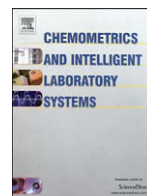




Contents lists available at ScienceDirect

Chemometrics and Intelligent Laboratory Systems

journal homepage: www.elsevier.com/locate/chemolab

How to build a robust model against perturbation factors with only a few reference values: A chemometric challenge at 'Chimimétrie 2007'

Juan Antonio Fernández Pierna^a, Fabian Chauchard^b, Sébastien Preys^b, Jean Michel Roger^c, Oswin Galtier^d, Vincent Baeten^a, Pierre Dardenne^{a,*}

^a Walloon Agricultural Research Centre (CRA-W), Valorisation of Agricultural Products Department, Henseval Building, Chaussée de Namur no.24, 5030 Gembloux, Belgium

^b Ondalys, 385, avenue des Baronnes, 34730 Prades-Le-Lez, France

^c Research Group Cemagref - SupAgro, BP 5095, 34196 Montpellier, France

^d University Paul Cézanne, 3 avenue Robert Schuman, 13628 Aix-en-Provence, France

ARTICLE INFO

Article history:

Received 5 November 2009

Received in revised form 20 May 2010

Accepted 21 May 2010

Available online xxxx

Keywords:

Robust

Chemometrics

Reference values

Challenge

ABSTRACT

Following up on the success of previous chemometric challenges arranged during the annual congress organised by the French Chemometrics Society, the organisation committee decided to repeat the idea for the Chimimétrie 2007 event (<http://www.chimimetrie.org/>) held in Lyon, France (29–30 November) by featuring another dataset on its website. As for the first contest in 2004, this dataset was selected to test the ability of participants to apply regression methods to NIR data. The aim of Challenge 2007 was to perform a calibration model as robust and precise as possible using a data set with only a few reference values available and submitted to different perturbation factors. The committee received nine answers; this paper summarizes the best three approaches, as well as the approach proposed by the organisers.

© 2010 Published by Elsevier B.V.

1. Introduction

For the fourth consecutive year and following on from the success of the chemometric contests organised during previous congresses [1–3], another dataset was proposed for the 'Chimimétrie 2007' meeting (<http://www.chimimetrie.org/>) held in Lyon, France (29–30 November 2007). As for the first contest in 2004, this dataset was selected to test the ability of participants to apply regression methods to NIR data. The aim of Challenge 2007 was to perform a calibration model as robust and precise as possible using a data set with only a few reference values available and submitted to different perturbation factors. In the context of robust statistics, the term robust is reserved for a group of methods that provide good estimates for the majority of data. [4] Robust approaches, which are able to neglect the outlier's presence and to represent the data majority, are strongly required in order to draw reliable conclusions about data [5–10]. But robustness is also of crucial importance in control system designs. Real systems are vulnerable to different perturbation factors and measurement noise and there are always discrepancies between mathematical models used for design and the actual system in practice [11–13]. In the early 1990 s, Dardenne et al. [14] studied the effects of various properties such as moisture, particle size and ambient temperature on NIR calibration models in order to reduce the wet chemistry needed for them. Artificial spectral

variations were created by changing the moisture and particle size of wheat samples when predicting the protein content. These samples, measured at different temperatures on a monochromator, produced wide spectral variations that helped to develop robust models. The data used in this challenge come from that paper.

Nine participants took up the challenge with the proposed data. The results were evaluated on the basis of the best validation criteria (R^2 and RMSEP) obtained for the predicted values of the test set, and also on the quality of the approach from a methodological perspective. The three best approaches were presented during the congress and are summarised here, together with the approach put forward by the organizers of Challenge 2007.

2. Materials and methods

Several datasets were provided to the participants: a calibration dataset, an experimental design dataset, a standard replicates dataset and a test dataset.

2.1. Calibration dataset

For this challenge, only 10 spectra of ground wheat acquired using a FOSS NIRSystems 4500, measured between 1300 nm and 2398 nm each 2 nm, were supplied, together with the protein content of the 10 samples measured in g/kg Dry Matter (DM). The aim was to build a calibration model as robust and precise as possible with the available 10 reference values for protein content.

* Corresponding author.

E-mail address: dardenne@cra.wallonie.be (P. Dardenne).

2.2. Experimental design dataset

There were also the spectra from an experimental design based on other 11 samples of ground wheat (with unknown reference protein values). These 11 whole grain samples were separated into two homogeneous sets: one was dried to reduce the moisture content to $\pm 9\%$, and the other was moistened to reach 13–14% humidity.

Each grain sample set was then divided again in two groups: the first subsample was ground finely (Cyclotec apparatus, level C) and the second one was ground more coarsely (Ika apparatus, level I). A total of 11×4 sample sets were thus available. These samples were measured by reflection NIR in duplicate at three room temperatures (18 °C, 23 °C and 27 °C). The experimental design database therefore contained $11 \times 4 \times 3 \times 2 = 264$ spectra.

2.3. Standard replicates dataset

Additional spectral information was supplied from one set of 10 samples (sealed cells) scanned on 31 different instruments of the same type and from a second set of 10 scanned on 17 instruments. A total of 480 spectra were thus available.

The reference protein values corresponding to the experimental design samples and to the 10 samples measured on the different instruments were unknown.

2.4. Test dataset

Some 2000 spectra of ground wheat from routine analyses were acquired from ± 10 different instruments with several levels of granulometry, humidity and temperature, and provided by the Requasud network (<http://www.requasud.be/>) from 1991 to 2007. The 2000 spectra each had a reference protein value obtained by the reference method, but these values were not communicated to the participants.

2.5. Notations

The calibration dataset was kept in matrices \mathbf{X}_{cal} , \mathbf{y}_{cal} whereas the data from the experimental design in a matrix \mathbf{X}_{ed} and the standard replicates dataset in two matrices \mathbf{X}_{instr1} (31 instruments) and \mathbf{X}_{instr2} (17 instruments). The Test dataset was kept in \mathbf{X}_{test} .

DM dry matter

EPO external parameter orthogonalisation

LS-SVM least-squares support vector machines

MSC multiplicative scatter correction

NIR near infrared

OSC orthogonal signal correction

PCA principal component analysis

PLS partial least squares

RMSEC root mean square error in calibration

RMSEP root mean square error in prediction

RPD ratio of standard error of prediction to sample standard deviation

SEC standard error of calibration

SECV standard error of cross-validation

SNV standard normal variate

3. Results

3.1. Participant 1

Two approaches were tested: exhaustive calibration on artificial data, and Orthogonal Signal Correction (OSC) [15] on these artificial data.

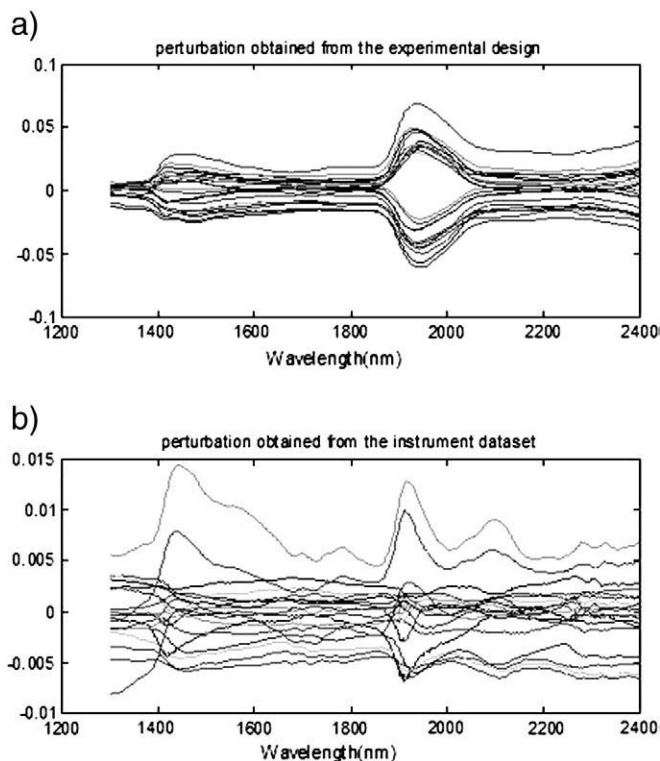


Fig. 1. Examples of perturbation spectra from a) the experimental design dataset and b) the standard replicates dataset.

3.1.1. Approach 1: exhaustive calibration on artificial data

This approach involved using information about the perturbation factors and 'injecting' it into the 10 spectra of the calibration dataset (\mathbf{X}_{cal}) in order to generate an artificial calibration database \mathbf{X}_{artif} . Then, PLS (Partial Least Squares) was applied in this new calibration database that included all the possible perturbation factors.

In order to do so, initially all the possible perturbing vectors were identified using the experimental design dataset (\mathbf{X}_{ed}) and the standard replicates dataset ($\mathbf{X}_{instr1} + \mathbf{X}_{instr2}$) for each of the 24 spectra ($\mathbf{X}_{sample i}$) of the same sample i (one sample i had 24 ' j ' combinations of possible perturbations). The 24 vectors of perturbations were calculated as follows:

$$-\delta(x_{i,j}) = x_{i,j} - \bar{x}_i \quad (1)$$

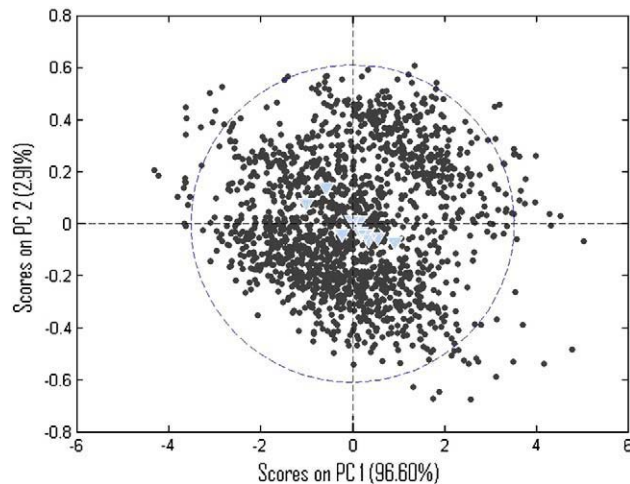


Fig. 2. PCA plot showing the spectral space of the increased calibration dataset (\mathbf{X}_{artif}) including the original \mathbf{X}_{cal} .

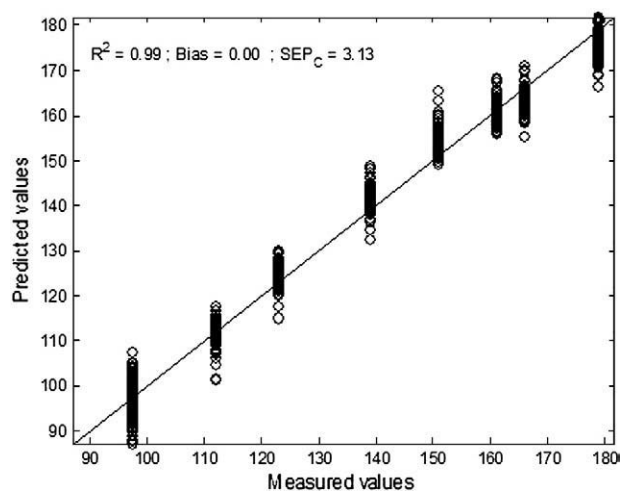


Fig. 3. Final model constructed using X_{artif} for participant 1.

where \bar{x}_i is the mean spectra of sample i and x_{ij} is the j th spectra of sample i . All the perturbations were put in a new matrix D_{ed} containing the 264 possible perturbations. The same procedure was followed in order to identify all possible perturbations from the standard replicates dataset, resulting in a matrix D_{instr} of 480 perturbation spectra. Fig. 1 shows an example of the perturbation spectra that could be obtained.

Finally, all the perturbations from D_{ed} and D_{instr} were combined (using addition), resulting in a matrix of 127,464 spectra ($(264 \times 480) + 264 + 480$). In order to generate an artificial database, 20,000 perturbation spectra were randomly taken and added to 2000 repetitions of the 10 spectra from X_{cal} , resulting in a matrix X_{artif} . The reference protein values corresponding to this matrix y_{artif} were obtained from the 2000 repetitions of the reference protein value for the calibration dataset y_{cal} , because the perturbations added to the spectra should not change the reference value.

Fig. 2 shows that the spectral space of the calibration database increased when comparing from X_{cal} to X_{artif} using PCA (Principal Component Analysis). This was due to the integration of perturbation into the database. The model was then calibrated using X_{artif} and a second derivative using the Savitsky and Golay algorithm (window 9, polynomial order 3) ([16]) as pre-processing and wavelength selection (from 1950 to 2250 nm). Some non-linearity remained, but the model in cross-validation gave interesting performances. (Fig. 3).

3.1.2. Approach 2: OSC approach on artificial data

Orthogonalisation techniques are an interesting alternative for building robust models, i.e. independent as much as possible to potential perturbation factors. They include the External Parameter Orthogonalisation (EPO) [17] approach (see participants 2 and 3). The Orthogonal Signal Correction (OSC) approach takes account of the reference values, since the orthogonalisation is performed against the

y value to be predicted [15]. OSC based on an experimental design was proposed to take advantages of both EPO and OSC [18], i.e. trying at the same time to be insensitive to identified perturbation factors by removing the corresponding noise, without lowering the y prediction ability. In this application, however, although perturbation variability was present in the experimental designs, there was no variability in the reference values. The OSC approach was therefore proposed as a pre-treatment based on the previously built artificial calibration database, where reference values and perturbation variability were present. The multivariate directions hence identified by applying OSC on the artificial calibration database (containing the perturbation variability) were used to perform orthogonalisation of the initial calibration database.

The other pre-treatments used were a second derivative using the Savitsky and Golay algorithm with a window of 9 and a polynomial of degree 3, followed by an SNV (Standard Normal Variate) on the whole wavelength range. Cross-validation optimised the PLS calibration model with only two latent variables and one OSC-factor removed, giving a parsimonious model.

3.2. Participant 2

The method chosen for calibration was PLS. The first optimisation of the model was carried out by testing different pre-processing methods, using the calibration error (SEC) and the cross-validation leave-one-out error (SECV) as selection criteria. Because the spectra were affected by a baseline and a multiplicative effect, the pre-treatment used was a second derivative using the Savitsky and Golay algorithm with a window of 9 and a polynomial of degree 3 ([16]), followed by SNV ([19]), as explained in [20]. Two models were retained: one was based on the raw spectra, and the other was based on the pre-processed spectra, as previously described.

Both models were then submitted to a test for each perturbation factor, independently from other factors, in order to determine its influence. To do this, the matrices of experiments were averaged according to the perturbation factors not involved, in order to retain the variability only for the matrix whose effect needed to be evaluated. Both models were then tested on this reduced matrix. For example, to study the influence of moisture, the experimental design was averaged in terms of granulometry, temperature and repetitions, which supplied 22 spectra: 11 for the dry products and 11 for the humid products. The predictions for each sample at two moisture contents were then compared.

The method used to improve the robustness of the calibration was based on EPO ([17,21]). As explained in the Section 3.1, this technique involves the orthogonalisation of the measurement space of the spectra, on the basis of the disturbances caused by the perturbation factors. The orthogonalisation removes the components of the subspace spanned by the differences between the spectra repeated for the same sample. On one hand, the more the number of removed components increases, the more spectra of each sample become similar and then the calculated predictions on these spectra become more similar. On the other hand, if too many components are removed, this could alter the calibration. To choose the dimension of the space to be

Table 1

Validation criteria of the models after test for each factor individually (the left side corresponds to the raw samples and right side to the pre-treated data) for participant 2.

	Raw data					Pretreated data				
	R^2	Bias	SEP	Slope	Offset	R^2	Bias	SEP	Slope	Offset
Humidity	0.97	-26.90	7.10	0.82	2.85	1.00	1.93	5.66	0.80	27.60
Granulometry	0.99	-11.70	3.00	1.02	-14.70	0.99	-2.77	2.34	0.97	1.75
Temperature	1.00	-0.12	1.61	0.97	3.46	1.00	-0.33	1.17	1.02	-2.86
Repetition	1.00	0.21	1.00	1.01	-1.42	1.00	-0.06	0.52	1.00	-0.68
S1	0.83	-	5.73	-	-	0.91	-	4.12	-	-
S2	0.95	-	2.95	-	-	0.96	-	2.55	-	-

Table 2
Validation criteria of the models after test for each factor individually after orthogonalisation (the left side corresponds to the raw samples and right side to the pre-treated data) for participant 2.

	Raw data					Pretreated data				
	R^2	Bias	SEP	Slope	Offset	R^2	Bias	SEP	Slope	Offset
Humidity	0.99	0.05	2.51	1.01	-1.60	1.00	0.13	1.47	1.00	0.75
Granulometry	0.99	-0.27	2.46	1.03	-4.85	1.00	0.00	1.40	1.01	-1.68
Temperature	1.00	0.01	1.37	0.98	2.02	1.00	0.00	0.58	1.00	-0.44
Repetition	1.00	-0.20	0.47	1.00	-0.73	1.00	-0.07	0.28	1.00	-0.25
S1	0.98	-	1.77	-	-	0.98	0.00	1.69	-	-
S2	0.99	-	0.96	-	-	1.00	-	0.66	-	-

removed by orthogonalisation, we used the Wilks's lambda, which represents the ratio of the variance inter samples and the variance intra samples. This indicator was calculated in accordance with the number of removed factors and the number of latent variables of the model. The examination of its evolution enabled us to choose the optimal dimension to be removed.

Finally, the orthogonalisation of the models was achieved in two ways:

- By individual orthogonalisation, i.e. bringing together the six projectors (4 for X_{ed} , 1 for X_{instr1} and 1 for X_{instr2}) calculated individually for every perturbation factor.

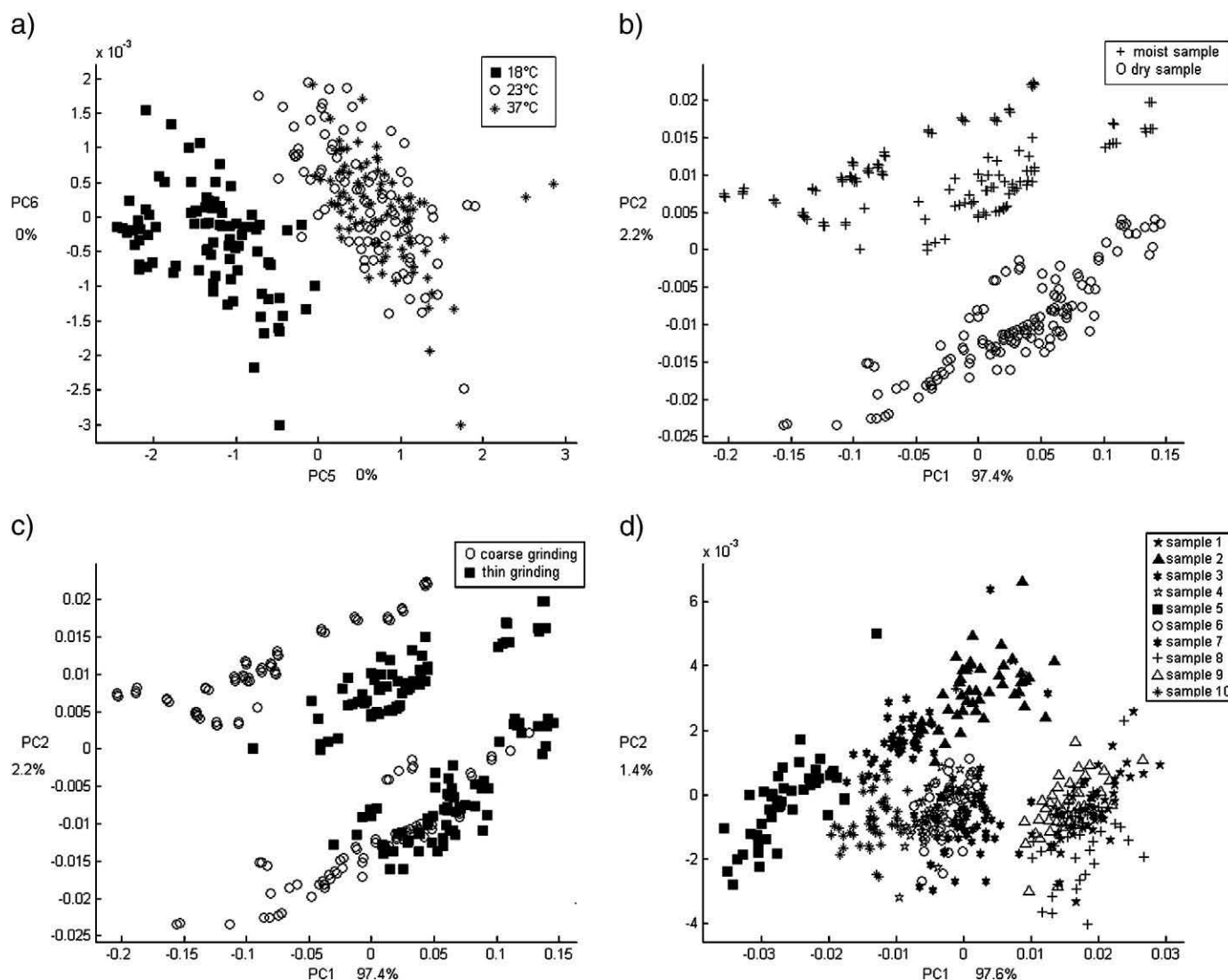


Fig. 4. PCA scores calculated without pre-processing, showing importance of external parameters: (a) temperature, (b) moisture, (c) grinding (d) repeatability on different NIR apparatus.

- By a global orthogonalisation, i.e. calculating a global projector, on all the X matrices brought together.

These two methods were applied to both models (raw and pre-treated), resulting in four different predictions for the test set. The three most different predictions were retained.

As previously explained, the models were submitted to a test for each perturbation factor independently to determine its influence. The sensitivity of these models to the different size factors is illustrated in Table 1 (the left side corresponds to the raw samples and right side to the pre-treated data). As shown in the table, the perturbation factors have varying influences on the robustness of the model:

- For moisture, there is clearly a positive effect of the pre-treatment: the bias is far less strong, as is the dispersion about the calibration line. However, the slope is far from 1 (± 0.8) in both cases.
- For granulometry, the pre-treatment reduces the bias and the dispersion about the calibration line.
- Temperature and sample repetition seem to be slightly influential. The pre-treatment improves the prediction, independently of the robustness.

The sensitivity of the models, after orthogonalisation, to the different size factors is illustrated in Table 2. It shows the clear advantage of orthogonalisation, especially for the most influential factors. The number of components that need to be removed is 3 for the factors of influence in the first experimental design, and 8 for the inter spectrometer repetitions.

3.3. Participant 3

As a first step, the goal was to show that the perturbation factors (i.e. moisture, grinding, temperature and repeatability on different NIR apparatus) were significant. To demonstrate this, PCA was performed after centering raw spectra. For each PCA, clusters were made up of each perturbation factor (Fig. 4).

A pre-processing method was then applied, its aim being to remove from the data space the part that was most influenced by the perturbation factor variations. As the previous participants, the method proposed was EPO [17], which estimates the parasitic subspace by computing a PCA on a set of spectra measured on the same objects, while the perturbation factor varies.

Different pre-treatments were applied: SNV [19], MSC [22] (multiplicative scatter correction), first derivative spectra using Savitsky Golay algorithm, and raw spectra. The first derivative spectra pre-treatment gave the best result.

As for participant 2, when applying EPO pre-processing, there are two possibilities: by individual orthogonalisation, i.e. eliminating perturbation factors on an ad hoc basis and by a global orthogonalisation. Both possibilities were tested and the second method was chosen because it gave the best PCA representation in that example. Fig. 5 shows the two first scores of PCA with EPO pre-processing on the first derivative spectra and on raw spectra. Comparing this with Fig. 4.d clearly shows the powerful performance of EPO pre-processing.

The PLS model was calibrated on the calibration matrix after EPO pre-treatment of the first derivative spectra. The Jack-Knife technique [23] was used to fix the required number of PLS factors for model construction. Cross-validation was applied in regression, so the optimal factor number was determined based on the prediction of the sample kept out of the individual model. A final model with two latent variables and an R^2 of 0.99 and RMSEC (root mean square error in calibration) of 3.30 was applied to the test dataset matrix.

3.4. Challenge organisers

Fig. 6 shows the spectra of the 10 available samples. The strategy proposed sought to create a huge calibration set by taking account of

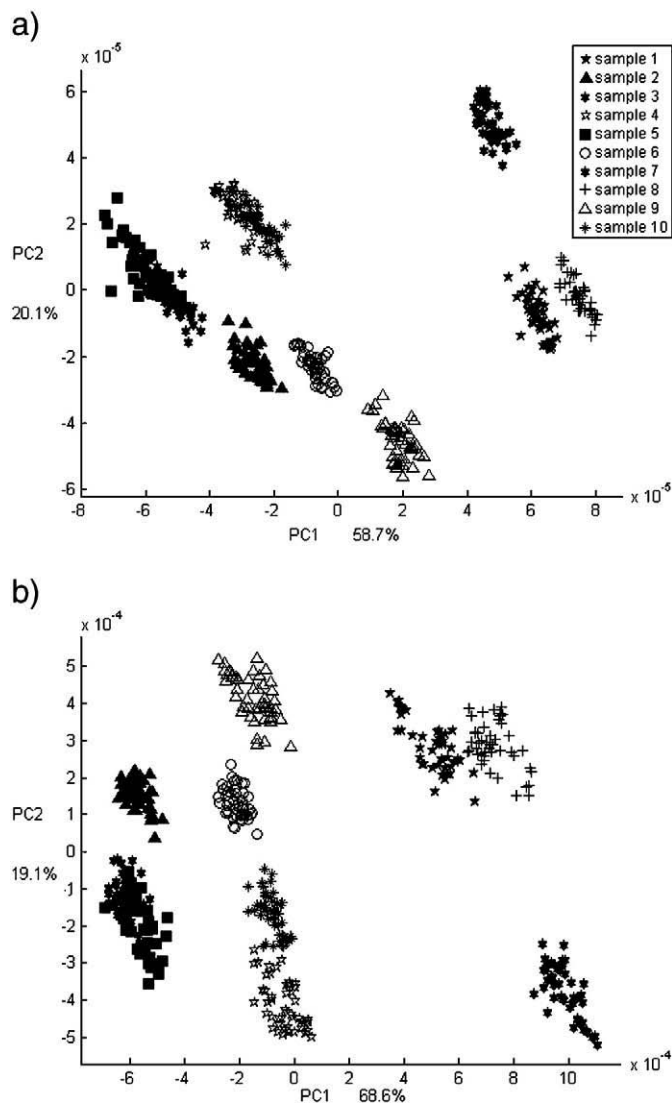


Fig. 5. Two first component scores of PCA with EPO pre-processing: (a) on 1st derivative spectra, (b) on raw spectra.

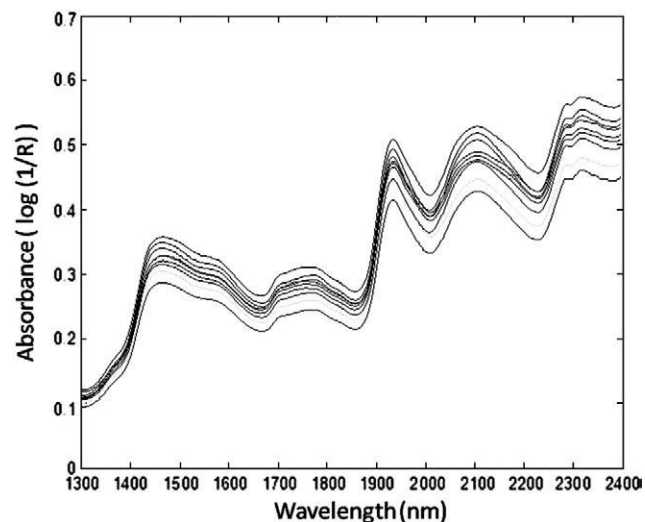


Fig. 6. Spectra of the 10 available samples (X_{cal}).

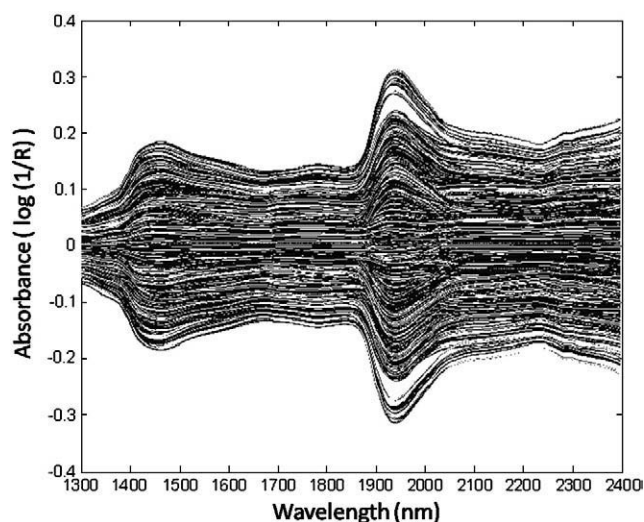


Fig. 7. Differences between the pairs of samples based on the experimental design dataset.

the perturbation factors, i.e. the experimental design and the repetitions between instruments. The main idea was to compute all the differences between all the pairs of the same sample in both perturbation factors, and then add these differences to the 10 spectra with known Y values. The procedure can be described in two steps:

Step 1 The experimental design (X_{ed}) comprised 11 samples \times 2 moistures \times 2 grinders \times 3 temperatures \times 2 replicates scans. In total, 264 spectra were available. All the differences between the pairs of samples were computed, giving a total of 552 ($24 \times (24 - 1)$) differences per sample. Then, for the 11 samples, a matrix containing 6072 differences was retained. Fig. 7 shows all these differences.

Step 2 The differences between the two series of 10 samples scanned on different instruments were computed. For set 1 (X_{instr1}) scanned on 31 instruments, 930 ($31 \times (31 - 1)$) differences were calculated for each sample. For set 2 (X_{instr2}) scanned on 17 instruments, 272 ($17 \times (17 - 1)$) differences for each sample were obtained. For the 10 samples, a total of 12,020 differences were computed, as shown in Fig. 8.

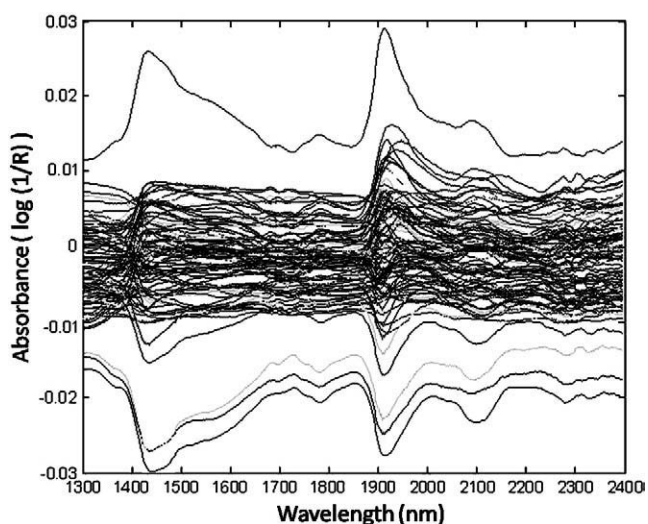


Fig. 8. Differences between the two series of samples scanned based on the Standard Replicates Dataset.

All these differences were added to the 10 available spectra (X_{cal}), giving a total of 180,920 spectra with only 10 reference protein values. Some of these calculated differences, mainly in the case of the standard replicates database, are larger than the rest, which was useful to include the variability present in those perturbation factors. Due to computer and time limitations, a random selection of spectra was made and the final model was constructed using 2548 spectra and 10 reference protein values. Fig. 9 shows the projection of the 2000 samples of the test dataset (X_{test}) on the 2548 spectra set constructed, indicating that both sets match perfectly.

The proposed model was constructed using least-squares support vector machines (LS-SVM) [24]. The results are shown in Fig. 10.

Fig. 10 provides a comparison in terms of RMSEC between the model using only the 10 available samples (Fig. 10 a) and the model constructed using the 2548 spectra obtained following the strategy described (Fig. 10 c). An additional parameter, RPD (ratio of the standard deviation of the population over the standard error of prediction), was also included. In both cases, the prediction for the test dataset is also shown (Fig. 10 b and d respectively).

4. Final results

The evaluation of the approaches was based on the best results obtained for the predicted values of the test dataset (2000 spectra). The reference protein values obtained by the reference method for these spectra were not communicated to the participants. For the evaluation, the R^2 , RMSEP (Root Mean Square Error in Prediction) and the RPD were used as validation criteria. The results for the different approaches are summarised in Table 3.

5. Conclusion

The challenge produced a wide diversity of results. However, when evaluating the 10 answers received (nine from the participants, one from the organisers) as a whole, Challenge 2007 showed that it is still possible to perform a robust and precise calibration when only a few reference values are available independently of the perturbation factors. The use of well-defined experimental designs, including repeated measurements under different conditions, will lead to the use of simple, easy and low-cost instruments and the construction of robust models.

During the congress the approaches summarised here were presented, together with the challenge organisers' approach. The

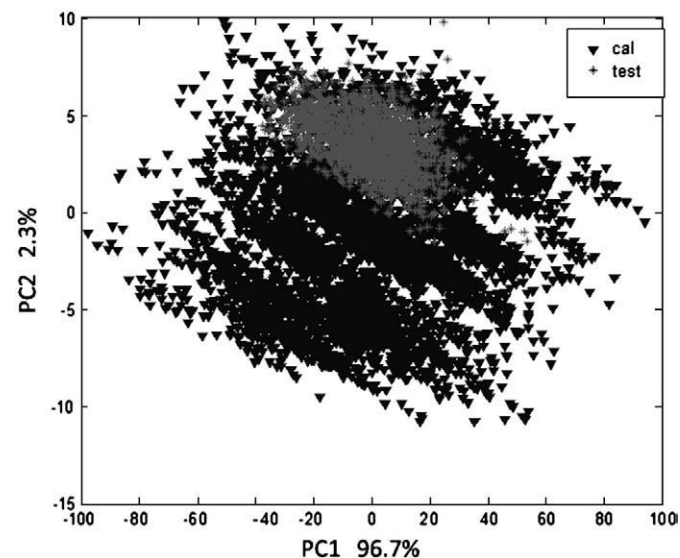


Fig. 9. PC1 vs. PC2 showing the projection of the 2000 samples of the X_{test} on the 2548 spectra dataset constructed by the challenge organisers.

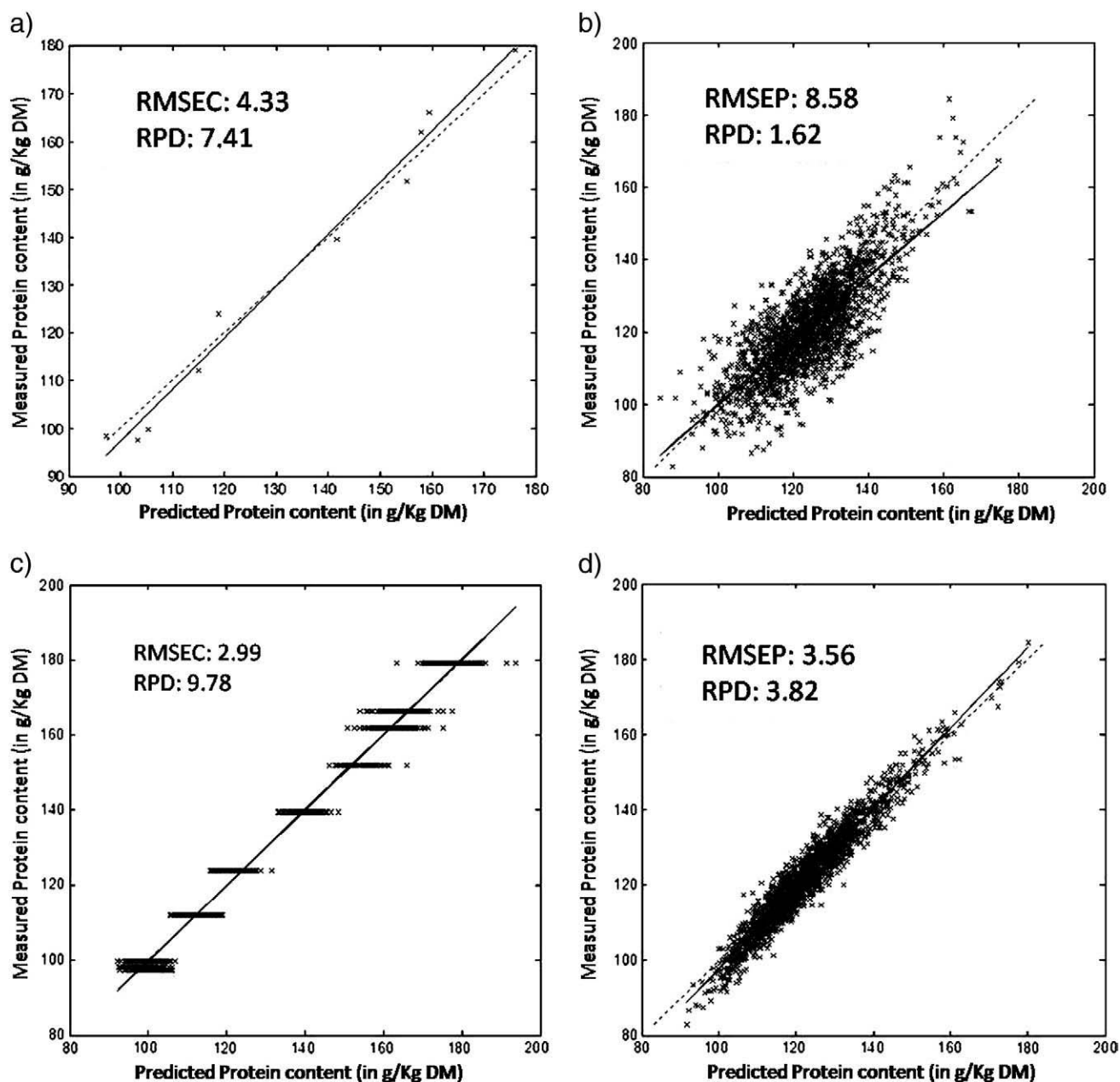


Fig. 10. Results obtained after the application of the SVM model for the challenge organisers. a) model using only the 10 available samples of X_{cal} ; b) prediction of X_{test} based on model represented in a; c) model constructed using the 2548 spectra obtained following the strategy described and d) prediction of X_{test} based on model represented in c.

participants found the results interesting and it was decided to include another challenge for the next congress.

The data of this challenge and previous challenges is available on the internet address of the French Chemometric Society (<http://www.chimiometrie.org/>).

Table 3

Summary of the results of the different approaches in terms of R^2 , RMSEP and RPD (ratio of the standard deviation of the population over the standard error of prediction).

	R^2	RMSEP	RPD
Participant 1 (approach 1)	0.91	3.91	3.32
Participant 1 (approach 2)	0.88	4.51	2.88
Participant 2	0.89	4.24	3.05
Participant 3	0.86	4.86	2.65
Challenge organizer	0.93	3.56	3.82

Acknowledgments

We would like to thank all the participants who spent time working on the data and presenting their results. Apart from the authors of this paper, the other participants in Challenge 2007 were Marion Cuny from Eurofins, Jean-Claude Boulet from INRA, Abdelaziz Faraj from IFP, Dominique Bertrand from INRA, Frédéric Estienne from Sanofi Aventis and Ludovic Duponchel from LASIR.

References

- [1] P. Dardenne, J.A. Fernández Pierna, A NIR data set is the object of a Chemometric contest at 'Chimiométrie 2004', Chemometrics and Intelligent Laboratory Systems 80 (2006) 236–242.
- [2] J.A. Fernández Pierna, P. Dardenne, Chemometric contest at 'Chimiométrie 2005': a discrimination study, Chemometrics and Intelligent Laboratory Systems 86 (2007) 219–223.

- [3] J.A. Fernández Pierna, P. Dardenne, Soil parameter quantification by NIRS as a Chemometric challenge at 'Chimiométrie 2006', *Chemometrics and Intelligent Laboratory Systems* 91 (2008) 94–98.
- [4] M. Daszykowski, K. Kaczmarek, Y. Vander Heyden, B. Walczak, Robust statistics in data analysis—a review. Basic concepts, *Chemometrics and Intelligent Laboratory Systems* 85 (2007) 203–219.
- [5] P.J. Huber (1981), *Robust Statistics*, John Wiley and Sons eds., New York.
- [6] I.N. Wakeling, H.J.H. MacFie, A robust PLS procedure, *Journal of Chemometrics* 6 (1992) 189–198.
- [7] M.I. Griep, I.N. Wakeling, P. Vankeerberghen, D.L. Massart, Comparison of semirobust and robust partial least squares procedures, *Chemometrics and Intelligent Laboratory Systems* 29 (1) (1995) 37–50.
- [8] S. Serneels, C. Croux, P. Filzmoser, P.J. Van Espen, Partial robust M-regression, *Chemometrics and Intelligent Laboratory Systems* 79 (1–2) (2005) 55–64.
- [9] J. González, D. Peña, R. Romera, A robust partial least squares regression method with applications, *Journal of Chemometrics* 23 (2) (2009) 78–90.
- [10] B. Liebmann, P. Filzmoser, K. Varmuza, Robust and classical PLS regression compared, *Journal of Chemometrics* 24 (3–4) (2010) 111–120.
- [11] D.W. Gu, P.H. Petkov, M.M. Konstantinov, *Robust Control Design with MATLAB®*, in: M.J. Grimble, M.A. Johnson (Eds.), *Advanced Textbooks in Control and Signal Processing*, Springer, Glasgow, UK, 2005.
- [12] M. Zeaiter, J.M. Roger, V. Bellon-Maurel, D.N. Rutledge, Robustness of models developed by multivariate calibration. Part I: The assessment of robustness, *Trends in Analytical Chemistry* 23 (2) (2004) 157–170.
- [13] M. Zeaiter, J.M. Roger, V. Bellon-Maurel, Robustness of models developed by multivariate calibration. Part II: The influence of pre-processing methods, *Trends in Analytical Chemistry* 24 (5) (2005) 437–445.
- [14] P. Dardenne, G. Sinnaeve, L. Bollen, R. Biston, Reduction of wet chemistry for NIR calibrations, in: G.D. Batten, P.C. Flinn, L.A. Welsh, A.B. Blakeney (Eds.), *Leaping ahead with Near Infrared Spectroscopy: Proceedings of the 6th ICNIRS*, NIR Spectroscopy Group, Royal Australian Chemical Institute, Melbourne, Victoria, Australia, 1994.
- [15] S. Wold, H. Antti, F. Lindgren, J. Öhman, Orthogonal signal correction of near-infrared spectra, *Chemometrics and Intelligent Laboratory Systems* 44 (1–2) (1998) 175–185.
- [16] P.A. Gorry, General least-squares smoothing and differentiation by the convolution (Savitzky–Golay) method, *Analytical Chemistry* 62 (1990) 570–573.
- [17] J.M. Roger, F. Chauchard, V. Bellon-Maurel, EPO-PLS external parameter orthogonalisation of PLS : application to temperature-independent measurement of sugar content of intact fruits, *Chemometrics and Intelligent Laboratory Systems* 66 (2) (2003) 191–204.
- [18] S. Preys, J.M. Roger, J.C. Boulet, Robust calibration using orthogonal projection and experimental design. Application to the correction of the light scattering effect on turbid NIR spectra, *Chemometrics and Intelligent Laboratory Systems*. 91 (2008) 28–33.
- [19] R.J. Barnes, M.S. Dhanoa, S.J. Lister, Standard normal variate transformation and de-trending of near-infrared diffuse reflectance spectra, *Applied Spectroscopy* 43–5 (1989) 772–777.
- [20] A.M.C. Davis, T. Fearn, Back to basics: removing multiplicative effects (1), *Spectroscopy Europe* 19 (4) (2007) 24–28.
- [21] A. Andrew, T. Fearn, Transfer by orthogonal projection: making near-infrared calibrations robust to between-instrument variation, *Chemometrics and Intelligent Laboratory Systems* 72 (1) (2004) 51–56.
- [22] T. Isaksson, T. Naes, The effect of multiplicative scatter correction (Msc) and linearity improvement in Nir spectroscopy, *Applied Spectroscopy* 42 (1988) 1273–1284.
- [23] J. Riu, R. Bro, Jack-knife technique for outlier detection and estimation of standard errors in PARAFAC models, *Chemometrics and Intelligent Laboratory Systems* 65 (2003) 35.
- [24] R.P. Cogdill, P. Dardenne, Least-squares support vector machines for chemometrics: an introduction and evaluation, *Journal of Near Infrared Spectroscopy* 12 (1) (2004) 93–100.