

A new approach for modelling sensor based data

Thomas Skov* & Rasmus Bro

Department of Food Science (IFV), The Royal Veterinary and Agricultural University (KVL)
Rolighedsvej 30, DK-1958 Frederiksberg C, Denmark

*E-mail: thsk@kvl.dk

Abstract

Data from sensor based analytical instruments (e.g. electronic nose) can be arranged in three dimensions as Sample \times Time \times Sensor. Due partly to limitations in software, partly to general practice, the dimensions of the data are often reduced in the time mode and traditional two-way chemometric models (e.g. principal component analysis – PCA) are used for the exploration of the data. Thus the internal relationship between the different modes (samples, time and sensors) is destroyed and the result can be an incomplete data analysis.

A new approach to handle the time information is introduced combining new and more advanced multi-way chemometric models with traditional pre-processing techniques as an alternative to PCA-like methods. As the pre-processing is an essential but often time-consuming part of the data analysis, a semi-automated approach has been used to make comparison of multiple analyses simple.

The main feature of this new approach is the exploration of the total time profiles such that potentially relevant information is not discarded by feature extraction before the actual data analysis. Multi-way modelling is evaluated and for sensor based data, it is shown that the so called PARAFAC2 multi-way model that handles shifted profiles offers some advantages compared to PCA and alternative multi-way methods.

The principles are exemplified with an example from the licorice industry. An algorithm and graphical user interface has been made available at www.models.kvl.dk to ease testing this new approach.

1

2 **Keywords:** Electronic nose, data analysis, feature extraction, PARAFAC2

3

1 **1. Introduction**

2 A number of instrumental analytical techniques are used in food research as well as in
3 laboratories of the food-associated industry (e.g. spectroscopy, gas chromatography,
4 electronic noses/tongues). The raw data from these analyses are rarely used directly. Rather
5 the data are manipulated using for example baseline correction, peak identification and
6 feature extraction (slope, area etc.). These operations take time and offer no guarantee that the
7 data is properly pre-processed and that the available and useful information is extracted.
8 Many of the analytical techniques also require a non-trivial and often expert based post-
9 processing of the data, which in many situations takes unnecessarily long time and uses an
10 excess of manpower (e.g. identification and validation of peaks in gas chromatography).
11 Often these techniques are combined with traditional univariate exploration of the data with
12 extensive focus on a small (local) part of the total variables.

13
14 The electronic nose is a more or less automated analytical method that is based on headspace
15 generation (volatile compounds) from a sample. The samples can often be analysed with little
16 or no prior sample preparation, which makes the electronic nose very suitable for at-line
17 analyses. The commercial electronic noses are often based on sensors with different
18 selectivity and sensitivity – a hybrid sensor array system, which gives a very powerful
19 analytical instrument [1]. The low selectivity of the sensors makes the electronic nose more
20 appropriate for classification and separation tasks, than more specific analyses on the volatile
21 compounds in the sample. This is especially so for complex food analyses where flavour
22 plays an important part; e.g. coffee powders [2], olive oils [3], apple juices [4] and cheddar
23 cheese [5].

24
25 The pre-processing step and pattern recognition system (e.g. principal component analysis –
26 PCA) is often an integrated part of the electronic nose and thus a fast (initial) evaluation of
27 data is possible. However, the choice of pre-processing steps and pattern recognition system
28 is in general very restricted. The raw data from electronic noses is often arranged in three
29 dimensions (e.g. *Samples* \times *Times* \times *Sensors*) (Figure 1) and thus can not be directly
30 analysed with the two-way based integrated software. As a result of these limitations
31 important information from the raw data may well be missed or changed leading to a
32 suboptimal final interpretation. In this paper it is shown that it is not necessary to transform
33 the informative measured time profile into univariate features before meaningful data
34 analysis can be performed. By maintaining the time information, it is anticipated that at least

1 in some situations more quantitative and qualitative information can be extracted from the
2 data.

3

4 The aim of this article is to show the steps in the data analysis of sensor based data and
5 provide new models to such data. The steps are divided into a *pre-processing* function
6 combined with a mathematical *processing* method e.g. multivariate data analysis
7 (chemometrics) and a *post-processing* step.

8

9 An illustration of the efficiency of selecting the appropriate combination of pre-processing
10 steps and a suitable chemometric model will be given using an example from a concrete
11 situation.

1 **2. Theory**

2 *2.1. Data analysis*

3 Data-analysis is usually done on pre-processed data. The rationale for pre-processing is to get
4 the most efficient representation of the data, especially efficient with respect to the
5 subsequent modelling of the data. The combination of pre-processing steps is evaluated
6 through appropriate multivariate models, which again are often evaluated by means of
7 graphical illustrations, quantitative classification or regression models. Thus these three steps
8 strongly interact in finding the appropriate analysis approach for a given data set.

9
10 The interchanging of pre-processing steps is not straightforward as one option may influence
11 the other – e.g. if the raw data needs to be scaled to give a reasonable solution this does not
12 necessarily imply the same need for data transformed by taking first derivatives of the time
13 profile.

14
15 *2.2. Pre-processing*

16 This section describes some of the most common pre-processing techniques that are used for
17 data from sensor based analytical instruments.

18
19 **Baseline correction**

20 Several baseline correction methods can be used to pre-processes the typical sensor response
21 curve [6-9]. A frequently used baseline correction is the difference between the maximum
22 signal (or minimum) and the baseline signal, ΔS or $S-S_0$. The following baseline correction
23 methods for MOS sensors (Table 1) have been suggested [9,10]. In Figure 2 an example of
24 baseline correction using the fractional approach - $(S-S_0) / S_0$ is shown together with the
25 corresponding raw data.

26
27 **Transformation**

28 The structure of sensor based data often makes a “transformation” step feasible. As examples
29 of transformations, for sensor based data, the use of logarithmic transformations has been
30 proposed and also the use of first or second derivatives may result in an easier interpretation
31 of the data by removing trivial but disturbing baseline offsets or slopes [7,8].

32
33 **Feature extraction**

1 When analysing the sensor response curve there are several ways to extract the information
 2 that is located in different parts of the response curve (e.g. adsorption, maximum/minimum or
 3 desorption parameters) [7]. A simple technique to reduce the number of features is to select a
 4 subset of the available features (e.g. 10 of the 240 measurements of the sensor signal). The
 5 discarded features should then either contain no valuable (or little) information in the specific
 6 application or have a large correlation to some of the chosen features [8].

7
 8 Several extraction features (individual/combinations) should be investigated to estimate their
 9 discrimination power. In general though, the quality of a feature extraction depends highly on
 10 both the application and the instrument.

11 **Centering and scaling**

12 Centering aims at removing constant terms in the data in order to make the data compatible
 13 with the model. Centering across the first mode (e.g. the sample mode) of a three-way array
 14 can be done by matricizing the array to an $I \times JK$ matrix and then center as in ordinary two-
 15 way analysis, where the column mean is subtracted from every column element [11]:

$$16 \quad x_{ijk}^{Centered} = x_{ijk} - \frac{\sum_{i=1}^I x_{ijk}}{I}$$

17
 18
 19
 20 Scaling provides that each variable is scaled to an equal variation and thus, each variable will
 21 have the same opportunity to affect the model [11]. Scaling within the third mode (e.g. the
 22 sensor mode) is done by scaling each of the sensor variables to a unit mean square [11]:

$$23 \quad x_{ijk}^{Scaled} = \frac{x_{ijk}}{\sqrt{\sum_{i=1}^I \sum_{j=1}^J \frac{x_{ijk}^2}{IJ}}}$$

24 **2.3. Processing/modelling step**

25 **2.3.1. Notation**

26 The notation and terminology to describe matrices (2-way arrays) and higher order arrays is
 27 adapted from [12] and [13]. Scalars are indicated with lower-case italics (e.g. x_{ijk}) and vectors
 28 with bold lower-case characters (e.g. \mathbf{y}). In chemometrics, the data is often described with the

1 symbol \mathbf{X} . Ordinary two-way arrays (matrices) are denoted \mathbf{X} (bold-face) whereas higher
2 order arrays are denoted $\underline{\mathbf{X}}$. The ijk th element of a three-way array $\underline{\mathbf{X}}$ is denoted x_{ijk} where the
3 indices run as follows: $i = 1, \dots, I; j = 1, \dots, J; k = 1, \dots, K$.

4
5 The three-way array will be denoted $\underline{\mathbf{X}} (I \times J \times K)$ where I is the number of samples, J the
6 number of time measurements and K the number of sensors. Thus the ijk th element in $\underline{\mathbf{X}}$
7 corresponds to the measured sensor signal from sample i , measured at time j , with sensor k .

8 9 2.3.2. PCA and PARAFAC models

10 PARAFAC1 (often called PARAFAC) was developed in 1970 [14] in psychometrics, but has
11 found widespread use in chemometrics. To explain the PARAFAC1 model it is feasible to
12 start out explaining some of the basic properties of the PCA model. PCA is a *bilinear* model.
13 Thus, in PCA each component is the outer product of two vectors (scores pertaining to
14 samples and loadings pertaining to variables). The model is linear in the scores for fixed
15 loadings and vice versa and therefore bilinear. If a given data set can be modelled by an R -
16 component PCA model, then the $I \times J$ matrix \mathbf{X} can be approximated as \mathbf{TP}^T where \mathbf{T} is the
17 score matrix ($I \times R$) and \mathbf{P} the loading matrix ($J \times R$). One particular view of PCA will be
18 helpful in the following. If the rows of \mathbf{X} are different samples and the columns are holding
19 the measurements at different (J) times then the PCA model specifically states that *any*
20 measured time profile – row of \mathbf{X} – can be modelled as a linear combination of a few (R)
21 loadings. Thus PCA tries to describe the different time profiles as different linear
22 combinations (defined by the scores) of a few underlying latent variables (the loadings). This
23 is an essential and important property of PCA in this context. It is illustrated in Figure 3A
24 where it can be seen that when the same phenomenon occurs in different proportions in
25 different samples, this is easily handled by a one-component PCA model.

26
27 Sometimes though, the same phenomenon does not only change in magnitude but also in
28 shape. E.g. the maximum of the peak shifts position (Figure 3B). This is a dramatically
29 different situation because the variables (times in this case) essentially change meaning from
30 sample to sample. PCA and most other multivariate methods can not handle such shifts in an
31 efficient manner. Every type of shift must basically be modelled by a separate component
32 which is not a meaningful and parsimonious way to model the data. The PARAFAC2 model
33 to be discussed will be shown to be able to handle this type of variations.

1 If data is arranged in a three-way array (e.g. Samples \times Times \times Sensors) it is necessary to
2 matricize (i.e. unfold) the data in order to fit a PCA model. This can be done as shown in
3 Figure 4. A PCA model of such unfolded data is often referred to as unfold PCA. However,
4 the unfolding may obscure the natural relationship between variables of different dimensions.
5 The use of multi-way chemometric methods (e.g. PARAFAC1 and PARAFAC2) is often
6 more appropriate and offers several advantages compared to PCA.

7
8 PARAFAC1 is similar to PCA. Instead of a bilinear model handling matrices, it is a trilinear
9 model handling three-way arrays (or even higher-way arrays). This is illustrated in Figure 5.
10 As PCA, a PARAFAC1 model provides scores for the sample mode and loadings for the time
11 mode, but also an additional set of loadings for the sensor mode. Such a model has interesting
12 properties in itself [12,14] but as for PCA, the PARAFAC1 model assumes that the
13 underlying basic variations (latent variables) remain constant in shape across different
14 samples. If this holds and the underlying phenomena are additive in nature, then the
15 PARAFAC1 model directly and uniquely estimates the underlying phenomena [15]. This
16 strict (tri-)linearity of PARAFAC1 is very advantageous when applicable [12] but the model
17 will often be too restricted to deal with severely shifted time profiles e.g. as also seen in batch
18 production data [16].

19
20 PARAFAC2 is a development of the original PARAFAC1 model which aims at handling
21 shifted or more generally varying profiles in a more efficient manner than PARAFAC1 (and
22 PCA) can. In Figure 6 the graphical description of the PARAFAC2 model is shown. In this
23 figure, scores (A) and sensor loadings (C) are of a similar structure as in PARAFAC1.
24 However, in the time mode where it is assumed that there are shifts, there is not only one
25 loading matrix as there would be in a PARAFAC1 model. Rather, for each sensor, k , an
26 individual time loading matrix is given. Hence, the PARAFAC2 model allows that every
27 sensor (or alternatively sample) can have its own distinct set of time loadings. This is clearly
28 different from the assumptions e.g. in the PCA model where all samples are assumed to be
29 governed by the same set of time loadings.

30
31
32 The matrix formulation of the PARAFAC1 model [17] and PARAFAC2 model [18-20] are
33 shown below. The multi-way models are somewhat difficult to grasp from linear algebraic
34 expressions because linear algebra is intrinsically related to two-way matrices. However,

1 these expressions are important for specifying the exact nature of the PARAFAC2 model.
 2 The models are mostly fitted in a least squares sense.

3

4 PARAFAC1: $\mathbf{X}_k = \mathbf{A}\mathbf{D}_k\mathbf{B}^T + \mathbf{E}_k$

5 PARAFAC2: $\mathbf{X}_k = \mathbf{A}\mathbf{D}_k\mathbf{B}_k^T + \mathbf{E}_k$

6

7 where \mathbf{X}_k is the k th frontal slab of the three-way array and \mathbf{D}_k is a diagonal matrix holding the
 8 k th row of \mathbf{C} in its diagonal. \mathbf{A} , \mathbf{B}/\mathbf{B}_k , and \mathbf{C} are parameters to be estimated and \mathbf{E}_k residuals.
 9 The major difference between the two models is that PARAFAC2 allows the loading matrix
 10 for mode B to be different for the k sensors (\mathbf{B}_k).

11

12 The PARAFAC1-solution cannot be rotated without a loss of fit and hence only one best-fit
 13 solution (unique) is possible. For the PARAFAC2 model to be valid and to retain uniqueness,
 14 all the cross-products matrices of the \mathbf{B}_k matrices are constrained to be constant for all k . This
 15 can be formulated as [19]:

16

17 $\mathbf{B}_k^T\mathbf{B}_k = \mathbf{B}^T\mathbf{B}, \quad k = 1, \dots, K$ (3)

18

19 This means that for every sensor k , a set of profiles \mathbf{B}_k (e.g. elution profiles [18]) is estimated
 20 under the constraint that the cross-products of the profile matrices are identical.

21

22 *2.4. Post-processing/graphical presentation*

23 The appropriate selection of pre-processing steps combined with the correct chemometric
 24 model provides a solution that is ready to be further analysed using post-processing
 25 techniques. The most obvious post-processing step is the graphical illustration of the model
 26 parameters, which will be described in the following section.

27

28 Data from sensor based analytical instruments can be arranged as *Sample/Sensor* \times *Time* \times
 29 *Sensor/Sample*. The results of a model fitted to such data directly can be presented in terms of
 30 parameters (scores or loadings) from a certain mode. The sample and sensor modes are
 31 illustrated with score and loading plots (scatter plots), respectively, whereas the time profiles
 32 are best shown as sensor signal vs. time. In Table 2 the number of elements in the plots from
 33 the different chemometric models are compared.

34

1 It can be seen that the use of the multi-way chemometric models gives a simpler loading plot
2 of the sensor mode as compared to a corresponding PCA model of the data rearranged into a
3 matrix. The plots of mode B differ for the different multi-way models and the special features
4 of the models are illustrated in this mode. PARAFAC2 that allows a certain freedom in one
5 mode (e.g. mode B) contains more loading elements, but the simple structure of the scores A
6 and loadings C has been preserved.

7

8 Another way to evaluate the models parameters is to use different classification technique
9 (e.g. K Nearest Neighbours (KNN) or Soft Independent Modelling of Class Analogy
10 (SIMCA) [21]. In this paper both visual interpretation of scatter plots and results from a KNN
11 classifier will be used to validate the chemometric models.

12

13

1 **3. Experimental**

2 *3.1. Data origin*

3 In the food industry great emphasis has been put on the standardization and control of the
4 food production. By minimizing the variations similar products can be produced even under
5 varying conditions. Deviations will occasionally occur and it is then important to be able to
6 recognize this and to take the proper actions. In the spring 2002, Pingvin Lakrids A/S
7 initiated an investigation to find a fast, easy, and more or less automated instrumental
8 analytical technique that could help them differentiate between normal licorices and deviating
9 licorices. For this purpose an electronic nose showed promising results.

10

11 *3.2. Instrumentation*

12 Analyses were conducted using an α -FOX-3000 electronic nose manufactured by Alpha
13 M.O.S (Multi Organoleptic System) (Alpha M.O.S, Toulouse, France, [http://www.alpha-](http://www.alpha-mos.com/)
14 [mos.com/](http://www.alpha-mos.com/)) with a sensor array system consisting of twelve Metal Oxide Semiconductive
15 sensors (MOS) distributed in two chambers. Chamber one: SY/LG, SY/G, SY/AA, SY/Gh,
16 SY/gCTI, and SY/gCT; chamber two: T30/1, P10/1, P10/2, P40/1, T70/2, and PA2. The
17 signal of sensor one (SY/LG) was very noisy and perturbed and pre-analyses showed that this
18 sensor disturbed the interpretation of the results. The data from sensor one (SY/LG) was
19 therefore not included in the data set.

20

21 The presence of two additional sensors in each chamber permitted control of the internal
22 temperature and the relative humidity. The temperatures in chamber one and two were 95°C
23 and 65°C, respectively, and the humidity was adjusted to 0 % RH in both chambers.

24

25 *3.3. Samples*

26 A commercial licorice product manufactured by Pingvin Lakrids A/S, Avedøre, Denmark
27 was used for the analysis. From this three different licorice characteristics were evaluated: 1)
28 Good licorices (GOOD), 2) Bad licorices (deviating licorices) (BAD) and 3) Fabricated bad
29 licorices (FBAD). The latter samples were prepared by drying the good licorices for a longer
30 time and were included because of an observed burned taste in the bad licorices and this was
31 tried resembled by a further drying of the good licorices.

32

33 Equilibrium was established between the headspace and the licorice sample and from the
34 headspace 1 ml gaseous sample was transferred to the sensors. An acquisition time of 120 sec

1 measuring the sensor signal (S) every $\frac{1}{2}$ sec was applied for each measurement. Between the
2 measurements the baseline of the sensors were re-generated with carrier gas and the baseline
3 signal (S_0) measured. The experiment was completely randomised and each licorice
4 characteristic was evaluated with six replicates.

5
6 Initial data analysis using unfold PCA and visual inspection of raw sensor signals revealed
7 that sample FBAD-4 was an outlier and thus this sample is removed from the data set. Thus
8 one element in $\underline{\mathbf{X}}$ corresponds to the measured sensor signal from sample i , measured at time
9 j , with sensor k , where $i = 1, \dots, 17; j = 1, \dots, 241; k = 1, \dots, 11$.

10

11 *3.4. Data handling*

12 The pre-processing and modelling steps were performed using a Graphical User Interface
13 (GUI) program - SENSABLE. SENSABLE is a newly developed processing tool for sensor
14 based data that can be downloaded at <http://www.models.kvl.dk/source/>. In Appendix A, a
15 screen dump of the SENSABLE program is shown. The current version of SENSABLE
16 requires MATLAB at least version 6.x.

17

18 The different chemometric models can also be calculated using N-way toolbox version 2.10
19 and the PARAFAC2 toolbox for MATLAB[®] version 6.5 (The MathWorks Inc., Natick, MA,
20 USA). The toolboxes can be downloaded at <http://www.models.kvl.dk/source/>.

21

22 The K Nearest Neighbours (KNN) classification is performed calculating the distance to the
23 three nearest neighbours using Euclidian distance in scores. Given the low number of samples
24 in each group, it was not tested to increase the number of neighbours. The scores of three test
25 samples – one from each licorice group randomly selected – are predicted (not included in the
26 calibration) and the distances to the calibration samples calculated yielding a classification of
27 each sample. This exclusion and classification of three samples was repeated several times.

4. Results

The effect of different pre-processing steps is evaluated and validated using an appropriate chemometric model. The combinations of pre-processing steps and the model will be assessed according to explained variance of the model and according to how well the samples are separated in a score plot (visually evaluated). Note though, that percentage fitted variation in itself does not provide a useful measure of appropriateness. For mathematical reasons, a model with more degrees of freedom (e.g. unfold PCA compared to PARAFAC1) will always fit better. The problem is that the better fit will mostly be fit to noise if the mathematically simpler model (PARAFAC1) is already fitting the systematic variation. The degree of overfit can be evaluated by comparing the fit with the cross-validated fit. The simplest possible yet adequate model is desired. Beforehand, PARAFAC1 can be anticipated to be adequate and hence sufficient, if the underlying time profiles do not change shape only magnitude across different sensors or samples. Unfold PCA on the other hand will be appropriate if the time profiles are unrelated across different sensors or samples. PARAFAC2 is an intermediate model that allows the time profiles to vary in shape but not to be absolutely unrelated.

The focus will be on separating the normal licorices (good) from the deviating licorices (bad), and not as much to separate the good or deviating licorices from the fabricated bad licorices as the latter two are expected to be more similar and of less interest to separate. Care must be taken, as e.g. removing outliers from the data will affect the pre-processing of the data. Upon removal of outliers the data must thus be re-preprocessed and a new model calculated before the actual model evaluation.

Initially the data are assumed to be well analysed with a multi-way chemometric model and the PARAFAC2 model is a likely candidate due to observed shift between the different sensors signals. In the PARAFAC2 model orthogonality constraints are applied for the sensor and the sample mode which helps avoiding numerical problems. All models presented in this article are partly validated using cross-validation by splitting the sample mode in seven randomly chosen intervals.

The number of components selected for each model is determined from a thorough model inspection using the following parameters: explained variance, size and structure of residuals compared to modelled data and raw data, split-half analysis, and visual interpretation of the

1 scores and loadings. While simple quantitative diagnostics such as cross-validation results
2 seem easier to use, it is important to realize that these only rarely provide correct answers
3 when used for exploratory models as those in this paper. It is always necessary to use
4 additional means such as the above for overall validation. In the following, the basic
5 workflow for achieving a good model is presented.

6 7 *4.1. Pre-processing steps*

8 4.1.1. Baseline correction and centering/scaling

9 In Table 3 PARAFAC2 models for different baseline corrections are evaluated. As shown in
10 Figure 2 the baseline subtracted (fractional pre-scaled) sensor signals are of widely different
11 scales. This indicates that scaling the sensor mode is useful to give all the sensors the same
12 opportunity to influence the model.

13
14 As seen from Table 3 several combinations of centering/scaling and baseline pre-scaling
15 methods provide an efficient separation of the samples. Thus one could argue that selecting
16 just one of these combinations will give a sufficient pre-processing. However, the pre-
17 processing steps interact and thus must be interchanged before the most efficient combination
18 can be selected. Depending on the data size and structure this process can be time-consuming
19 and thus a standard pre-processing is often selected. With good insight and experience such
20 prior determined pre-processing will often be useful, but problems may arise when new types
21 of samples or sensors are included in the analysis.

22
23 The use of integrated software facilitates comparison of many different models. This is
24 helpful when exploring the usefulness of new models as those presented here. In the
25 following a good combination of pre-processing, transformation and model is selected to
26 show the effect of transformation and feature extraction. All combinations of pre-processing
27 steps have been evaluated to make sure that the optimal selection from Table 3 is still optimal
28 after the evaluation of the last two pre-processing steps.

29
30 An example of the separation using the fractional baseline correction and centering and
31 scaling is shown in Figure 7A. Figure 7A shows that the good samples are well separated
32 from the deviating samples, but also that it is possible to separate the deviating samples from
33 the fabricated bad licorices. As expected the latter two licorice characteristics appear more
34 similar to each other than to the good licorices.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

Transforming the data with logarithmic and similar transformations are not likely relevant for these data (initial results not shown confirm this).

4.1.2. Feature extraction

The concept of feature extraction is to extract the relevant information from the data and to reduce the number of parameters in the data set, ultimately leading to more robust and simple models. In Table 4 an evaluation of some feature extraction methods is shown.

From Table 4 it can be seen that using a subset of the time variables gives a model that is explaining the data well and separates the samples as efficient as the model using all time variables (score plot is not illustrated, but is very similar to Figure 7A). The use of only the maximum of each sensor response curve gives a separation of the samples that is quite similar to the three-way models (Figure 7B). Other types of feature extraction makes the data array smaller and simpler, but offer a less sufficient separation of the samples.

To further evaluate the separation of the three groups a KNN classification is performed and the results are shown in Table 5.

Table 5 indicates that despite an almost perfect classification of the test samples some information about the samples is lost using only the maximum. But the use of the maximum makes it possible to use ordinary PCA and gives a good indication that the samples can be separated; however, the results using PARAFAC2 on the full data set are better though. This concept of using only one variable from the time mode (i.e. reducing the dimension of the data array) is often seen in analyses using commercial electronic noses because the integrated software limits the data to two dimensions.

4.2. Processing step (modelling of pre-processed data)

As was seen from Figure 7 both two-way and three-way chemometric models are capable of separating the samples into the three groups.

The major difference between the PARAFAC1 and PARAFAC2 is that PARAFAC2 is modelling each sensor separately (B_k loadings), whereas PARAFAC1 estimate a common time profile loading (B loading) that is used for all sensors, as illustrated in Figure 8.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

The scale showing the size of the B_k loadings is arbitrary and can not be directly compared to the real time profiles. Both component one and two of the PARAFAC2 model contributes to the B_k loadings, but it is seen that the contribution from component one is rather noisy and much smaller in magnitude than the loadings of the time profiles given in component two.

PARAFAC2 includes eleven different time profile loadings and this gives a superior separation of the samples as seen from Figure 7A. However, for each component some of the time profiles seem to differ only on the magnitude level and thus a model like PARAFAC1 might be adequate to describe the fewer common time profile shapes.

From pre-analyses (e.g. split-half and core consistency analyses) it was found that the more restricted PARAFAC1 model also needs two components to give the optimal (unique) solution. However, as shown in Figure 9 the samples are not separated in the same efficient way as for PARAFAC2 (and PCA).

The time profile loadings for PARAFAC1 in Figure 8 also shows that the two components have shapes that are very similar to the real sensor response curves (Figure 2). This suggests that the information of the shifted time profiles is not included in the model, but that only an average sensor time profile is allowed for each component. This could indicate that more components are needed to describe the shifted behaviour of the sensor signals. However, the model gets worse including more than two components - thus PARAFAC1 seems to be a too restricted model for this data.

The insufficient separation of the samples using PARAFAC1 indicates that this kind of data is not strictly tri-linear, but rather that the time profile of each sensor should be modelled independently, as is the case using PARAFAC2. The result from the PCA model also suggests that most information of the different samples is located in the peak, and in particular in the maximum of the signal. However, ignoring the three-way structure of the original data tends to give a less sufficient separation of the samples and suggests that important information also can be found in the rest of the time profile curve.

1 **5. Conclusion**

2 In this article the processing of sensor based data has been outlined and PARAFAC2 has
3 been suggested for enabling analysis of the raw data even when there are time shifts.
4 Refraining from extracting features is expected to make it possible to extract more
5 information in situations where more subtle details of the data are needed for proper
6 quantification or classification. The extent, to which this is feasible, in general, remains to be
7 explored in various contexts.

8
9 Choosing appropriate pre-processing and model requires knowledge of the data structure and
10 the model. For sensor based data the PARAFAC2 model that handles shifted profiles offers
11 some advantages compared to related methods such as PCA. However, the structure of the
12 PARAFAC2 model requires a more careful interpretation of the model parameters and an
13 extensive post-processing procedure could be essential to the concluding interpretation.

14
15 The electronic nose combined with the correct data pre-processing, processing and post-
16 processing was found to have the potential to be a fast, easy to use and suitable analytical
17 methods in the application, but some important hurdles must be passed before this method
18 can be introduced as an at-line method. Most notably, more experience is needed in using the
19 multi-way approach for analyzing sensor data.

20 21 22 **Acknowledgements**

23 The authors are grateful to the Danish Technical University (DTU), Lyngby, Denmark for
24 placing the electronic nose at our disposal and for providing expertise on this field. We also
25 wish to thank Pingvin Lakrids, Toms Group A/S, Avedøre, Denmark for letting us use their
26 facilities to produce the samples.

27

1 **References**

- 2 [1] P.N. Bartlett, J.W. Gardner, *Sensors and Sensory Systems for an Electronic Nose*. Kluwer
3 Academic Publishers, London, UK, 1992.
- 4
- 5 [2] C. Gretsch, J. Delarue, A. Toury, P. Visani, R. Liardon, Detection of aroma above a coffee
6 powder: limits and perspectives of electronic sensors. ASIC Conference - 17th International
7 Scientific Colloquium on Coffee, Nairobi, Kenya, 1997.
- 8
- 9 [3] Y.G. Martín, J.L.P. Pavon, B.M. Cordero, C.G. Pinto, Classification of vegetable oils by linear
10 discriminant analysis of Electronic Nose data, *Analytica Chimica Acta* 384 (1999) 83-94.
- 11
- 12 [4] R.N. Bleibaum, H. Stone, T. Tan, S. Labreche, E. Saint-Martin, S. Isz, Comparison of sensory
13 and consumer results with electronic nose and tongue sensors for apple juices, *Food Quality and*
14 *Preference* 13 (2002) 409-422.
- 15
- 16 [5] H.-Y. Chung, J.A. Partridge, B.R. Harte, Discrimination of light oxidized off-flavors in milk and
17 cheddar cheese using solid-phase microextraction - gas chromatography and the olfactory
18 sensing technique. *Proceedings of the 13th IAPRI Conference on Packaging* (2002) 1144-1155.
- 19
- 20 [6] J.W. Gardner, M. Craven, C. Dow, E.L. Hines, The prediction of bacteria type and culture
21 growth phase by an electronic nose with a multi-layer perception network, *Measurement Science*
22 *& Technology* 9 (1998) 120-127.
- 23
- 24 [7] S. Roussel, G. Forsberg, V. Steinmetz, P. Grenier, V. Bellon-Maurel, Optimisation of electronic
25 nose measurements. Part I: Methodology of Output Feature Selection, *Journal of Food*
26 *Engineering* 37 (1998) 207-222.
- 27
- 28 [8] T. Eklöv, P. Mårtensson, I. Lundström, Selection of Variables for Interpreting Multivariate Gas
29 Sensor Data. *Analytica Chimica Acta* 381 (1999) 221-232.
- 30
- 31 [9] C. Distante, M. Leo, P. Siciliano, K.C. Persaud, On the study of feature extraction methods for an
32 electronic nose, *Sensors and Actuators B* 87(2) (2002) 274-288.
- 33
- 34 [10] R. Gutierrez-Osuna, H.T. Nagle, B. Kermani, S.S. Schiffman, *Handbook of Machine Olfaction*,
35 WILEY-VCH Verlag GmbH & Co., Weinheim, Germany, 2003.
- 36

- 1 [11] R. Bro, A.K. Smilde, Centering and scaling in component analysis, *Journal of Chemometrics* 17
2 (2003) 16-33.
3
- 4 [12] R. Bro, PARAFAC – Tutorial and Applications, *Chemometrics and Intelligent Laboratory*
5 *Systems* 38 (1997) 149-171.
6
- 7 [13] H.A.L. Kiers, Towards a standardized notation and terminology in multiway analysis, *Journal of*
8 *Chemometrics* 14 (2000) 105-122.
9
- 10 [14] R.A. Harshman, Foundations of the PARAFAC procedure: Model and conditions for an
11 ‘explanatory’ multi-mode factor analysis, *UCLA Working Papers in Phonetics* 16 (1970) 1-84.
12
- 13 [15] N.D. Sidiropoulos, R. Bro, On the uniqueness of multilinear decomposition of N-way arrays,
14 *Journal of Chemometrics* 14 (2000) 299-240.
15
- 16 [16] B.M. Wise, N.B. Gallagher, E.B. Martin, Application of PARAFAC2 to fault detection and
17 diagnosis in semiconductor etch, *Journal of Chemometrics* 15 (2001) 285-298.
18
- 19 [17] R.A. Harshman, M.E. Lundy, PARAFAC: Parallel Factor Analysis, *Computational Statistics &*
20 *Data Analysis* 18 (1994) 39-72.
21
- 22 [18] R. Bro, H.A.L. Kiers, C.A. Andersson, PARAFAC2 - Part II. Modelling chromatographic data
23 with retention time shifts, *Journal of Chemometrics* 13 (1999) 295-309.
24
- 25 [19] H.A.L. Kiers, J.M.F ten Berge, R. Bro, PARAFAC2 - Part I. A direct fitting algorithm for the
26 PARAFAC2 model, *Journal of Chemometrics* 13 (1999) 275-294.
27
- 28 [20] R.A. Harshman, PARAFAC2: Mathematical and technical notes. *UCLA Working Papers in*
29 *Phonetics* 22 (1972) 30-44.
30
- 31 [21] S. Wold, C. Albano, W. J. Dunn, III, U. Edlund, K. H. Esbensen, P. Geladi, S. Hellberg, E.
32 Johansson, W. Lindberg, and M. Sjöström, Multivariate data analysis in chemistry, in: B. R.
33 Kowalski (Ed.), *Chemometrics - Mathematics and Statistics in Chemistry*, D. Reidel Publishing
34 Company, Dordrecht, Netherlands, 1984.
35
36
37

1 **Figure legends**

2 Figure 1. Example of three-way data set, X from an electronic nose.

3

4 Figure 2. Left: Raw/Real sensor signal and Right: Baseline corrected sensor signal - (S-
5 S0)/S0: of one sample and eleven MOS (Metal Oxide Semiconductive) sensors.

6

7 Figure 3. Illustration of the principles of data that can be modelled with PCA (A) and of data
8 that show shifted profiles and thus can not be modelled using PCA (B). X is the two-way
9 data; T is the score vector and P the loading vector.

10

11 Figure 4. Three-way nose data of dimension I (samples) \times J (times) \times K (sensors). The data
12 can be matricized to two-way data by concatenating every frontal slab next to each other.
13 Note that sensors in the second (J) mode and time in the third (K) mode is then confounded in
14 the resulting matrix. The data can also alternatively be unfolded in the other modes.

15

16 Figure 5. Graphical explanation of the principles of the PCA model (A) and PARAFAC1
17 model (B) using a two-component model as an example.

18

19 Figure 6. Graphical description of the PARAFAC2 model. Data: Sensors \times Time \times Samples.
20 Mode A: sensors, mode B: time profiles, and mode C: samples.

21

22 Figure 7. Score plots for PARAFAC2/PCA models. (A): PARAFAC2 model of the fractional
23 baseline corrected data centering the sample mode and scaling the sensor mode (Data size:
24 $17 \times 241 \times 11$). (B) PCA model of data where the maximum of each sensor signal is extracted
25 (Data size: 17×11).

26

27 Figure 8. B_k loadings (time profile loadings) for each sensor of the PARAFAC2 model
28 component 1 (A) and component 2 (B). B loadings for a PARAFAC1 model with two
29 components (C). From pre-analyses the number of components is found to be the most
30 optimal describing the pre-processed data for the given model.

31

32 Figure 9. Score plot for samples for a PARAFAC1 model with 2 components.

33

1 **Table legends**

2 Table 1. Baseline correction formulas for MOS sensors (modified from [9])

3

4 Table 2. Characteristics of plots from the different chemometric models (in brackets are
5 shown the number of elements of an $I \times J \times K$ array in the different plots). Unfold PCA is
6 PCA performed on the three-way array rearranged (unfolded) to a matrix.

7

8 Table 3. Baseline correction methods and use of centering and/or scaling evaluated from
9 explained variance and score plots from PARAFAC2 models. Two components are included
10 in all models. Note that explained variance is in percentage of pre-processed data and hence
11 not directly comparable for different models.

12

13 Table 4. Feature extractions based on fractional pre-scaled data evaluated from PARAFAC2
14 model ($N > 2$) and PCA models ($N = 2$). Two components are used in each model and the
15 sample mode is centered and sensor mode scaled.

16

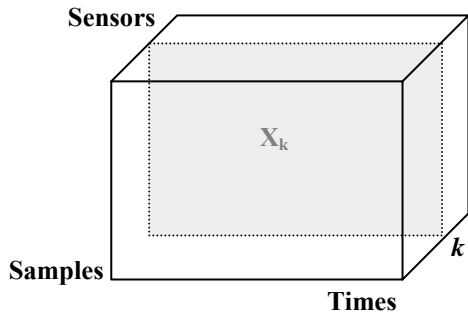
17 Table 5. Results from KNN classification calculating the distance from each test sample to
18 the three nearest neighbours. The test set consists of one sample randomly selected from each
19 licorice group. Sixty samples were evaluated for each model.

20

1 **Figures**

2 Figure 1

3



4

5

6 Figure 2

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

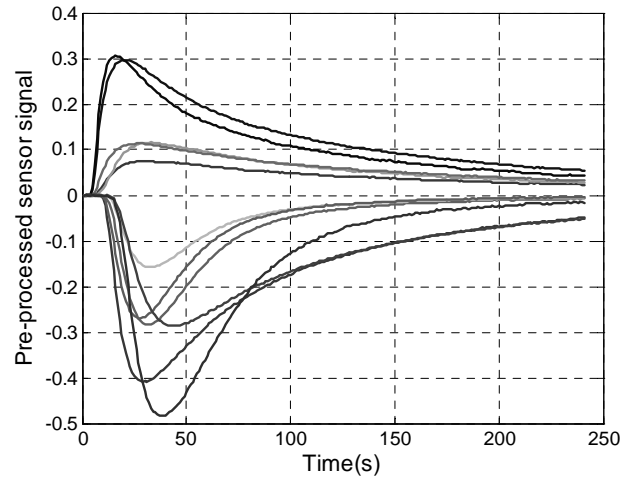
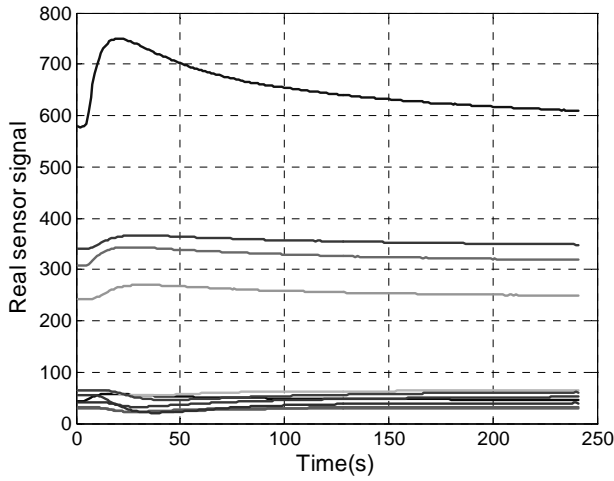


Figure 3

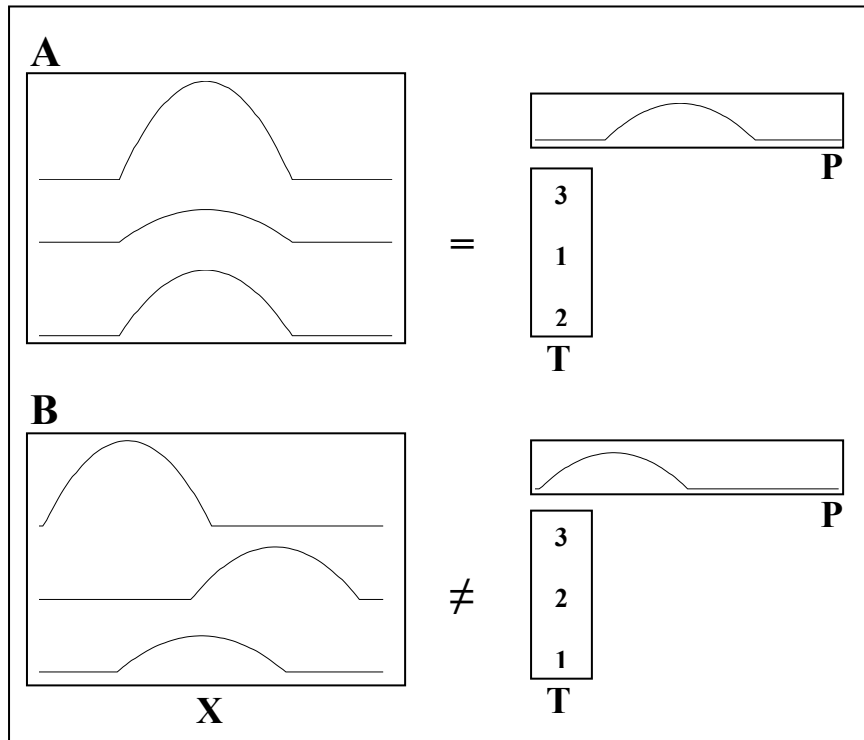
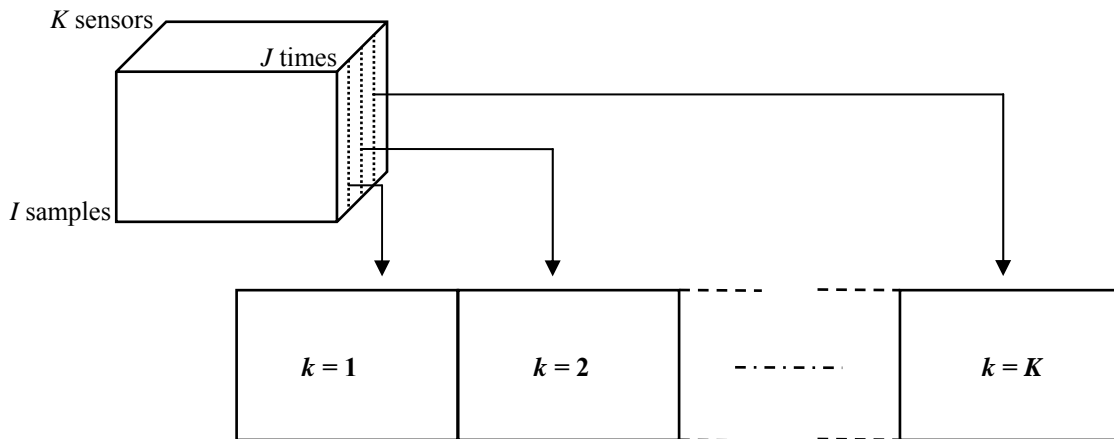
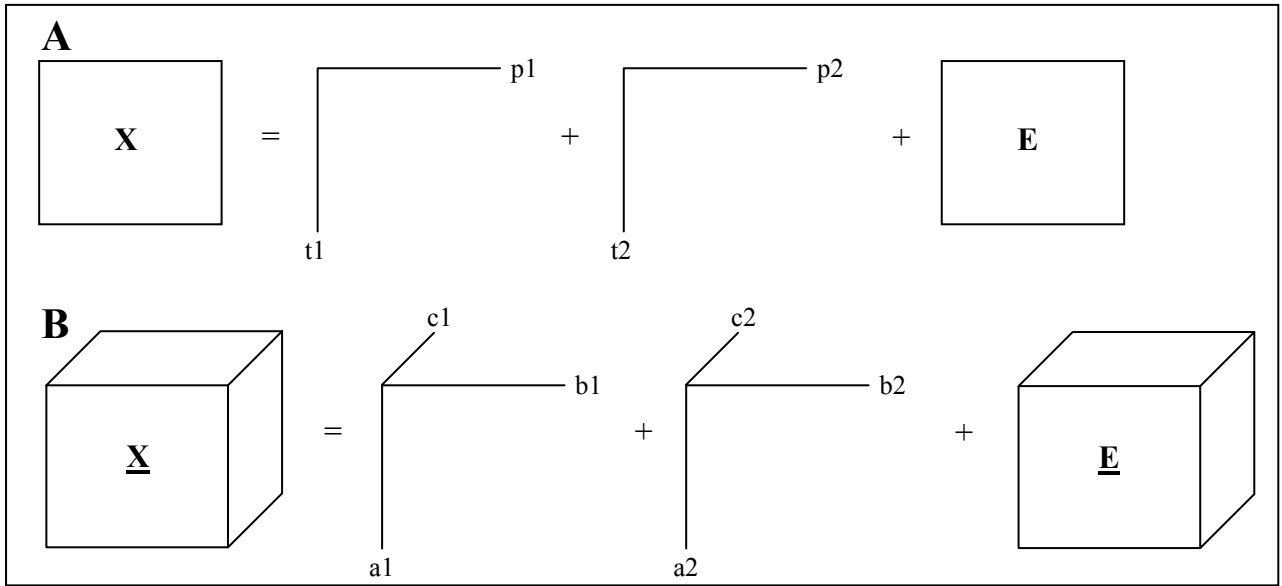


Figure 4

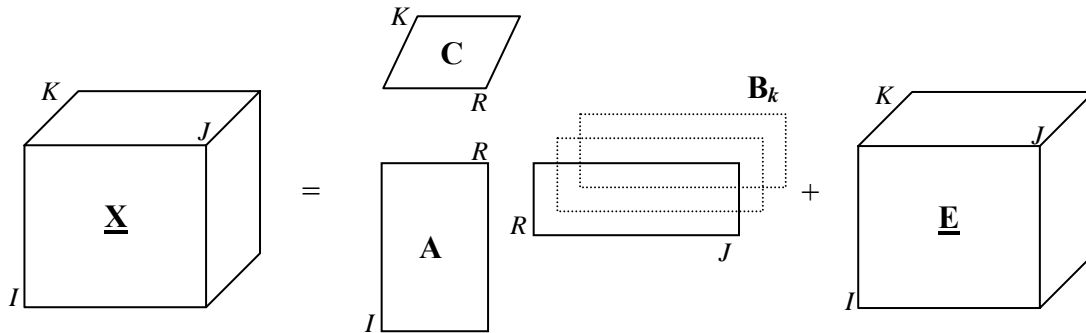


25
26
27
28
29
30
31
32
33
34
35

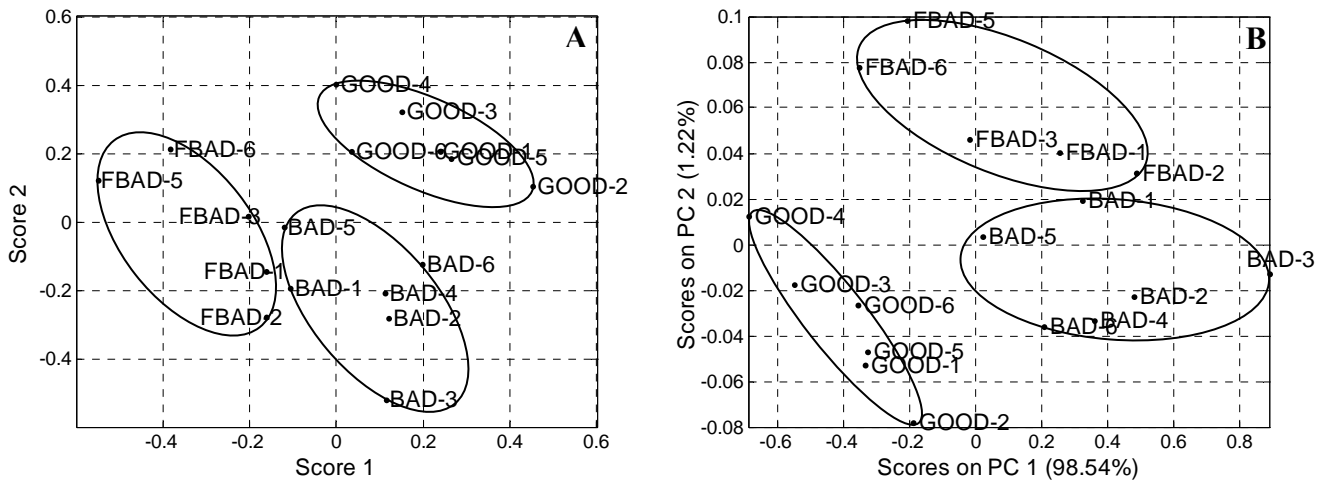
1 Figure 5
2



3
4
5 Figure 6



6
7
8 Figure 7
9



10
11
12
13
14
15
16
17
18
19
20
21
22

Figure 8

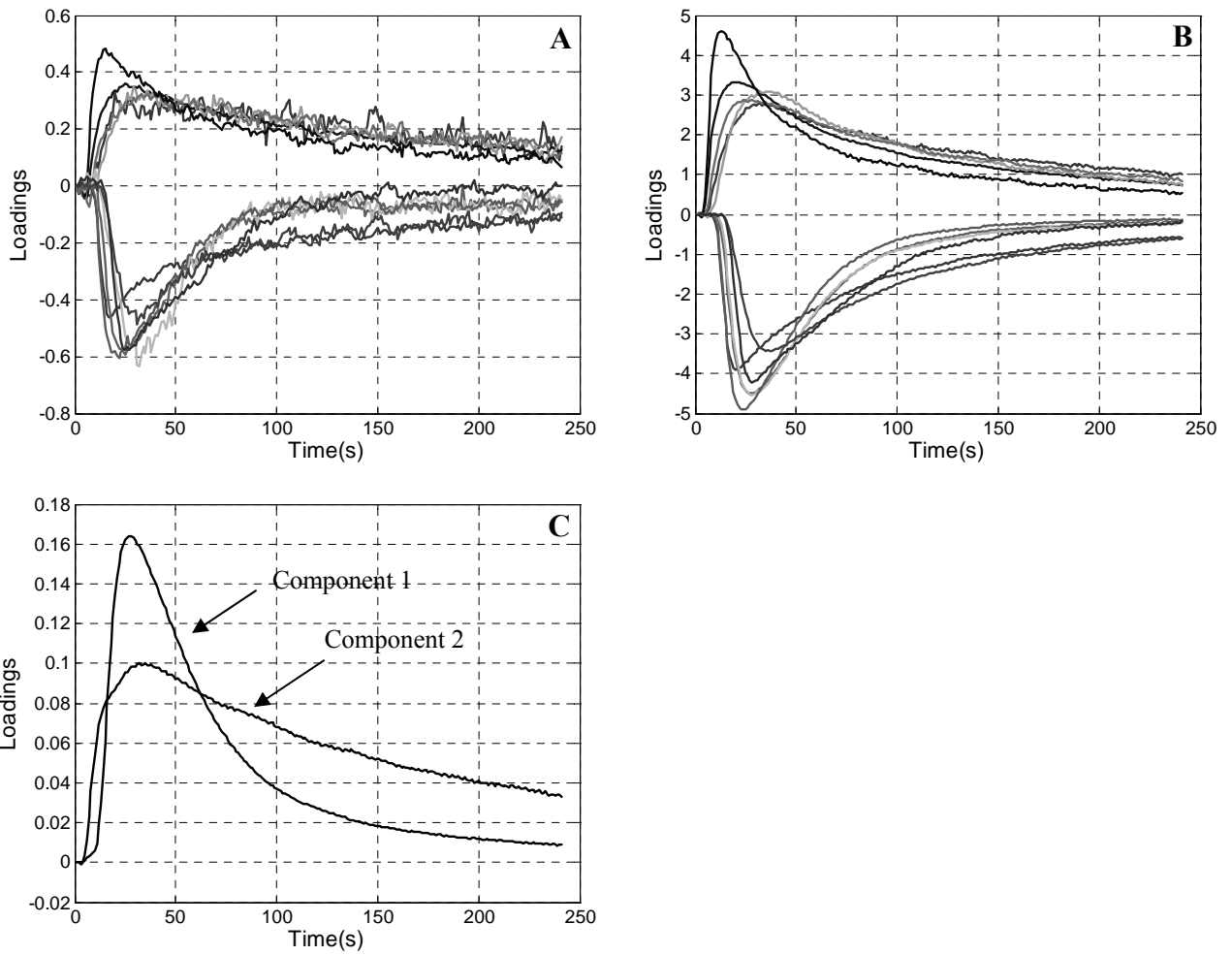
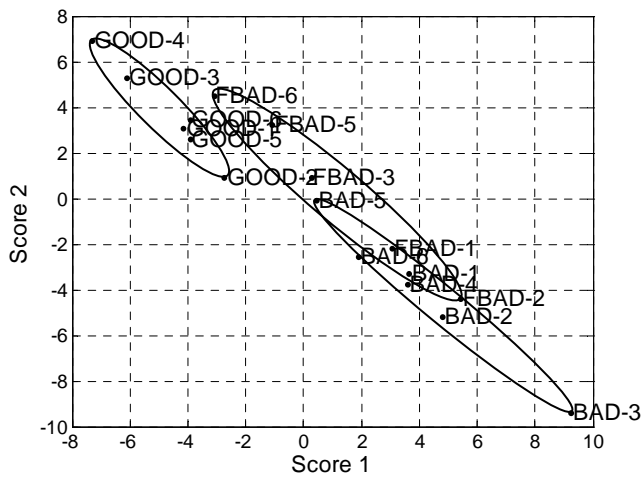


Figure 9



Tables

Table 1

Method	Formula ⁽¹⁾
Difference	$S(t) - S_0$
Relative	$S(t) / S_0$
Fractional	$(S(t) - S_0) / S_0$

⁽¹⁾ S_0 is the baseline signal (signal of the carrier gas) and $S(t)$ the signal after a certain time.

Table 2

Model	Mode A	Mode B (time) ⁽¹⁾	Mode C
Unfold PCA	Score plot (I)	Loading plot ($J \times K$) ⁽²⁾	
PCA ⁽³⁾	Score plot (I)	Loading plot (K)	
PARAFAC1	Score plot (I)	Sensor signal vs. Time ($J \times R$)	Loading plot (K)
PARAFAC2	Score plot (I)	Sensor signal vs. Time ($J \times K \times R$)	Loading plot (K)

⁽¹⁾ R is the number of components in the model

⁽²⁾ Mode B is nested within mode C giving a complex loading plot with many loading elements

⁽³⁾ Extracting one feature of mode B reduces the data size to $I \times K$

Table 3

Method	Centering and/or scaling				Model characteristics		
	None	Centering	Scaling	Centering + scaling	Explained variance	Separation of samples ⁽²⁾	
Difference ⁽¹⁾ $S-S_0$			X		99.96	++	B-FB
				X	99.06	++	B-FB
Fractional $(S-S_0)/S_0$	X				99.97	-	
		X			99.30	++	B-FB
			X		99.96	-	
				X	99.20	++	
Relative S/S_0	X				99.99	+(+)	B-FB
		X			99.30	+(+)	B-FB
			X		99.99	+(+)	B-FB
				X	99.30	+(+)	B-FB

⁽¹⁾ Scaling is essential for this baseline correction method

⁽²⁾ Separation of normal and deviating licorices evaluated from score plots with the first component plotted against the second. -: no separation, +: some separation, and ++: samples are well separated. Symbols in brackets indicate a less significant separation than the same symbol with no brackets. Letters - G: good, B: bad, FB: Fabricated bad – are used to indicate which groups that are not ‘perfectly’ separated.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

Table 4

Feature extraction	Size of array	Model	Explained variance	Separation of samples ⁽⁴⁾	
All time variables ⁽¹⁾	17 × 241 × 11	PARAFAC2	99.20	++	
Subset of time variables ⁽²⁾	d = 3	17 × 81 × 11	PARAFAC2	99.21	++
	d = 7	17 × 35 × 11	PARAFAC2	99.21	++
	d = 20	17 × 13 × 11	PARAFAC2	99.21	++
	d = 40	17 × 7 × 11	PARAFAC2	99.15	++
Absolute maximum	17 × 11	PCA	99.76	++	
Max adsorption slope	17 × 11	PCA ⁽³⁾	92.31	-	
Max desorption slope	17 × 11	PCA ⁽³⁾	89.11	-	

⁽¹⁾ PARAFAC2 model from Table 3 that is fractional pre-scaled, sample mode centered and sensor mode scaled

⁽²⁾ Subset using every d variable in the time mode (data array \underline{X} : 17 × 1:d:241 × 11). If the size of the time mode is not divisible with d, the number of features selected will be rounded to the lowest integer.

⁽³⁾ Five components in the PCA model

⁽⁴⁾ See Table 3.

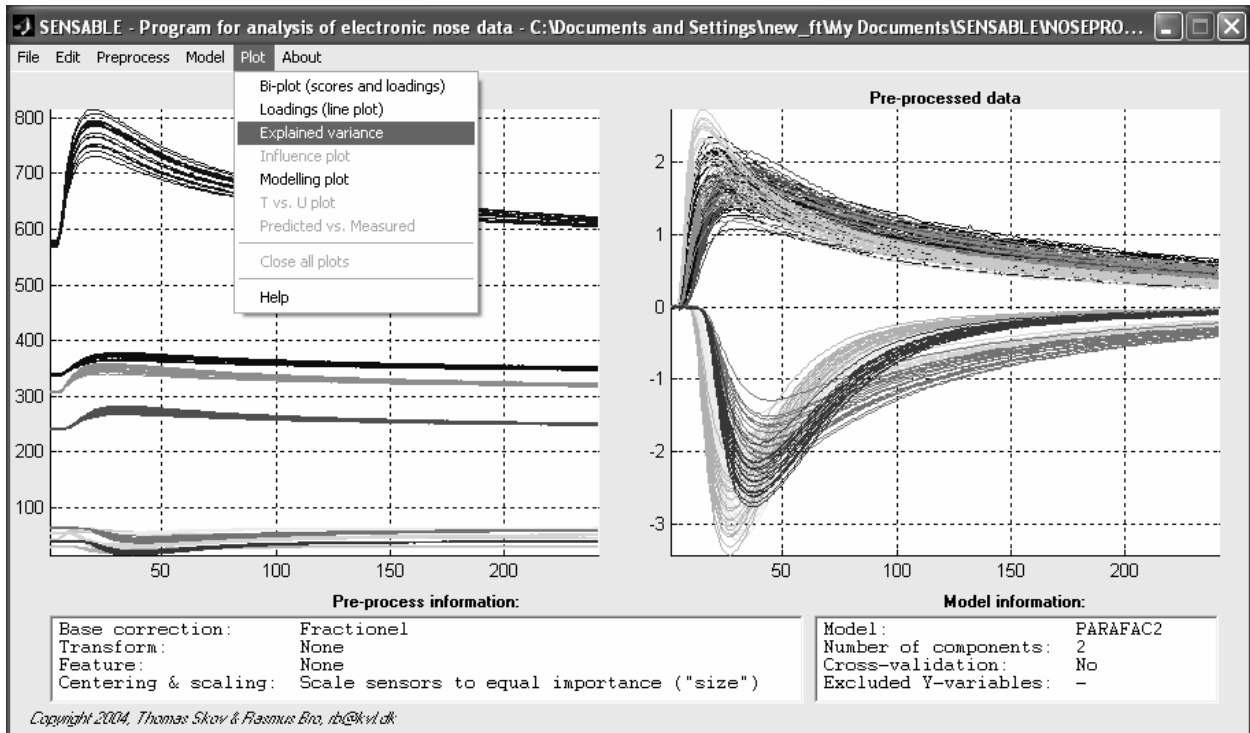
Table 5

Model	Classification	
	Correct	Incorrect
PCA	57	3
PARAFAC2	59	1

1 **Appendix**

2 Screen dump from the SENSABLE program

3



4

5