
5 part is due to the repeatability error.

6 In spectroscopy, measurement error may be due to both the variation of
7 the measurement conditions (e.g. angle of view, reference, sensor tempera-
8 ture, etc) and the physical properties of the measured sample (e.g. particle
9 size, roughness, physical texture, etc). In addition, a change in sample chem-
10 ical composition modifies the light diffusion process. Measuring non-uniform
11 samples adds a major difficulty compared to homogeneous powdered samples
12 [7]. For example, variations in oil, water, fiber, protein, and mineral content
13 of the same individual change the nature of the light-scattering medium.
14 These chemical and physical variations can lead to unwanted variations in
15 spectra. These variations can for example modify the spectra by adding
16 them with baselines, which can be horizontal lines or slopes, even curves.
17 When the changes are proportional to the intensity value, the effect is called
18 multiplicative effect. These phenomena are well known in spectroscopy [8, 9]
19 and are one of the sources of repeatability error.

20 Some sources (sensors, operator, etc.) of errors can be described as fac-
21 tors. If we consider these factors should have negligible effects on the total
22 variance, we don't account for these factors in the analysis of variance and
23 they are referred to as hidden. Some of these factors can be described as
24 a random factor. In the case of univariate data, these random factors can
25 be accounted for, using adapted methods such as Generalized Linear Mixed

1 Models (GLMMs) [10]. In the case of multivariate data, solutions lie in (1)
2 reducing dimensionality or (2) modelling multiple responses (e.g., longitudi-
3 nal analysis by using GLMMs; [11]). However, GLMMs are unsuitable in
4 most cases, especially with spectral data [12].

5 For multivariate data, chemometric methods can be used to take advan-
6 tage of the DoE [13]. Among those methods, the Analysis of Variance -
7 Simultaneous Component Analysis (ASCA) [14] is the most used, but the
8 repeatability error is not taken into account in ASCA. This method system-
9 atically averages the repetitions of measurements on an individual. Aver-
10 aging observations helps to reduce the repeatability error but requires in-
11 creasing the number of measurements. Climaco-Pinto et al. [15] propose a
12 method consisting in reducing the residual variance after each step of the
13 ANOVA-PCA. This method makes the separation of factors clearer by re-
14 ducing within-variance. This method does not focus on the repeatability
15 error.

16 In this study, a method named Reduction of the error of repeatability-
17 Analysis of variance-Simultaneous Component Analysis (REP-ASCA) is pro-
18 posed to reduce repeatability error by adding additionnal dataset. REP-
19 ASCA is evaluated on NIR spectral data collected on coffee beans.

20 The objectives of this paper are: (1) to suggest a method (REP-ASCA)
21 to reduce the repeatability error (2) to apply REP-ASCA to real data from
22 NIR spectroscopy (3) to study the influence of the number of measurement
23 in the additionnal set on the reduction of the repeatability error.

1 **2. Theory**

2 *2.1. ASCA*

3 Multivariate data can be represented as a matrix \mathbf{X} of size (N, P) where
4 N is the number of observations and P the number of the measured variables.
5 The analysis of variance establishes a model in order to decompose the matrix
6 \mathbf{X} into different sources of variation. By doing so, relationships between the
7 factors and the measured variables can be studied.

8 The ASCA method includes of three main steps: The first one consists
9 of partitioning the observation matrix \mathbf{X} into matrices associated with the
10 effect of the studied factors and their interactions. For example, in a DoE
11 involving two studied factors A and B, \mathbf{X} can be written according to the
12 ASCA approach as follows:

$$\mathbf{X} = \mu + \mathbf{X}_A + \mathbf{X}_B + \mathbf{X}_{AB} + \mathbf{E} \quad (1)$$

13 With μ the mean of all observations, \mathbf{X}_A and \mathbf{X}_B the terms associated
14 respectively to the main effects of the factors A and B, \mathbf{X}_{AB} , the interaction
15 term between A and B and \mathbf{E} the residuals. Matrices \mathbf{X}_A , \mathbf{X}_B and \mathbf{X}_{AB}
16 contain identical replicates of the average spectra collected at each level of
17 the particular factor (A, B and AB, respectively). For example, for a given
18 level of the factor A, all the rows of \mathbf{X}_A corresponding to this level will contain
19 the same spectrum.

20 For a complete and balanced design, factors A, B and AB are indepen-
21 dent of one another [16]. More specifically, matrices \mathbf{X}_A , \mathbf{X}_B and \mathbf{X}_{AB} are
22 all orthogonal to one another [17]. Thus: $var(\mathbf{X}) = var(\mathbf{X}_A) + var(\mathbf{X}_B) +$
23 $var(\mathbf{X}_{AB}) + var(\mathbf{E})$. The term $var(\mathbf{E})$ contains the within-level variance and

1 the other terms the between-level variance for the corresponding factor. In
 2 practice, the variances of the factors are replaced by the sum of the squares of
 3 the elements of the associated matrix. For example, with $\|\mathbf{X}_A\|^2$ correspond-
 4 ing to the square Frobenius' norm of \mathbf{X}_A , the variance of \mathbf{X}_A is replaced by
 5 [18, 17, 19]:

$$\text{SSQ}(\mathbf{X}_A) = \|\mathbf{X}_A\|^2 \quad (2)$$

6 The second step of ASCA tests if factors have a significant effect on
 7 the total variance. To do that, a permutation test is performed and a p-
 8 value is estimated by comparing variance value on the original data with the
 9 distribution of variance obtained by permutations. The permutation test can
 10 be performed for data which does not necessarily respect the conditions of
 11 normality [12]. For a factor i and its associated matrix \mathbf{X}_i , the p-value is
 12 calculated according to the following equation [20]:

$$\text{p-value}(\mathbf{X}_i) = \frac{\text{nbr}(\text{SSQ}(\mathbf{X}_{i,perm}) \geq \text{SSQ}(X_i))}{k} \quad (3)$$

13 Where $\text{nbr}()$ calculates the number of occurrences, k is the number of
 14 permutations and $\mathbf{X}_{i,perm}$ the matrix obtained after a random row permuta-
 15 tions. The p-value is then calculated by counting the number of cases where
 16 the variance of the studied factor is lower than the variances resulting from
 17 the permutations. Indeed, by doing so, the effect of the studied factor is
 18 compared to its distribution under the null hypothesis as estimated by the
 19 permutations.

20 For the last step, a Simultaneous Component Analysis (SCA) [21, 22] is
 21 performed on all the terms of the matrices belonging to equation 1. Thus,

1 for a given factor i , the decomposition of the corresponding matrix \mathbf{X}_i is
2 written:

$$\mathbf{X}_i = \mathbf{T}_i \mathbf{P}_i^t \quad (4)$$

3 Where \mathbf{T}_i and \mathbf{P}_i correspond respectively to the scores and the loadings of
4 the principal components. \mathbf{T}_i and \mathbf{P}_i are respectively matrices of size (N,L)
5 and (P,L) with L is the number of levels of factor i minus one.

6 This analysis reduces the representation space. These loadings define a
7 subspace spanned by \mathbf{X}_i and highlight the spectral directions related to the
8 factors studied. The scores are the new coordinates of the observations on
9 these principal components.

10 *2.2. Repeatability error*

11 When a sample measurement is repeated several times, measurements
12 are never identical. Observed variation between measurements is due to the
13 repeatability error. When measurements are performed on different samples,
14 variance differences are due to the variance between samples but also con-
15 tain repeatability error. The greater the repeatability error, the greater the
16 within-variance increases at the expense of between-variance. In fact, the
17 variance of the factor becomes smaller relatively to the total variance. This
18 tends to decrease the significance of the factors.

19 Some of the error can be due to a random phenomenon whose average is
20 null. In ASCA, the repeatability error is reduced because repeated measures
21 are systematically averaged over every sample. As a consequence, when the
22 repeatability error is high, one solution is to reduce its variance by increasing

1 the number of observations. As a result, the number of acquisitions increases
 2 considerably. In addition, it is impossible to link the repeatability error
 3 to a specific factor in ASCA. Indeed, it is a nested factor [23]. The solution
 4 proposed in this study is to use an additional set of repeated measures to take
 5 away the repeatability error of the observations before analysis of variance.

6 [24]

7 2.3. Proposed method: REP-ASCA

8 The method proposed in this study consists in identifying the spectral
 9 directions responsible for the repeatability error and removing them from
 10 \mathbf{X} . By definition, \mathbf{X} can be represented in two dual spaces: the column-
 11 space of dimension R^N spanned by the observations and the row-space of
 12 dimension R^P spanned by the variables. In ASCA, the decomposition part
 13 of the variance is performed in the column-space. Equation 1 can be written
 14 as a series of orthogonal projections [17]:

$$\mathbf{X} = \mathbf{M}_1\mathbf{X} + \mathbf{M}_A\mathbf{X} + \mathbf{M}_B\mathbf{X} + \mathbf{M}_{AB}\mathbf{X} + \mathbf{E} \quad (5)$$

15 Where each matrix \mathbf{M}_i is an orthogonal projection matrix in the column-
 16 space (R^N) defined by the following equation:

$$\mathbf{M}_i = \mathbf{D}_i(\mathbf{D}_i^t\mathbf{D}_i)^{-1}\mathbf{D}_i^t \quad (6)$$

17 Where \mathbf{D}_i called the dummy-matrix or design matrix is a binary matrix
 18 encoding the class belonging of the factor i .

19 The variance of repeatability error can be included into all factors. This
 20 variance cannot be expressed into a unique factor and therefore represented

1 by an orthogonal projection. This error cannot be reduced in the column-
2 space.

3 The spectral directions of the repeatability error can be described in the
4 row-space through the space W spanned by the repeatability error. Let \mathbf{X}_s
5 be a matrix containing the repetitions of measurements performed on a set
6 of samples representative of \mathbf{X} . Each repeated set is centered and all the
7 centered sets are stored in a matrix \mathbf{W} . The matrix \mathbf{W} contains only the
8 effects related to the repetition of measurements. A Principal Component
9 Analysis (PCA) is performed on \mathbf{W} and provides loadings \mathbf{P} which define the
10 repeatability error subspace. The matrix \mathbf{X} is then projected orthogonally
11 to \mathbf{P} using the following formula:

$$\mathbf{X}_\perp = \mathbf{X}(\mathbf{I} - \mathbf{P}\mathbf{P}^t) \quad (7)$$

12 ASCA is then performed on \mathbf{X}_\perp instead of \mathbf{X} . This method contains two
13 parameters to adjust: the number k of dimensions removed by the projection,
14 i.e. the number of rows of \mathbf{P} and the minimum number of spectra required
15 in \mathbf{W} , i.e. the number of rows of \mathbf{W} .

16 The setting of the parameter k can be done according to a set of criteria.
17 First, the scree plot of the PCA eigenvalues can be examined, and the value
18 of k that corresponds to a net break is the selected. Another criterion is
19 to study the evolution of the ratio between-variance/total-variance for the
20 design factors and identify a break, a plateau. When the data are signals as it
21 is the case with spectral data, studying shape of loadings (break, baseline) can
22 provide information about the nature of repeatability errors and thus allow
23 to adjust the k parameter accordingly. The REP-ASCA method assumes

1 that the repeatability error is structured, i.e. a large part of this error has a
2 subspace of limited size, which can be removed by orthogonal projection into
3 the row-space. When the number k is too high, the loadings become noisy.

4 When building \mathbf{W} , that is to say acquiring \mathbf{X}_s , it is necessary to check
5 that enough repeated measures are available to describe the repeatability
6 error. This can be verified by monitoring the evolution of the structure
7 of the subspace spanned by \mathbf{W} . Let \mathbf{J}_n be the inertia matrix of \mathbf{W} when
8 it contains n spectra. The RV coefficient [25, 26] between \mathbf{J}_n and \mathbf{J}_{n-1} is
9 calculated for each addition of elements. The evolution of this coefficient as
10 a function of n reflects the evolution of the subspace structure generated by
11 \mathbf{W} . In practice, the filling of \mathbf{W} is stopped when this coefficient no longer
12 changes.

13 **3. Materials**

14 *3.1. Experimental data*

15 The REP-ASCA was tested on data resulting from near infrared spec-
16 troscopy of coffee beans [18]. The aim of the experiment was to study species
17 and roasting times influences on spectra. Spectra were recorded over a spec-
18 tral range of 1000 nm to 2500 nm with a resolution of 1.2 nm and trans-
19 formed in pseudo-absorbance ($-\log 1/R$). Initially, ten geographical origins
20 were compared representing a total of 800 spectra. For this study, only one
21 origin (80 spectra) was selected.

22 Two species of coffee (Robusta and Arabica) were roasted during four
23 different periods (0mn, 25mn, 50mn, 75mn) at a constant temperature of
24 180 °C. Seven successive spectral measurements were collected on each sam-

1 ple. This DoE yielded a matrix \mathbf{X} of 56 spectra (2 species x 7 spectral
2 measurement x 4 duration times). In addition, 3 repetitions of measure-
3 ments have been made on each of the 8 samples leading to 24 spectra in the
4 matrix \mathbf{X}_s .

5 3.2. REP-ASCA model

6 The matrix \mathbf{X} was projected orthogonally to the k first loadings extracted
7 from X_s yielding \mathbf{X}_\perp . The value of k was varied from 0 (classical ASCA)
8 to 12. In accordance with ASCA modelling, \mathbf{X}_\perp was decomposed into the
9 following terms:

$$\mathbf{X}_\perp = \mu + \mathbf{X}_{species} + \mathbf{X}_{time} + \mathbf{X}_{speciesXtime} + \mathbf{E} \quad (8)$$

10 With: $\mathbf{X}_{species}$ and \mathbf{X}_{time} factors related respectively to species and roast-
11 ing duration time. $\mathbf{X}_{speciesXtime}$ is the interaction term and \mathbf{E} the residuals.

12 4. Results and discussion

13 4.1. Preliminary analysis of spectra

14 :

15 The \mathbf{X} spectra are shown in 1a. We can observe some characteristic
16 peaks of coffee spectra according to the literature [27, 28, 29]. Positions and
17 characteristics of these peaks are presented in Table 1:

18 However, some differences between spectra can be observed. There is
19 a shift in the baseline. This shift increases with the wavelength. This is
20 frequently observed in NIR spectroscopy when the spectra are acquired in

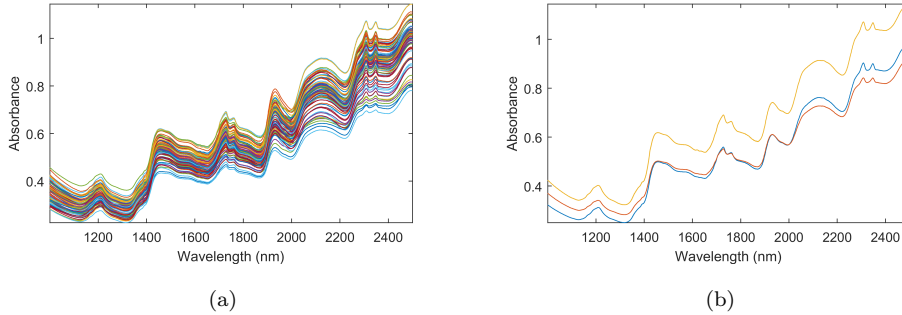


Figure 1: (a) All spectra forming by \mathbf{X} , (b) 3 spectra measured on one sample

Table 1: Position and peak characteristic of coffee spectra

Wavelength (nm)	Bond vibration	Assignment
1212	C–H stretching 2 nd overtone	Caffeine
1450	C–H stretching 1 st overtone	H ₂ O and Cellulose
1725/1765	C-H stretching 1 st overtone CH ₂	Fatty acids
1940	O–H stretching + O–H deformation	H ₂ O
2080-2150	C=O (1 st overtone) and O–H (combination)	Protein/Fatty acids/Carbohydrates
2309/2350	C–H + C=C (combination)	Carbohydrates : Cellulose

1 reflectance [30, 31, 9, 32]. This phenomenon is due to differences in the light
 2 scattering.

3 Figure 1 b shows three spectra of \mathbf{X}_s measured on the same sample and
 4 illustrating the repeatability error. These spectra present the same kind of
 5 variation observed on \mathbf{X} spectra (Fig. 1a). A vertical translation as a shift
 6 in a constant baseline appears between spectra at all wavelengths. This is
 7 because each spectrum was measured on a different area of the same sam-
 8 ple. In each area, the reflectance intensity varies according to the geometric
 9 configuration of the measured grains in relation to the sensor. These aspects

1 are also observed over the whole databases \mathbf{X} and \mathbf{X}_s .

2 The baselines also vary in slopes (see Fig. 1b). This is due to a variation
3 of particle size and surface roughness of the product [30]. It can therefore be
4 assumed that the internal structure of coffee beans varies within the same
5 sample.

6 A more detailed analysis of 1b shows three distinct spectral regions for
7 baseline variations:

- 8 • Between 1000 nm and 1400 nm, a negative slope is observed, which
9 varies for each repetition.
- 10 • Between 1400 nm and 2000 nm, baselines are increasing, but appear to
11 have identical slopes to ones observed at other spectral ranges.
- 12 • Between 2000 and 2500, baselines are increasing, but appear to have
13 different slopes. It can be hypothesized that all of these baseline vari-
14 ations are due to several scattering regimes (e.g., Mie and Rayleigh).
15 These scattering regimes are governed by the relationship between par-
16 ticle size and light wavelength [7].

17 All these baselines deform the information contained in the spectra. These
18 undesired effects should be removed to better describe the factors.

19 4.2. ASCA on \mathbf{X}

20 The Explained variance and the p-value for each factor are obtained from
21 ASCA on \mathbf{X} and summarized in Table 2. We can then observe that the
22 largest source of variance is explained by the residuals. 82.01% of the total
23 variance could not be related to a studied factor. 5.93% of the variance is

Table 2: Explained variances and p-values of the factors

Factor	Explained Variance (%)	p-value
$\mathbf{X}_{species}$	5.93	0.056
\mathbf{X}_{time}	8.61	0.128
$\mathbf{X}_{species \times time}$	3.45	0.548
\mathbf{E}	82.01	

1 explained by the species, 8.61% of the total variance is explained by the
 2 roasting duration time and 3.45% by the interaction term.

3 In addition, the p-values obtained after the permutation test, show that
 4 none of the factors is statistically significant at 5%. These results can be
 5 explained in part by a large part of variance not explained by the model.
 6 The factor variances are too low in comparison to the residual variance.
 7 Here, the case study shows that the repeatability error is so large that none
 8 of the identified factors are expressed. Within-variance of repeated measures
 9 are more important than variances between-level of factors. It is then very
 10 difficult in this particular case to draw conclusions with ASCA.

11 4.3. REP-ASCA

12 4.3.1. Selection of the number k of component removed

13 The parameter k , which corresponds to the number of components of \mathbf{W}
 14 to remove, must be set. For this purpose, consequences on ASCA performed
 15 on \mathbf{X}_\perp in terms of explained variances and p-values are studied according to
 16 k .

17 Fig. 2a shows, that part of variance carried by \mathbf{E} decreases up to the

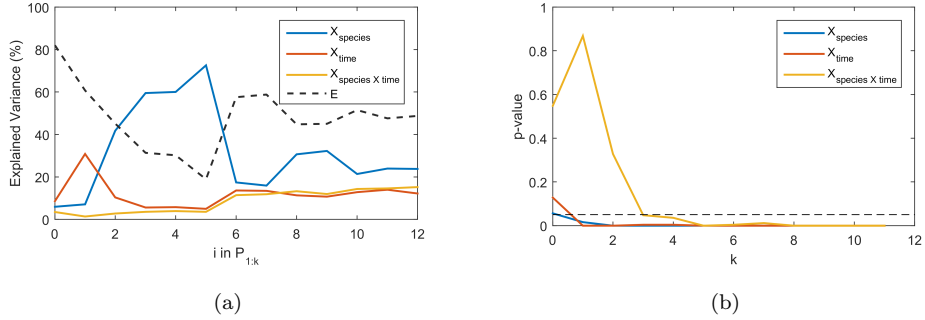


Figure 2: Evolution of: (a) explained variances by factors, (b) p-value for each factor ; according to the projection number.

1 5th component. Then, it subsequently increases. Inversely, the term $\mathbf{X}_{species}$
 2 explains more and more variance and that up to the 5th component. The ef-
 3 fects related to the repeatability error hide the expression of the term $\mathbf{X}_{species}$
 4 to the 5th component. For the \mathbf{X}_{time} term, projecting orthogonally to the
 5 1st component increases the variance part from 8.61% to 30.81%. This large
 6 increase shows that the 1st component of \mathbf{W} mainly masked variance due to
 7 roasting time. Nevertheless, the variance explained by this term decreases
 8 thereafter. This means that from the second component, the repeatability
 9 error is very much related to roasting. For the interaction term $\mathbf{X}_{species \times time}$
 10 explained variance is very small. It increases very slowly as the repeatability
 11 error is removed.

12 By reducing the variance due to the repeatability error, the variances of
 13 the factors can be modified in different proportions. Information is extracted
 14 in different proportions from the different factors. The significativity of each
 15 factor can then potentially change. The evolution of p-value (Fig. 2b) shows
 16 that the repeatability error masks the effects of the factors. By removing the

1 information carried by the first component of \mathbf{W} (i.e., projection orthogo-
2 nally to the 1st component of the repeatability error), the \mathbf{X}_{time} and $\mathbf{X}_{species}$
3 factors become both significant, in agreement with what reported in the orig-
4 inal study [18]. The interaction term becomes significant after a projection
5 orthogonal to the first three components.

6 *4.3.2. Analysis of the repeatability error*

7 Let's examine in detail the results of REP-ASCA for $k=5$. The sources
8 of repeatability can be explored through the analysis of the loadings of the
9 first 5 component (Fig 3).

1 Fig. 3a shows the percentages of explained variance by component. The
2 first component explains 91.6% of the total variance. The loadings of the first
3 component (Fig. 3b) are all positive. This indicates a systematic variation
4 across all variables. The shape of these loadings is the same as the aver-
5 age spectrum of the individuals of \mathbf{X}_s . This reflects the vertical translation
6 observed in (Fig. 1a).

7 The second component explains 6.53% (Fig. 3a) of the total variance.
8 Surprisingly, the loadings of the second component (Fig. 3c) are similar to
9 the loadings of the first component (Fig. 3b). However, slight differences
10 can be observed: first, the loadings of the second component are centered
11 on zero. This indicates an overall slope variation. Secondly, these loadings
12 show a zero slope between 1000 nm and 1400 nm, whereas it is negative on
13 the loadings of the first component of this same spectral range. As a result,
14 the deflection is not homogeneous over all wavelengths.

15 The third component explains 1.13% (Fig. 3a) of the total variance. The
16 loadings of the third component (Fig. 3d) show two slopes: one between
17 1000 nm and 1400 nm and another between 2000 nm and 2500 nm. This
18 corresponds to different scattering regimes related to the size of the scattering
19 media in respect to the wavelength. This scattering phenomenon is finer than
20 what we saw before. This phenomenon results in a non-constant curving
21 baseline effect. The fourth component explains 0.2% of the total variance
22 (Fig. 3a). On the loadings (Fig. 3e), a first slope is visible between 1000 nm
23 and 1300 nm and a second between 1900 nm and 2300 nm. Unlike previous
24 components, slopes are less visible. This component is therefore less related
25 to physical phenomena. This gives information related to the chemistry of the

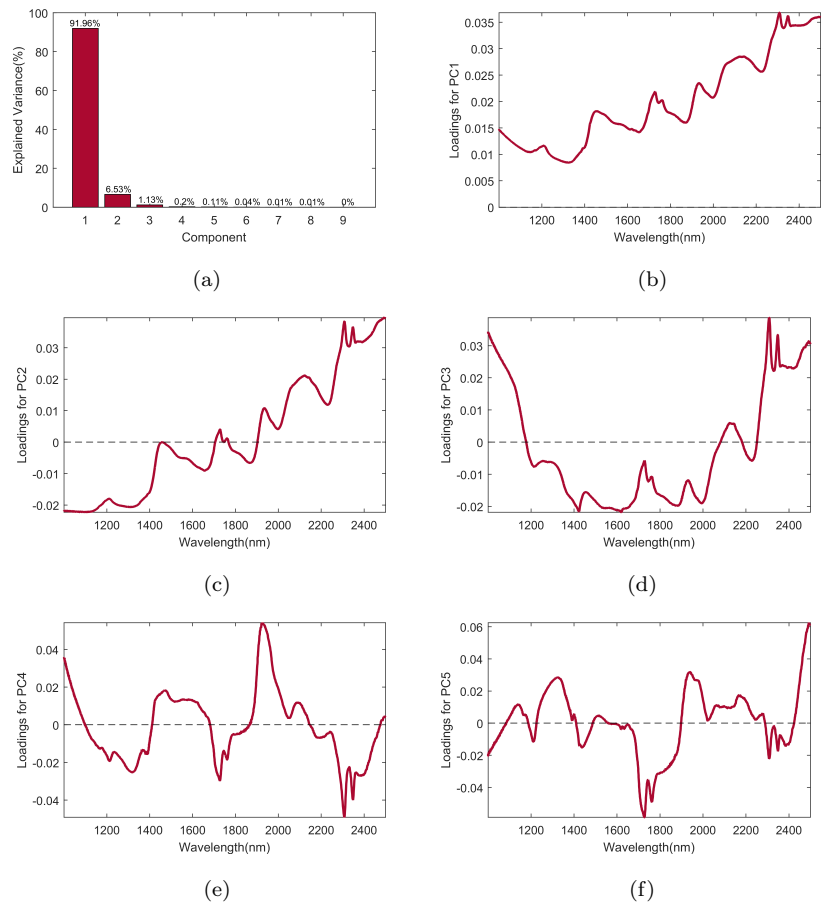


Figure 3: PCA performed on W: (a) Part of variances explained by principal components (b) Loadings on PC1, (c) Loadings on PC2 (d) Loadings on PC3, (e) Loadings on PC4, (f) Loadings on PC5.

1 medium. A characteristic peak at 1940 nm is related to the presence of water
2 [33] as a result of the combination of the O-H stretch band and the O-H₂
3 torsion band. Two peaks are located at 1710 nm and at 2307 nm respectively
4 linked to lipids and proteins content. The fifth component explains 0.11%
5 of the total variance (Fig. 3a). As shown in (Fig. 3f), the loadings of this
6 component show slopes between 1000 nm and 1300 nm and 2400 nm and
7 2500 nm. Peaks of 1730 nm and 1762 nm correspond to the lipids.

8 The description of the loadings of the first five components highlights
9 that the variance of the repeated measures is related to physical phenomena.
10 These physical phenomena induce essentially shifts in baselines on spectra.
11 In the next components, chemical information appears. The repeatability
12 error is therefore mainly of a physical nature, and to a lesser extent also due
13 to a chemical difference, changing the diffusion of light in the medium.

14 4.4. Factor analysis after REP-ASCA (with $k = 5$)

15 4.4.1. Species

16 Without or after orthogonal projection (Fig. 4a and 4b), the water con-
17 tent differentiates the two species by the peak of water located at 1930 nm.
18 A second smaller peak located at 1398 nm corresponds to the hydrogen bond.
19 It is also called weakly-hydrogen-bonded water [34]. Note that without or-
20 thogonal projection the loadings are all negative except the peak located at
21 1930 nm (Fig. 4a). After orthogonal projection according to the first 5 com-
22 ponents of \mathbf{W} , the loadings (Fig. 4b) are centered on zero. The orthogonal
23 projection therefore removed the baseline.

24 Some information that did not appear before is now apparent (Fig. 4b).
25 The three peaks at 2309 nm, 2348 nm and 2400 nm show a difference in fat

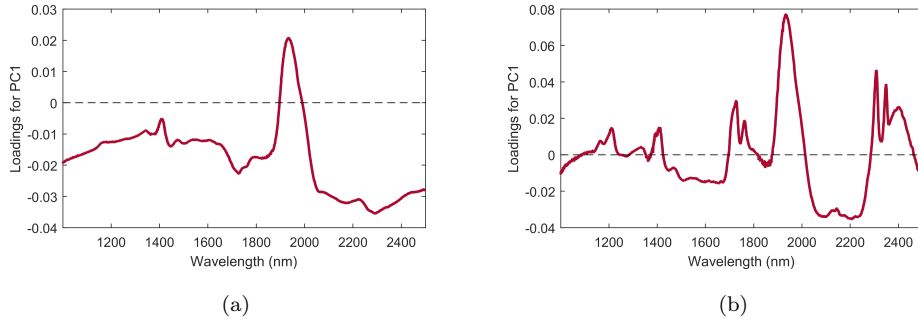


Figure 4: Loadings on PC1 of the term $\mathbf{X}_{species}$: (a) Without pretreatment (b) by using the projection orthogonally to the 5 first components of \mathbf{W} .

1 content [7]. The peaks at 1728 nm and 1763 nm reveal a strong presence
 2 of fatty acids [28]. It's known that Arabica beans has higher lipid content
 3 than Robusta ones [28]. The two species can be differentiated by the lipid
 4 content. A less important peak at 1212 nm also shows a difference in caffeine
 5 and carbohydrates. A negative plateau between 2070 nm and 2238 nm shows
 6 a difference in proteins and chlorogenic acids content.

7 4.5. Impact of the number of spectra in \mathbf{W}

8 4.5.1. Number required for \mathbf{W}

9 In this study, \mathbf{W} contains 24 spectra. This number is more than enough to
 10 describe the repeatability error. The impact of the number of observations
 11 on the structure of the subspace of \mathbf{W} is studied using the RV coefficient
 12 between the matrices \mathbf{J}_n and \mathbf{J}_{n-1} with n ranging from 4 to 24 spectra.
 13 Bootstrap is performed to check that the results are not due to the random
 14 selection of observations available.

15 The graph 5 shows the evolution of the mean and standard deviation of
 16 this coefficient according to n . Mean values are higher than 0.90 for all the

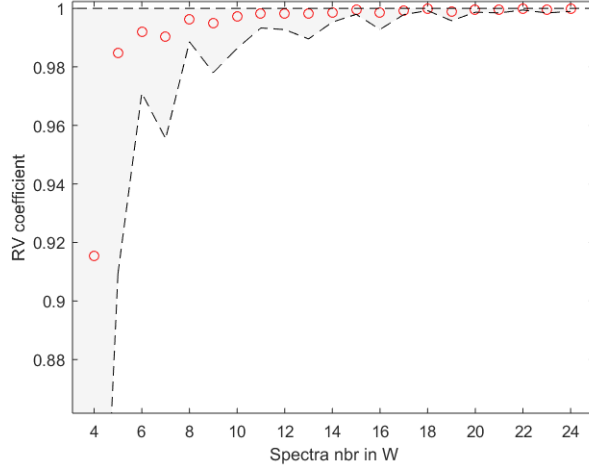


Figure 5: Evolution of the mean (red circle) and the standard deviation (grey area) of the RV coefficient between \mathbf{J}_n and \mathbf{J}_{n-1} according to the number n of spectra contained in \mathbf{W}

1 values of n . A plateau is visible with values close to 1 when n is greater than
 2 10. The standard deviation displayed is about 0.18 when $n < 6$. When n
 3 increases ($n > 11$), standard deviation gradually decreases to values below
 4 0.01. For a given n , a high value of the standard deviation means that
 5 the RV coefficient varies according to the n -th spectrum added in \mathbf{W} . The
 6 opposite is true, with a low standard deviation, the RV coefficient random
 7 varies slightly.

8 For $n < 8$, RV coefficients have low values meaning that the repeatability
 9 error is not correctly described. Additional observations must then be added.
 10 The increase of RV coefficient values with n is expected. Indeed, adding an
 11 observation has less impact in a dataset containing many observations. The
 12 plateau located from $n = 10$ shows that the structure of subspace \mathbf{W} is
 13 stabilizing. From 10 observations, adding spectra in \mathbf{W} does not provide

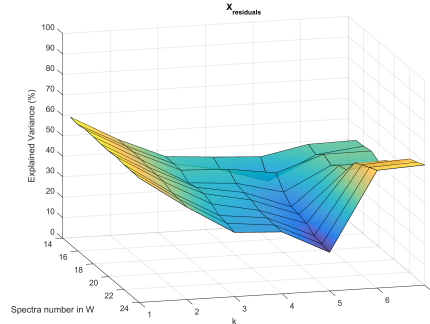


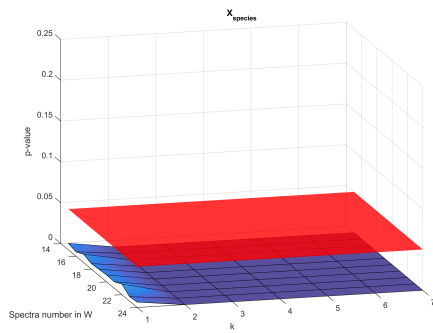
Figure 6: Explained variance of residuals according to number of projections and spectra number

1 more information on the repeatability error.

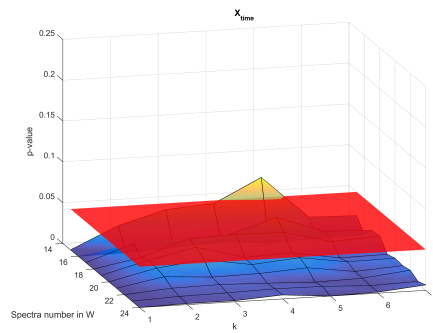
2 4.5.2. Impact on the selection of k

3 In this study, repeated measures are used to construct the \mathbf{X}_s matrix.
 4 This matrix is used to describe the repeatability error. The impact of the
 5 number of spectra used for \mathbf{X}_s matrix and the number of components used
 6 for orthogonal projection on the explained variance of residuals is illustrated
 7 in Fig. 6. The variance decreases up to the 5th component for a spectra
 8 number ranging from 15 to 24 in \mathbf{X}_s . The part of variance explained by the
 9 residuals is 37.31% and is 18.88% respectively for a number of 15 and 24
 10 spectra in \mathbf{X}_s . The fewer spectra in \mathbf{X}_s , the greater residual variance. As a
 11 result, the repeatability error is then less described.

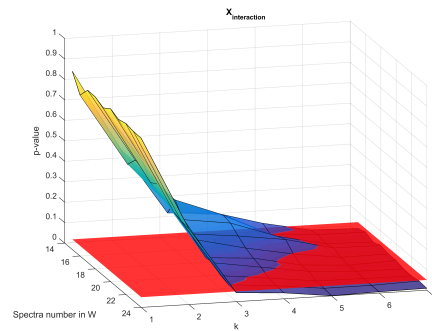
12 The p-values for each factor are presented in Fig. 7 according to the
 13 number of spectra used for \mathbf{X}_s matrix and the number of components used
 14 for orthogonal projection. For the species factor (Fig. 7a), reducing the
 15 number of spectra from 24 spectra to 15 does not have an impact on the
 16 significance even after projection because species is the predominant factor



(a)



(b)



(c)

Figure 7: Impact of the spectra number in \mathbf{X}_s and of the number k of components removed on p-values obtained for (a) Species, (b) Roasting time, (c) Interaction; The surface in red represents the value of 0.05 for the p-value.

1 in terms of portion of variance explained. For the roasting time (Fig. 7b)
2 and the interaction terms (Fig. 7c), the number of spectra in \mathbf{X}_s and the
3 number of components used for orthogonal projection affect the results of
4 the permutation tests. As presented in Fig. 7b, the roasting time term is not
5 significant after a projection orthogonal to the first five components of \mathbf{W}
6 when \mathbf{X}_s contains less than or equal to 15 spectra. For the interaction term
7 (Fig. 7c), significance may change depending on the number of spectra used
8 to build \mathbf{X}_s . Permutation tests may be strongly influenced when reducing the
9 number of spectra in \mathbf{X}_s . Here, to study the interaction term and the roasting
10 time term, it is necessary to have at least 20 spectra in \mathbf{X}_s . As a general,
11 the number of observations must be sufficient to reduce the repeatability
12 error. The more components to be removed, the higher the number of spectra
13 required to be included in \mathbf{X}_s .

14 In the previous study [18], the Standard Normal Variate transform (SNV)
15 [35] was used to reduce additive effects. This pretreatment is common in
16 spectroscopy. It corrects the constant baselines over the whole dataset but
17 SNV cannot correct non-constant (sloping or curving) baseline effects [36, 9].

18 The advantage of the proposed method is to reduce the repeatability error
19 with more accuracy. Based on the analysis of repeated measures, the com-
20 ponents of the repeatability error contain non-constant (sloping or curving)
21 baselines. Indeed, they reflect diffusion and absorption phenomena of both
22 physical and chemical nature.

1 5. Conclusion

2 In this study, a method to reduce the repeatability error (REP-ASCA)
3 was tested on spectral data containing large repeatability error. Indeed, the
4 interpretation of the factors through ASCA on such raw data (without correc-
5 tion) could mislead to wrong conclusions. The loadings of factor components
6 contain then only the information related to the repeatability error. It was
7 also illustrated how the significance of the factors can be influenced by the
8 nature of the repeatability error.

9 The proposed method reduces these effects with the help of an additional
10 set. The description of the repeatability error is then made through the study
11 of repeated measurement introduced in the additional set. This description
12 highlights the spectral regions related to the repeatability error. While re-
13 peatability error description is carried out in the column-space, reducing this
14 error is realized in the row-space.

15 The choice of the number of components to be removed is done by study-
16 ing the impact of projections orthogonally to these components on the various
17 factors of the analysis of variance. Moreover, this choice can be improved by
18 the inspection of the loadings of the main components of the repeatability
19 error. While removing the error due to repeatability, the results obtained
20 with REP-ASCA allow to highlight the spectral regions due to the factors of
21 interest.

22 REP-ASCA was tested on a dataset containing a repeatability error
23 mainly caused by samples with highly variable physical specificities and as-
24 sociated measurement errors (particle size, angle of view, etc). It would be
25 interesting to apply REP-ASCA to another dataset containing other kinds of

1 repeatability errors. In addition to reducing repeatability error, REP-ASCA
2 could remove other undesired effects (temporal effect, specific chemical con-
3 tent) by incorporating other components related to these effects. Futhermore,
4 the REP-ASCA methodology could correct any type of unwanted effects, pro-
5 vided that the spectral subspace generated by these effects can be identified.
6 A lot of methods can be used to carry out this identification [37].

7 On the practical side, adding an additional set to reduce repeatability
8 error may considerably reduce the number of repetitions per modality in the
9 set for the analysis of variance. Using this method, acquisition protocols
10 should include observations to describe these errors independently of the
11 observations for analysis of variance. This is required to make the measures
12 for the analysis of variance more robust. However, a sufficient number of
13 observations should be kept to study the significance of the factors.

14 **References**

- 15 [1] F. E. Grubbs, On estimating precision of measuring instruments and
16 product variability, *Journal of the American Statistical Association* 43
17 (1948) 243–264.
- 18 [2] J. N. Miller, J. C. Miller, *Statistics and Chemometrics for Analytical*
19 *Chemistry*, Pearson Education, 2005. Google-Books-ID: Efx77dxOC3sC.
- 20 [3] R. A. Fisher, *The design of experiments*, Oliver And Boyd; Edinburgh;
21 London, 1937.
- 22 [4] G. W. Oehlert, *A first course in design and analysis of experiments*,
23 W.H. Freeman, New York, 2000.

- 1 [5] R. A. Fisher, Statistical methods for research workers, Genesis Publish-
2 ing Pvt Ltd, 2006.
- 3 [6] H. Scheffé, The analysis of variance, The analysis of variance, Wiley,
4 Oxford, England, 1959.
- 5 [7] P. Williams, K. Norris, Near-infrared technology in the agricultural and
6 food industries., American Association of Cereal Chemists, Inc., 1987.
- 7 [8] J. Engel, J. Gerretzen, E. Szymańska, J. J. Jansen, G. Downey,
8 L. Blanchet, L. M. Buydens, Breaking with trends in pre-processing?,
9 TrAC Trends in Analytical Chemistry 50 (2013) 96–106.
- 10 [9] Å. Rinnan, F. v. d. Berg, S. B. Engelsen, Review of the most common
11 pre-processing techniques for near-infrared spectra, TrAC Trends in
12 Analytical Chemistry 28 (2009) 1201–1222.
- 13 [10] B. M. Bolker, M. E. Brooks, C. J. Clark, S. W. Geange, J. R. Poulsen,
14 M. H. H. Stevens, J.-S. S. White, Generalized linear mixed models: a
15 practical guide for ecology and evolution, Trends in Ecology & Evolution
16 24 (2009) 127–135.
- 17 [11] S. L. Zeger, K.-Y. Liang, Longitudinal data analysis for discrete and
18 continuous outcomes, Biometrics (1986) 121–130.
- 19 [12] M. Anderson, C. T. Braak, Permutation tests for multi-factorial analysis
20 of variance, Journal of statistical computation and simulation 73 (2003)
21 85–113.

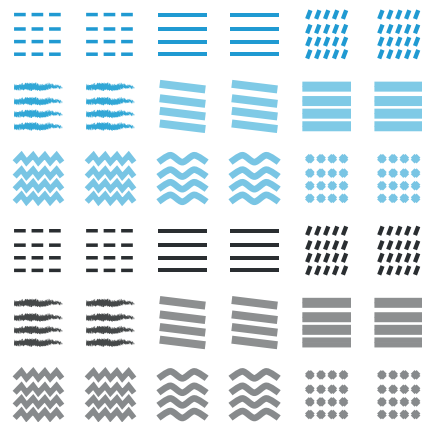
- 1 [13] R. G. Brereton, J. Jansen, J. Lopes, F. Marini, A. Pomerantsev, O. Ro-
2 dionova, J. M. Roger, B. Walczak, R. Tauler, Chemometrics in analytical
3 chemistry—part I: history, experimental design and data analysis tools,
4 Analytical and Bioanalytical Chemistry 409 (2017) 5891–5899.
- 5 [14] A. K. Smilde, J. J. Jansen, H. C. J. Hoefsloot, R.-J. A. N. Lamers,
6 J. van der Greef, M. E. Timmerman, ANOVA-simultaneous component
7 analysis (ASCA): a new tool for analyzing designed metabolomics data,
8 Bioinformatics 21 (2005) 3043–3048.
- 9 [15] R. Climaco-Pinto, A. Barros, N. Locquet, L. Schmidtke, D. Rutledge,
10 Improving the detection of significant factors using ANOVA-PCA by
11 selective reduction of residual variability, Analytica Chimica Acta 653
12 (2009) 131–142.
- 13 [16] R. G. Shaw, T. Mitchell-Olds, Anova for Unbalanced Data: An
14 Overview, Ecology 74 (1993) 1638–1645.
- 15 [17] A. K. Smilde, H. C. Hoefsloot, J. A. Westerhuis, The geometry of ASCA,
16 Journal of Chemometrics 22 (2008) 464–471.
- 17 [18] S. De Luca, M. De Filippis, R. Bucci, A. D. Magrì, A. L. Magrì,
18 F. Marini, Characterization of the effects of different roasting condi-
19 tions on coffee samples of different geographical origins by HPLC-DAD,
20 NIR and chemometrics, Microchemical Journal 129 (2016) 348–361.
- 21 [19] D. J. Vis, J. A. Westerhuis, A. K. Smilde, J. van der Greef, Statistical
22 validation of megavariate effects in ASCA, BMC Bioinformatics 8 (2007)
23 322.

- 1 [20] M. J. Anderson, Permutation tests for univariate or multivariate analysis
2 of variance and regression, *Canadian Journal of Fisheries and Aquatic*
3 *Sciences* 58 (2001) 626–639.
- 4 [21] H. A. Kiers, SCA: A Program for Simultaneous Components Analysis
5 of Variables Measured in Two Or More Populations: user’s Manual, iec
6 ProGamma, 1990.
- 7 [22] K. Van Deun, A. K. Smilde, M. J. van der Werf, H. A. Kiers,
8 I. Van Mechelen, A structured overview of simultaneous component
9 based data integration, *Bmc Bioinformatics* 10 (2009) 246.
- 10 [23] F. Marini, D. de Beer, E. Joubert, B. Walczak, Analysis of variance
11 of designed chromatographic data sets: The analysis of variance-target
12 projection approach, *Journal of Chromatography A* 1405 (2015) 94–102.
- 13 [24] ISO-5725-2 1994(en), International Organization for Standardization,
14 1994, Accuracy (trueness and precision) of measurement methods and
15 results— Part 2: Basic method for the determination of repeatability
16 and reproducibility of a standard measurement method, Standard, ????
- 17 [25] Y. Escoufier, Le Traitement des Variables Vectorielles, *Biometrics* 29
18 (1973) 751.
- 19 [26] H. Abdi, RV coefficient and congruence coefficient, *Encyclopedia of*
20 *measurement and statistics* 849 (2007) 853.
- 21 [27] L. Alessandrini, S. Romani, G. Pinnavaia, M. D. Rosa, Near infrared
22 spectroscopy: An analytical tool to predict coffee roasting degree, *An-*
23 *alytica Chimica Acta* 625 (2008) 95–102.

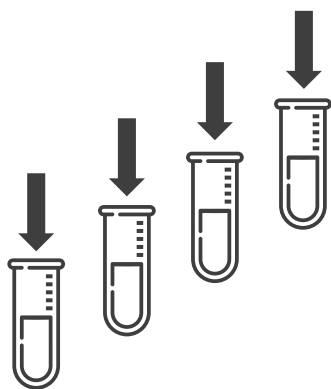
- 1 [28] I. Esteban-Díez, J. González-Sáiz, C. Pizarro, Prediction of sensory
2 properties of espresso from roasted coffee samples by near-infrared spec-
3 troscopy, *Analytica Chimica Acta* 525 (2004) 171–182.
- 4 [29] J. Ribeiro, M. Ferreira, T. Salva, Chemometric models for the quantita-
5 tive descriptive sensory analysis of Arabica coffee beverages using near
6 infrared spectroscopy, *Talanta* 83 (2011) 1352–1358.
- 7 [30] K. A. Bakeev (Ed.), *Process analytical technology: spectroscopic tools*
8 *and implemented strategies for the chemical and pharmaceutical indus-*
9 *tries*, Wiley, Chichester, West Sussex, 2nd ed edition, 2010. OCLC:
10 ocn473478595.
- 11 [31] D. L. Massart, B. G. M. Vandeginste, L. M. C. Buydens, S. de Jong, P. J.
12 Lewi, J. Smeyers-Verbeke, C. K. Mann, *Handbook of Chemometrics and*
13 *Qualimetrics: Part A, Applied Spectroscopy* 52 (1998) 302A.
- 14 [32] J. Workman, A. W. Springsteen (Eds.), *Applied spectroscopy: a com-*
15 *pact reference for practitioners*, Academic Press, San Diego, 1998.
- 16 [33] W. A. Luck, *Structure of water and aqueous solutions*, Verlag Chemie,
17 1974.
- 18 [34] A. A. Gowen, R. Tsenkova, C. Esquerre, G. Downey, C. P. O'Donnell,
19 Use of near Infrared Hyperspectral Imaging to Identify Water Matrix
20 Co-Ordinates in Mushrooms (*Agaricus Bisporus*) Subjected to Me-
21 chanical Vibration, *Journal of Near Infrared Spectroscopy* 17 (2009)
22 363–371.

- 1 [35] R. J. Barnes, M. S. Dhanoa, S. J. Lister, Standard Normal Variate
2 Transformation and De-Trending of Near-Infrared Diffuse Reflectance
3 Spectra, *Applied Spectroscopy* 43 (1989) 772–777.
- 4 [36] R. P. Cogdill, C. R. Hurburgh, G. R. Rippke, S. J. Bajic, R. W. Jones,
5 J. F. McClelland, T. C. Jensen, J. Liu, others, Single-kernel maize
6 analysis by near-infrared hyperspectral imaging, *Transactions of the*
7 *ASAE* 47 (2004) 311.
- 8 [37] J.-M. Roger, J.-C. Boulet, A review of orthogonal projections for cali-
9 bration, *Journal of Chemometrics* (2018) e3045.

EXPERIMENTAL DESIGN



REPEATED MEASURES



P_{\perp}



ASCA



PCA



LOADINGS

