

How to build a robust model with only a few reference values: a chemometric challenge at 'Chimiométrie 2007'

Juan Antonio Fernández Pierna¹, Fabian Chauchard², Sébastien Preys², Jean Michel Roger³,
Oswin Galtier⁴, Vincent Baeten¹, Pierre Dardenne^{1*}

¹Walloon Agricultural Research Centre (CRA-W), Quality of Agricultural Products Department,
Chaussée de Namur n°24, 5030 Gembloux, Belgium

²Ondalys, Montpellier, France

³Cemagref, Montpellier, France

⁴University Paul Cézanne, Marseille, France

* Corresponding author. E-mail address: dardenne@cra.wallonie.be

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 Chimiométrie 2007 event (<http://www.chimiometrie.fr/>) 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 with only a few reference values available. The committee received nine answers; this paper summarizes the best three approaches, as well as the approach proposed by the organisers.

1. Introduction

For the fourth consecutive year and following on from the success of the chemometric contests organised during previous congresses [1, 2, 3], another dataset was proposed for the 'Chimiométrie 2007' meeting (<http://www.chimiometrie.fr/>) 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 when only a few reference values were available. Regression algorithms, as PLS (Partial Least Squares) are very sensible to outlying observations, which are typically expected to be present in experimental data. One of the solutions for this drawback of the classical regression techniques is the possibility of constructing robust versions of the regression algorithms [4-9]. In the early 1990s, Dardenne et al. [10] 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. Material and Methods

Several datasets were provided to the participants: a calibration dataset, an experimental design dataset, a standard replicate 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 perform a calibration model as robust and precise as possible with the available 10 reference values for protein content.

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 values). These 11 whole grain samples were separated into two homogeneous sets: one was dried to reduce the moisture content to +/-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 was 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 was thus available. The reference values corresponding to the experimental design samples and to the 10 samples measured on the different instruments were unknown.

2.4 Test Dataset

Some 2,000 spectra 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 2,000 spectra each had a reference protein value obtained by the reference method, but these values were not communicated to the participants.

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 $\mathbf{X}_{\text{instr1}}$ (31 instruments) and $\mathbf{X}_{\text{instr2}}$ (17 instruments). The Test dataset was kept in \mathbf{X}_{test} .

3. Results

3.1. Participant 1

Two approaches were tested: exhaustive calibration on artificial data, and Orthogonal Signal Correction (OSC) [11] on these artificial data.

Approach 1: Exhaustive calibration on artificial data

This approach involved using information about perturbation and ‘injecting’ it into the 10 spectra of the calibration dataset (\mathbf{X}_{cal}) in order to generate an artificial calibration database $\mathbf{X}_{\text{artif}}$. As this new calibration database included all possible perturbations, PLS would automatically be able to find the most robust direction for prediction.

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 ($\mathbf{X}_{\text{instr1}} + \mathbf{X}_{\text{instr2}}$) for each of the 24 spectra ($\mathbf{X}_{\text{sample } i}$) of the same sample i (one sample i had 24 ‘ j ’ combinations of possible perturbations). The 24 vectors of perturbations were calculated as follows:

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

where $\bar{\mathbf{x}}_i$ is the mean spectra of sample i and $\mathbf{x}_{i,j}$ is the j^{th} spectra of sample i . All the perturbations were put in a new matrix \mathbf{D}_{ed} containing the 264 possible perturbations. The same procedure was followed in order to identify all possible perturbations from the instruments, resulting in a matrix $\mathbf{D}_{\text{instr}}$ of 480 spectra. Figure 1 shows an example of the perturbation spectra that could be obtained.

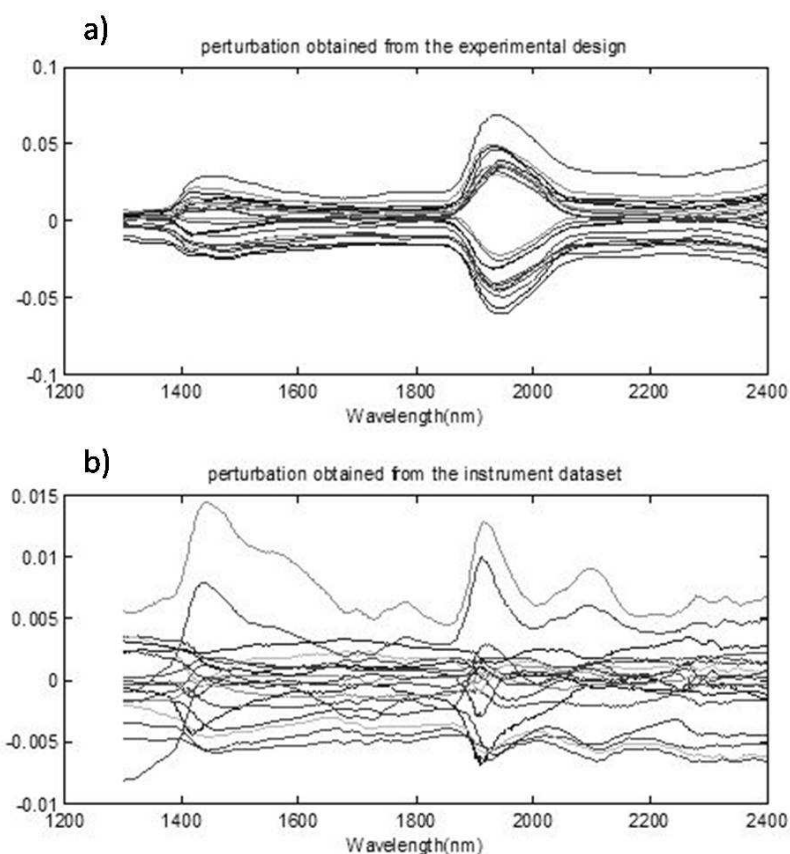


Figure 1: Examples of perturbation spectra from a) the experimental design dataset and b) the standard replicated dataset.

Finally, all the perturbations from \mathbf{D}_{ed} and $\mathbf{D}_{\text{instr}}$ were combined (using addition), resulting in a matrix of 127,464 vectors ($((264 \times 480) + 264 + 480)$). In order to generate an artificial database, 20,000 perturbations were randomly taken and added to 2,000 repetitions of the 10 spectra from \mathbf{X}_{cal} , resulting in a matrix $\mathbf{X}_{\text{artif}}$. The reference values corresponding to this matrix $\mathbf{y}_{\text{artif}}$ were obtained from the 2,000

repetitions of the reference value for the calibration dataset y_{cal} , because the perturbations added to the spectra should not change the reference value.

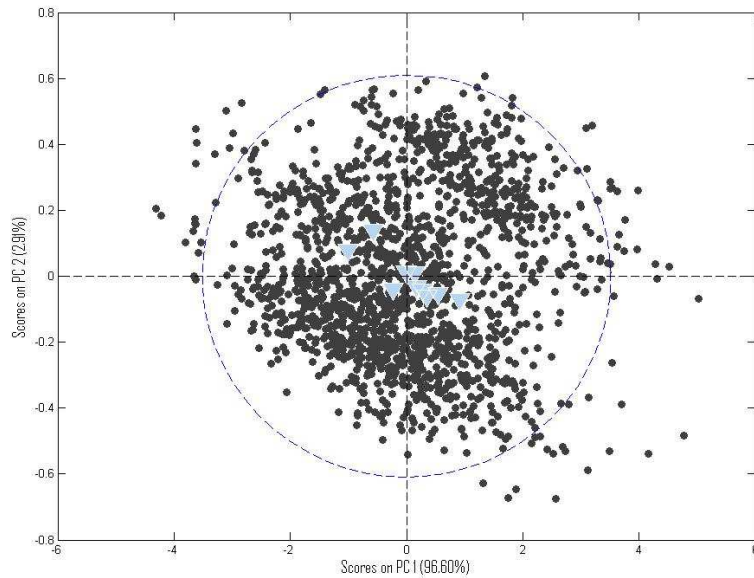


Figure 2 : PCA plot showing the spectral space of the increased calibration dataset (X_{artif}) including the original X_{cal} .

Figure 2 shows that the spectral space of the calibration database increased when comparing from X_{cal} to X_{artif} using PCA. This was due to the integration of perturbation into the database. The model was then calibrated using X_{artif} and a second derivative (window 9, polynomial order 3) as pre-processing and wavelength selection (from 1950 to 2250 nm). Some non-linearity remained, but the model in cross-validation gave interesting performances. (Figure3).

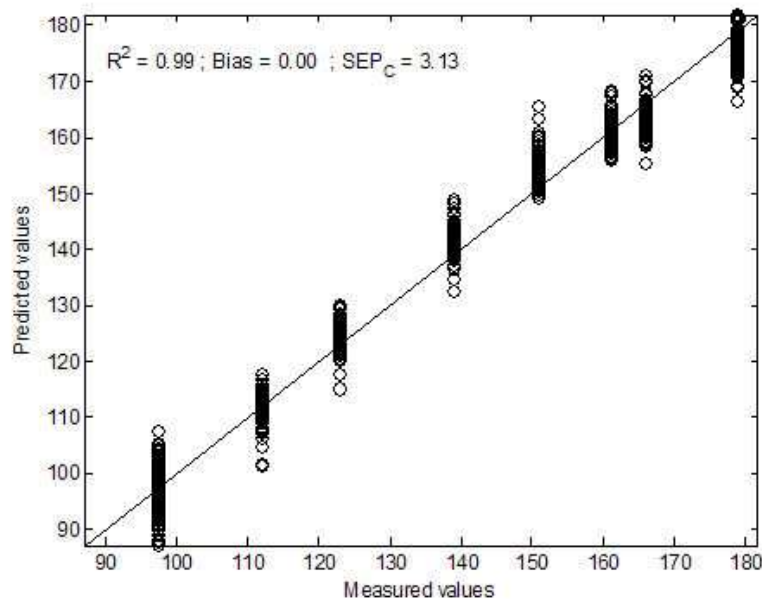


Figure 3 : Final model constructed using X_{artif} for participant 1.

Approach 2: OSC approach on artificial data

Orthogonalisation techniques are an interesting alternative for building robust models. They include the External Parameter Orthogonalisation (EPO) [12] approach (see participants 2 and 3). The Orthogonal Signal Correction (OSC) approach takes account of the reference values [11]. OSC based on an experimental design was proposed to take advantages of both EPO and OSC [13]. 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 on the previously built artificial calibration database, where reference values and perturbation variability were present.

The other pre-treatments used were a second derivative using the Savitsky and Golay algorithm ([14]) with a window of 9 and a polynomial of degree 3, followed by an SNV on the whole wavelength range. Cross-validation optimised the 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 ([14]), followed by SNV ([15]), as explained in [16]. 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 factor, independently from other factors, in order to determine its influence. To do this, the matrices of experiments were averaged according to the 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 ([12, 17]). This involves the orthogonalisation of the measurement space of the spectra, on the basis of the disturbances caused by the factors. The orthogonalisation removes the components of the sub-space 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 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 size of influence
- 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 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 factors have varying influences on the robustness of the model:

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				
	R2	Bias	SEP	Slope	Offset	R2	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	-	-

- 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.

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				
	R2	Bias	SEP	Slope	Offset	R2	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	-	-

3.3. Participant 3

As a first step, the goal was to show that the external parameters (e.g., moisture, grinding, temperature, and repeatability on different NIR apparatus) were

significant. To demonstrate external parameter importance, PCA was performed after centering raw spectra. For each PCA, clusters were made up of each external parameter (Figure 4).

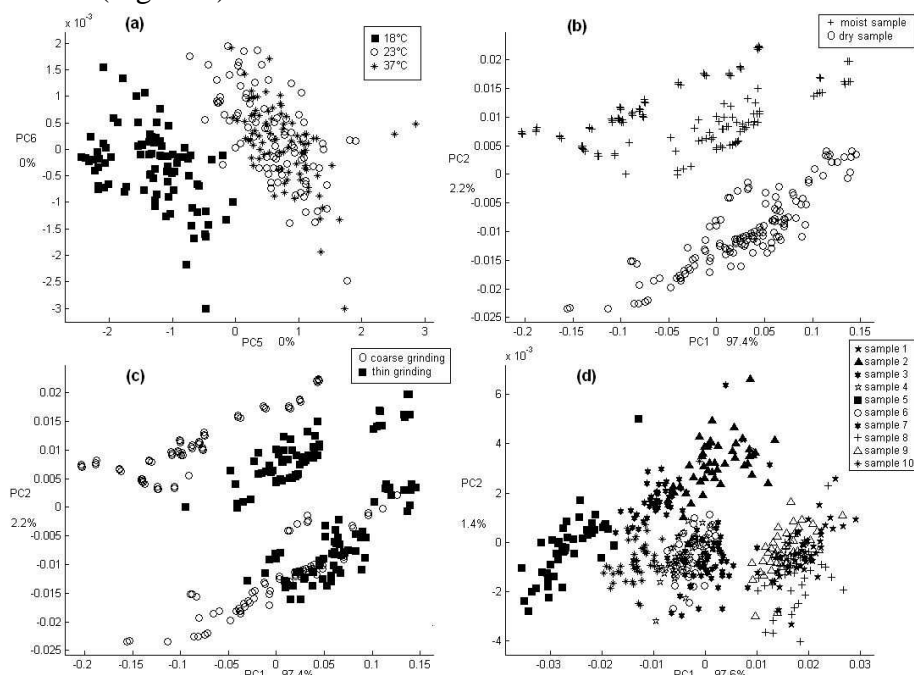


Figure 4 : PCA scores calculated without pre-processing, showing importance of external parameters: (a) temperature, (b) moisture, (c) grinding (d) repeatability on different NIR apparatus.

A pre-processing method was then applied, its aim being to remove from the data space the part that was most influenced by the external parameter variations. The method proposed was EPO [12], which estimates the parasitic subspace by computing a PCA on a set of spectra measured on the same objects, while the external parameter varies.

Different pre-treatments were applied: SNV [15] (standard normal variate), MSC [18] (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 external parameters 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. Figure 5 shows the two first scores of PCA with EPO pre-processing on the first derivative spectra and on raw spectra. Comparing this with Figure 4.d clearly shows the powerful performance of EPO pre-processing.

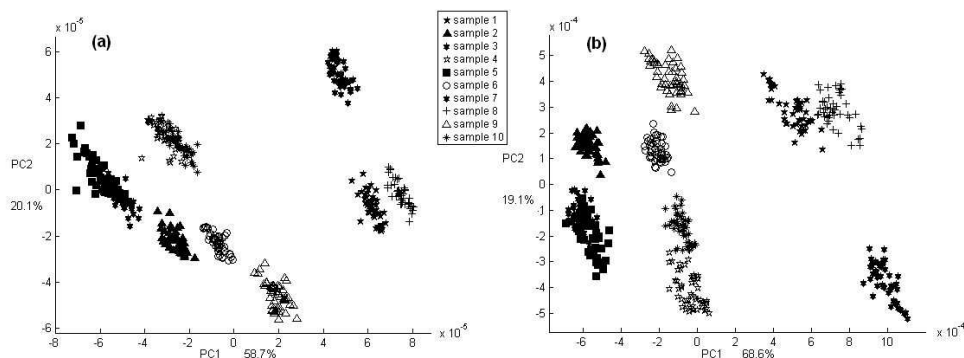


Figure 5 : Two first component scores of PCA with EPO pre-processing: (a) on 1st derivative spectra, (b) on raw spectra.

The PLS model was calibrated on the calibration matrix after EPO pre-treatment of the first derivative spectra. The Jack-Knife technique [19] was used to fix the required number of 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

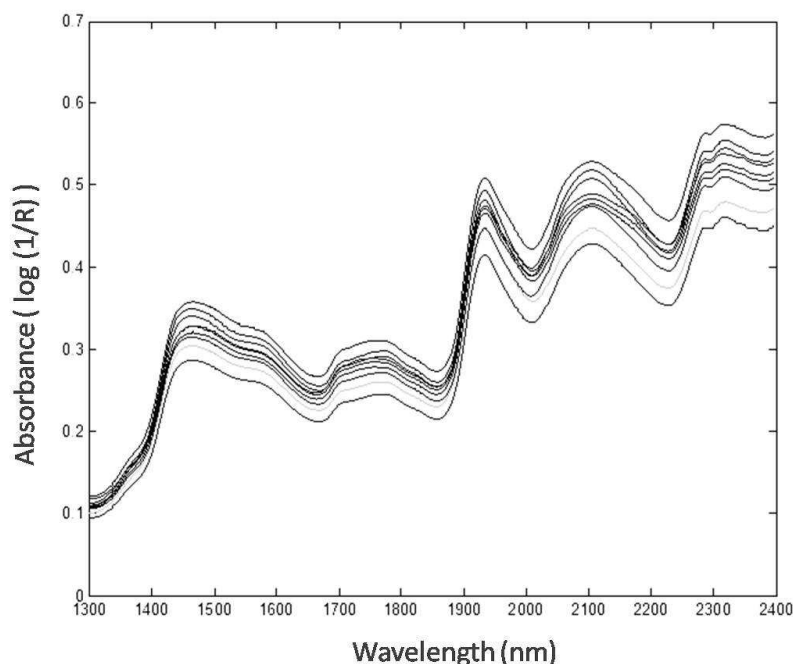


Figure 6 : Spectra of the 10 available samples (X_{cal}).

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

and the repetitions, 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 (\mathbf{X}_{ed}) comprised 11 samples x 2 moistures x 2 grinders x 3 temperatures x 2 replicated scans. In total, 264 spectra were available. All the differences between the pairs of samples were computed, giving a total of 552 ($24 * (24-1)$) differences per sample. Then, for the 11 samples, a matrix containing 6,072 differences was retained. Figure 7 shows all these differences.

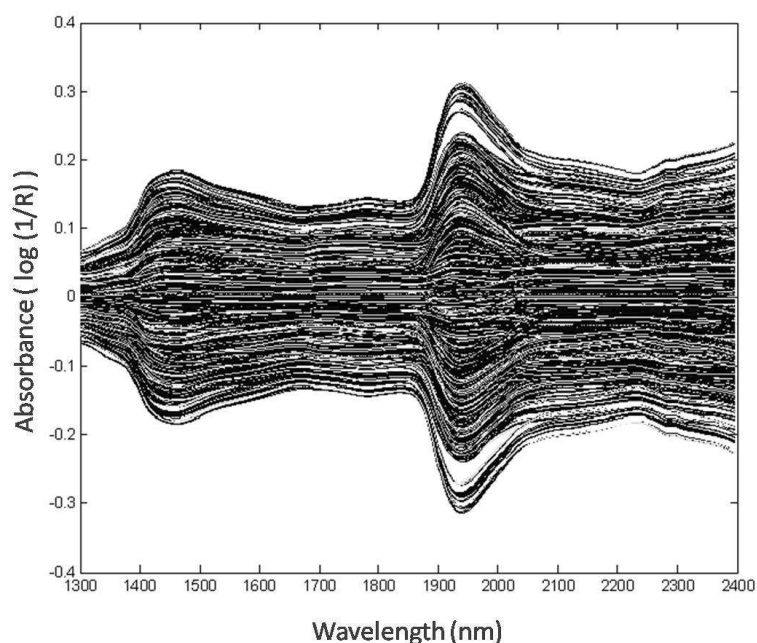


Figure 7 : Differences between the pairs of samples based on the experimental design dataset.

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

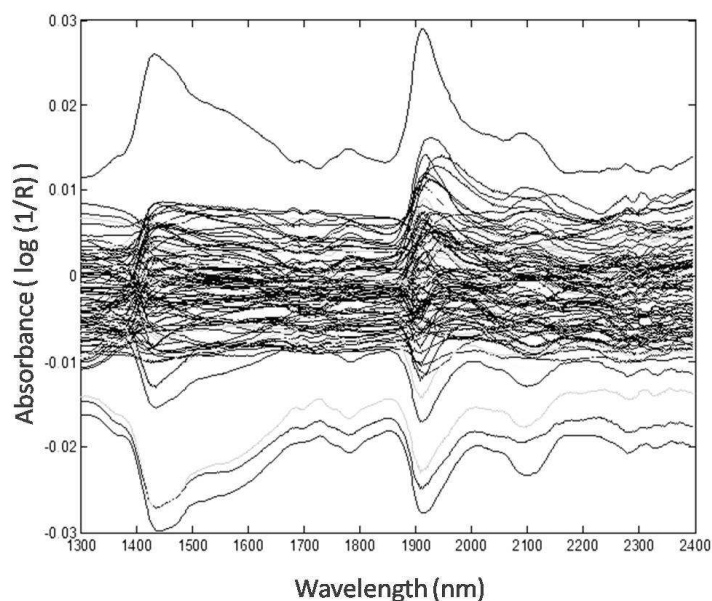


Figure 8 : Differences between the two series of samples scanned based on the Standard Replicates Dataset.

All these differences were added to the 10 available spectra (\mathbf{X}_{cal}), giving a total of 180,920 spectra with only 10 reference values. Due to computer and time limitations, a random selection of spectra was made and the final model was constructed using 2,548 spectra and 10 reference values. Figure 9 shows the projection of the 2,000 samples of the test dataset (\mathbf{X}_{test}) on the 2,548 spectra set constructed, indicating that both sets match perfectly.

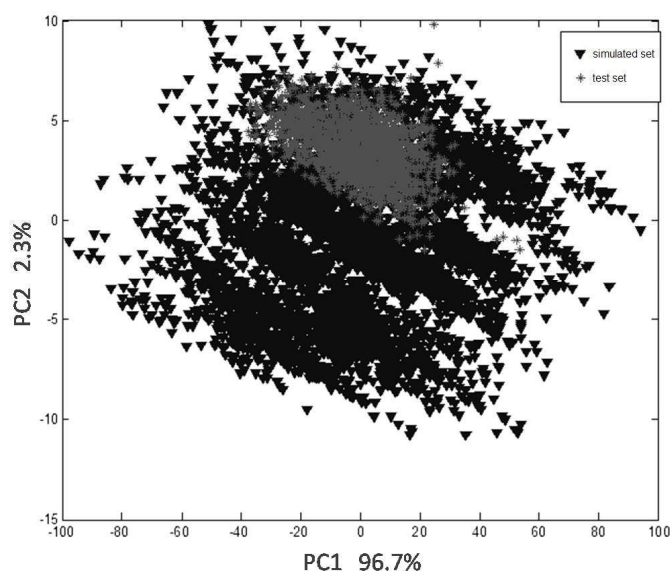


Figure 9 : PC1 vs. PC2 showing the projection of the 2,000 samples of the \mathbf{X}_{test} on the 2,548 spectra dataset constructed by the challenge organisers.

The proposed model was constructed using least-squares support vector machines (LS-SVM) [20]. The results are shown in Figure 10.

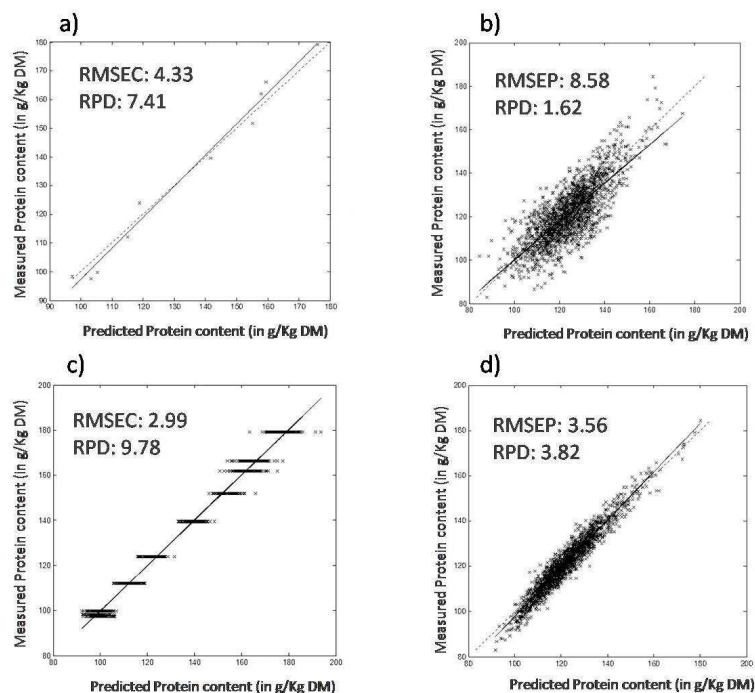


Figure 10 : Results obtained after the application of the SVM model for the challenge organisers. a) model using only the 10 available samples of X_{cal} ; b) prediction of X_{test} based on model represented in a; c) model constructed using the 2,548 spectra obtained following the strategy described and d) prediction of X_{test} based on model represented in c.

Figure 10 provides a comparison in terms of RMSEC and RPD between the model using only the 10 available samples (Figure 10 a) and the model constructed using the 2,548 spectra obtained following the strategy described (Figure 10 c). In both cases, the prediction for the test dataset is also shown (Figures 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 (2,000 spectra). The reference protein value obtained by the reference method for these spectra was not communicated to the participants. For the evaluation, the R^2 and RMSEP (Root Mean Square Error in Prediction) were used as validation criteria. An additional parameter, RPD (ratio of the standard deviation of the population over the standard error of prediction), was also included. The results for the different approaches are summarised in Table 3.

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).

	R2	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

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. 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 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.fr/>).

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.

6. References

- [1] P. Dardenne & J.A. Fernández Pierna (2006). 'A NIR data set is the object of a Chemometric contest at 'Chimiométrie 2004''. *Chemometrics and Intelligent Laboratory Systems*, 80, pp. 236-242.
- [2] J.A. Fernández Pierna & P. Dardenne (2007). 'Chemometric contest at 'Chimiométrie 2005': a discrimination study'. *Chemometrics and Intelligent Laboratory Systems*, 86, pp. 219-223.
- [3] J.A. Fernández Pierna & P. Dardenne (2008). 'Soil parameter quantification by NIRS as a Chemometric challenge at 'Chimiométrie 2006''. *Chemometrics and Intelligent Laboratory Systems*, 91, pp. 94-98.
- [4] P.J. Huber (1981), *Robust Statistics*, John Wiley and Sons, New York.
- [5] I.N. Wakeling & H.J.H. MacFie (1992), *Journal of Chemometrics*, 6, pp. 189-198.
- [6] M. I. Griep, I. N. Wakeling, P. Vankeerberghen & D. L. Massart (1995), 'Comparison of semirobust and robust partial least squares procedures', *Chemometrics and Intelligent Laboratory Systems*, 29 (1), pp. 37-50.
- [7] I.N. Wakelinc, H.J.H. Macfie, (2005) 'A robust PLS procedure', *Journal of Chemometrics* 6 (4), pp. 189-198.
- [8] S. Serneels, C. Croux, P. Filzmoser & P.J. Van Espen (2005) 'Partial robust M-regression'. *Chemometrics and Intelligent Laboratory Systems*, 79 (1-2), pp.55-64.
- [9] M. Daszykowski, K. Kaczmarek, Y. Vander Heyden & B. Walczak, (2007) 'Robust statistics in data analysis – a review. Basic concepts', *Chemometrics and Intelligent Laboratory Systems*, 85, pp. 203-219.
- [10] P. Dardenne, G. Sinnaeve, L. Bollen & R. Biston (1994). 'Reduction of wet chemistry for NIR calibrations'. In G. D. Batten, P. C. Flinn, L. A. Welsh and 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.

- [11] S. Wold, H. Antti, F. Lindgren & J. Öhman (1998). 'Orthogonal signal correction of near-infrared spectra'. *Chemometrics and Intelligent Laboratory Systems*, 44 (1-2), pp. 175-185,
- [12] J.M. Roger, Chauchard F. & V. Bellon-Maurel (2003). '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), pp. 191-204.
- [13] S. Preys, J.M. Roger, J.C. Boulet (2008). '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, pp. 28-33.
- [14] P.A. Gorry (1990). 'General least-squares smoothing and differentiation by the convolution (savitzky-golay) method'. *Anal. Chem.*, 62, pp. 570-573.
- [15] R. J. Barnes, M. S. Dhanoa & S. J. Lister (1989). 'Standard normal variate transformation and de-trending of near-infrared diffuse reflectance spectra'. *Appl. Spectrosc.*, 43-5, pp. 772-777.
- [16] A.M.C. Davis & T. Fearn (2007). 'Back to basics: removing multiplicative effects (1)'. *Spectroscopy Europe*, 19(4), pp. 24-28.
- [17] A. Andrew & T. Fearn (2004). 'Transfer by orthogonal projection: making near-infrared calibrations robust to between-instrument variation'. *Chemometrics and Intelligent Laboratory Systems*, 72(1), pp. 51-56.
- [18] T. Isaksson & T. Naes (1988). 'The Effect of Multiplicative Scatter Correction (Msc) and Linearity Improvement in Nir Spectroscopy'. *Applied Spectroscopy* 42, pp. 1273-1284.
- [19] J. Riu & R. Bro (2003). 'Jack-knife technique for outlier detection and estimation of standard errors in PARAFAC models'. *Chemometrics and Intelligent Laboratory Systems*, 65, pp. 35.
- [20] R. P. Cogdill & P. Dardenne (2004). 'Least-squares support vector machines for chemometrics: an introduction and evaluation'. *Journal of Near Infrared Spectroscopy* 12 (1), pp. 93-100.