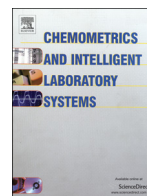




Contents lists available at ScienceDirect

Chemometrics and Intelligent Laboratory Systems

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

Use of a multivariate moving window PCA for the untargeted detection of contaminants in agro-food products, as exemplified by the detection of melamine levels in milk using vibrational spectroscopy

J.A. Fernández Pierna, D. Vincke, V. Baeten, C. Grelet, F. Dehareng, P. Dardenne *

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

ARTICLE INFO

Article history:

Received 22 July 2015

Received in revised form 13 October 2015

Accepted 30 October 2015

Available online xxxxx

Keywords:

Untargeted detection

Contaminant

Moving window

PCA

ABSTRACT

In this study, the concept of a Local moving window along the wavelength range in vibrational spectroscopic data was used to build reduced PCA models for characterizing agro-food products and detecting the presence of unusual ingredients or contaminants in an untargeted way. For each selected wavelength window in a locally reduced calibration set, a PCA analysis was performed and score residuals were extracted and used as to define thresholds to be applied to the spectral score residuals of the sample being investigated. When a residual at a certain wavenumber exceeded defined thresholds, the sample was suspected of being abnormal, indicating the possible presence of unusual ingredients and allowing non-targeted analysis. The method was applied to liquid UHT milk samples spiked with varying levels of melamine. Samples spiked at levels higher than 100 ppm were easily detected using this method, which would not have been possible using classical techniques such as Mahalanobis distance, usually applied as an outlier detection method.

© 2015 Elsevier B.V. All rights reserved.

1. Introduction

The concept of a moving window along one dimension in numerical data has long been used for different objectives. One of the main applications is the Savitzky–Golay filter, which is applied to a set of digital data points in order to smooth the data (i.e., to increase the signal-to-noise ratio without greatly distorting the signal) [1]. This technique is widely used as preprocessing tool in vibrational spectroscopy in order to avoid increasing the noise and complexity of the spectrum typically found when performing derivatives in the data. This moving window averaging technique fits successive subsets of adjacent data points with a low-degree polynomial using the linear least squares method. It uses information from a localized segment of the spectrum to calculate the derivative at a particular wavelength, rather than the difference between adjacent data points. In most cases, this avoids the problem of noise enhancement from the simple difference method.

The concept of a moving window is popular in the area of evolving data (e.g., environmental data). Techniques such as Evolving Factor Analysis (EFA) provide an estimation of the regions or windows where the concentration of the different components is evolving [2]. One example is the Fixed Size Moving Window EFA (FSMWFEFA), where the eigenvalues are calculated for sub-matrices of equal size moving in the same direction as the experiment is performed [3]. The moving window concept has also been applied in correlation

spectroscopy and time series analysis using methods such as Singular Spectrum Analysis (SSA) [4,5], also called Dynamic Principal Component Analysis (DPCA) [6], where the aim is to make a decomposition of the original series into the sum of a small number of independent and interpretable components such as a slowly varying trend, oscillatory components and structureless noise [7]. This is done by treating the signal as a set of repeated overlapping windows along the variable range. In the case of correlation spectroscopy, the dataset is split into a series of relatively small windows where their covariance maps in succession [8,9] or correlation coefficients [10] are calculated.

The concept of a moving window has been also used in combination with classical chemometric tools, such as Principal Component Analysis (PCA) or Partial Least Squares (PLS) regression, in the field of multivariate statistical process control (MSPC) where models possess the ability to automatically change their properties during online operations (adaptive models) [11], usually applied to vibrational spectroscopic data in order to select an optimal range of wavelengths, among others. Moving Window PCA (MWPCA) was first developed by Liu et al. [12] within the context of complex nonlinear time-varying processes with an algorithm called Moving Window Kernel PCA (MWKPCA), based on an iterative procedure for adapting the data mean and covariance matrix in the feature space and then approximating the eigenvalues and eigenvectors of the Gram matrix. Another version of MWPCA was proposed by Ryu et al. [13] to detect the presence of peak shift in spectra. With this technique, a moving window is constructed from a small data segment along the wavenumber axis where a PCA is performed in order to detect peak shifts and interpret highly correlated spectra.

* Corresponding author.

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

In the regression domain, methods such as Moving Window Partial Least Squares Regression (MWPLSR) [14] and, most recently, Moving Window Variable Importance in Projection (MW-VIP) [15] have been proposed and compared, with regard to reducing model errors, with classical methods for wavelength selection.

In this study, a MWPCA method was studied and applied to the characterization of various agro-food products using vibrational spectroscopic analysis tools such as mid-infrared (MIR). The objective was to exploit the huge amount of information contained in the data generated by such techniques, which could support the concept of data-driven discovery or untargeted analysis [16,17]. New crises of adulteration/contamination with illegal ingredients other than known ones continue to occur from time to time. By relying only on targeted analysis methods, adulteration could get out of control and analysis would become trapped in a cycle of 'adulteration, targeted analysis, and new adulteration', and so on [18]. In contrast to targeted analysis, which uses information from known possible unusual ingredients, an untargeted experiment registers all information within a certain correlation/similarity, including data from new products. Untargeted detection methods are therefore required for screening products for a range of known and unknown adulterants [19]. Untargeted analysis will mean alerts can be given more rapidly and fraud detected more easily. Until now, untargeted analysis has been associated mainly with direct analysis techniques, such as mass spectrometric-based metabolomics or isotope-assisted methods. Only a few studies have linked untargeted analysis with vibrational spectroscopic methods. Moore et al. [19] developed non-targeted screening tools to detect adulteration in skimmed milk powder using NIR spectroscopy; Xu et al. [20] investigated the feasibility of using FT-NIR spectroscopy and chemometrics for the rapid analysis of poplar balata in Chinese propolis and Lu et al. [18] developed a method for the untargeted detection of protein adulteration in yogurt by removing unwanted variations in pure yogurt. In all these cases, the approach involved building statistical models based on the measured fingerprints of a large representative set of normal and abnormal samples, and then applying these models to unknown samples in order to characterize them. More recently, the FOSS company (Foss, Hillerød, Denmark) has developed an Abnormal Spectrum Screening Module (ASM) where new milk samples are automatically compared to the spectra of the natural (not contaminated) historical dataset obtained with the MilkoScan™ FT120 (<http://www.foss.fr/industry-solution/products/milkoscan-ft1/>), then outliers are detected by a combination of the residuals from the PCA on natural samples and the Mahalanobis distance [21].

For this study, a moving window was selected along the wavelength axes in vibrational spectroscopic data. For each selected window in the calibration stage, PCA was performed by fixing the number of principal components and applying them to a validation or test set. The spectral score residuals in the calibration set were extracted and used to define thresholds to be applied to the spectral score residuals of the validation set. When a residual at a certain wavenumber exceeded the defined thresholds, the sample was suspected of being abnormal, indicating the possible presence of unusual ingredients and allowing untargeted analysis. A key challenge in all studies in this area is to define 'normal' and 'abnormal' according to fingerprint properties. In this study, this was solved by using a local technique that allowed, for each sample, the most spectroscopically similar samples in the calibration set to be selected before the application of MWPCA.

The study used, as an example, is the case of milk contaminated with melamine. Melamine (2,4,6-triamino-1,3,5-triazine) is a chemical compound rich in nitrogen. When combined with formaldehyde, it produces melamine resin, which is widely used in textiles, plastics, adhesives, flame-resistant products and some cleaning agents. Melamine has been illegally added to food/feed to artificially elevate the protein content value of products [22–24]. Since the discovery of melamine contamination in infant milk formula in China, strict regulations have been enforced throughout the world and many papers have been

published on the use of such methods as wet chemistry, chromatography, mass spectrometry and vibrational spectroscopy to detect melamine in both raw and powdered milk [25,26]. In this study, liquid ultra-high temperature (UHT) milk was contaminated with melamine at various levels ranging from 0.01% to 1% (100–10,000 ppm) and measured using Fourier Transform Mid-Infrared (FT-MIR) spectroscopy in order to test the performance of the proposed Local MWPCA method and determine its limits of detection.

2. Local moving window PCA

The general principle of the proposed methodology is to compute, for each sample, the PCA residuals at different windows throughout the spectrum range and compare them with those in the calibration set. In order to get better performances, as an initial step, the concept of a local method is applied to each unknown sample [27]. This involves using a nonparametric method based on a subset of n spectra selected by the correlation coefficient between the spectrum of the unknown and the 'clean' spectra of the calibration database. This enables a 'local' PCA model to be built at each window of the procedure. The next step involves specifying a window through the wavelength range where only the n calibration spectra previously selected are used to build a PCA model. The PCA scores and loadings are extracted to reconstruct each wavelength of the window and compute the residuals for each sample. The residual of the central point is computed and stored. This step is repeated for each window until the whole spectral wavelength range has been covered. For each PCA model built for calibration, the corresponding scores and loadings are saved for the prediction of the unknown spectrum and the computation of its central point residuals. Once the whole spectral range has been covered, the residuals of each window are set at an absolute value and scaled by mean and standard deviation. A standard deviation amplifier is then set in order to define the residual limits. Residuals are the noisy part of PCA and correspond to the part of the spectra that is not explained by the model, and they usually have small values. The residuals of each wavelength of the unknown spectrum are therefore examined in order to compute the number of residuals surpassing the limit, known as the 'hit number'. If the hit number is greater than a defined threshold, the sample is considered to be lying outside the local calibration and is considered to be 'abnormal' or 'unusual'.

In summary, the procedure is as follows:

- 1) Set up a library of 'clean' samples (reference set) of a given product.

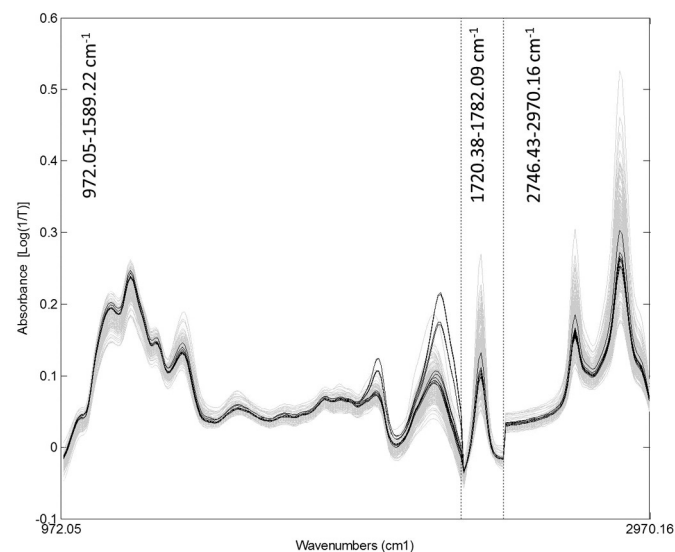


Fig. 1. FT-MIR spectra of the 'clean' liquid milk dataset (in grey) and the 12 samples spiked with melamine (in black).

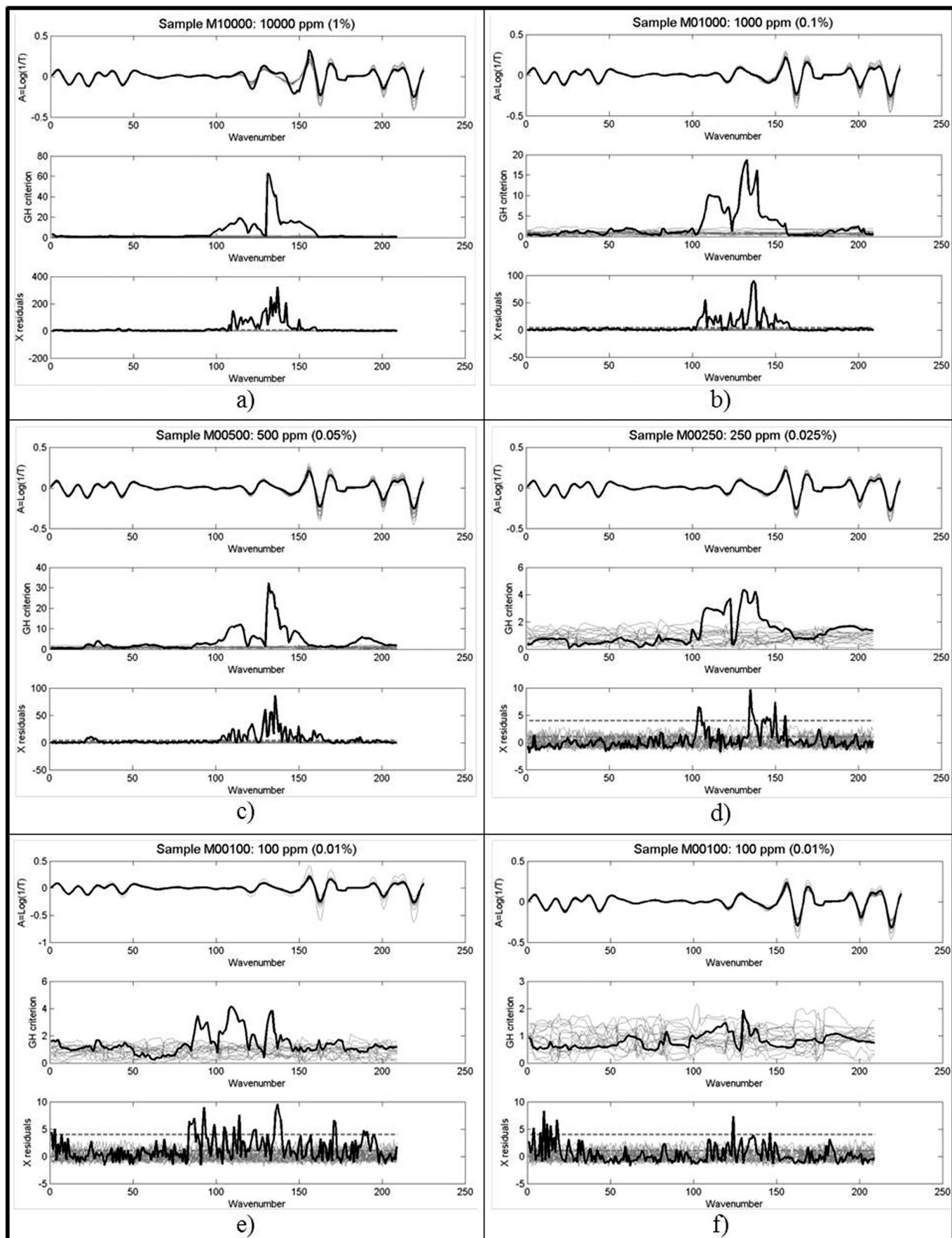


Fig. 2. First derivative spectra, the GH criterion and the LWPCA results, respectively, for each locally selected sample in the calibration set (in grey) and the spectrum of the sample to be characterized (in black).

- 2) For each unknown new spectrum, search the n closest spectra from the clean dataset based on the highest correlation with each spectrum of the library. The n value depends on the variability covered by the samples in the clean dataset. A large range of n values should be tested and externally validated to select the optimal one.
- 3) Compute a PCA model with the selected spectra for each window along the wavelength range. The choice of the width of the window is a crucial step, however as far as we know there are no general guidelines. In Press et al. [28] a rough guideline is proposed, which indicates that the maximum number of data points that should be used can be calculated based on the wavelength interval used and the peak width of the narrowest peak or FWHM (full width at half of maximum). The number of PCA components has been selected according to a previous global PCA performed in the clean dataset. In such PCA the number of components is selected taking into account only those that explain more than 99% of the variability of the clean dataset.
- 4) Compute the residual limits of the PCA model for each central point of the windows; absolute values of the residuals are standardized at each wavelength in order to reach a unique limit of 3 for all the wavelengths. This value corresponds to the 99.7% confidence interval in the clean dataset considering the hypothesis of Gaussianity of the residuals.
- 5) Use the PCA model to project the unknown new spectrum and compute its residuals.
- 6) Check if these residuals are within PCA residual limits by counting the number of times the limits are exceeded and the sum of the exceeded values above the threshold.
- 7) Restart the procedure at step 2 for the next unknown new spectrum.

In addition to the residuals, the Global H (GH) criterion [29] was also computed for the central point of calibration and prediction sets at each interval. This GH criterion is a modification of the Mahalanobis distance, H , in which $H \cdot H$ (H squared) is divided by the number of dimensions, used to derive H . The practical significance of the GH is that it indicates when a predicted value for a given sample is outside the limits defined by the population that make up the calibration set. For the GH, the limit of reliable results is usually 3.00, but not a general value exists as it can (slightly) decrease with the number of dimensions [30]. Samples containing variables with GH higher than 3.00 should be regarded as possible outliers or a possible case of contamination/adulteration.

3. Experiment

3.1. Materials and method

In order to assess the performance of this technique, a dataset of 300 samples of raw milk and processed liquid UHT milk was used as a 'clean' dataset. Another 12 liquid UHT milk samples were contaminated with melamine at various levels ranging from 0.01% to 1% (100–10,000 ppm). All these samples were measured in triplicate using a Standard Lactoscope FT-MIR automatic (Delta Instruments, Drachten, The Netherlands) in the 397.31–4000 cm^{-1} range with a resolution of 8 cm^{-1} . The spectra were cleaned up in order to remove irrelevant information and the final wavenumber ranges used (in cm^{-1}) were 972.05–1589.22, 1720.38–1782.09 and 2746.43–2970.16. Fig. 1 shows all the FT-MIR spectra of the 'clean' milk (in grey) and the 12 spiked samples. No clear evidence of contamination could be observed directly from the raw spectra.

The procedure described in the previous section was applied and used to detect adulteration in the 12 spiked samples. For this specific case, the optimal parameters for the Local moving window PCA have been determined as previously explained. For the local step, after optimization, a subset of $n = 15$ spectra from the clean dataset has been selected according to its correlation with each unknown spectrum. When applying the full width at half of maximum criterion, the width of the

window has been fixed to 11 data points. The number of PCA components in each model has been fixed to 2 (99% of variance) based on a global PCA performed on the clean dataset.

3.2. Software

For all the computations, chemometric analyses and graphics, Matlab v2007b programs (The Mathworks, Inc., Natick, MA, USA) were used. PCA models were derived using the svd algorithm [31] included in the PLS Toolbox (Eigenvector Research, Inc., Manson, WA, USA).

4. Results

Fig. 2 shows the first derivative spectra, the GH criterion and the Local MWPCA results, respectively, for each locally selected sample in the calibration set (in grey) and the spectrum of the sample to be characterized (in black). For the first sample, contaminated with 10,000 ppm of melamine, the three criteria allow abnormal behavior to be detected. For samples with contamination lower than 1,000 ppm, the results show that no clear conclusions can be drawn when looking directly at the spectra. Applying the GH criterion allows abnormalities at levels higher than 500 ppm to be detected. The Local MWPCA procedure enables contamination at levels up to 100 ppm to be detected, but at that level these results show that the detection of melamine in milk becomes unstable, indicating that the technique has probably reached its limit of detection. This contradicts findings reported in previous publications, which suggested that NIR spectroscopy could be used to detect melamine at levels lower than 1 ppm [32,33]. Norris [34] wrote that "those people responsible for developing calibrations for constituents at ppm levels must demonstrate the impact of the possible noise sources on their results before suggesting possible limits of detection", which was not done in the publications previously cited. In this study, in line with the recommendation made by Norris [34], three spectral replicates for each sample were collected. Then, for each clean sample, the differences between replicates were calculated two by two and the noise level was estimated as the average standard deviation of all differences. As example, Fig. 3 shows the three repetitions for a clean sample, when, different patterns were clearly visible. When applying this to all the available samples, the noise level (differences between replicates) reached an average estimated standard deviation of 200 μlog .

In Fig. 4, melamine contamination from 0.01% (100 ppm) to 1% (10,000 ppm) produced an optical density (OD) (or difference at the main peak (1558 cm^{-1})) of 0.15 Absorbance Units (UA). As indicated in Table 1, a concentration of 0.005% (50 ppm) did not produce

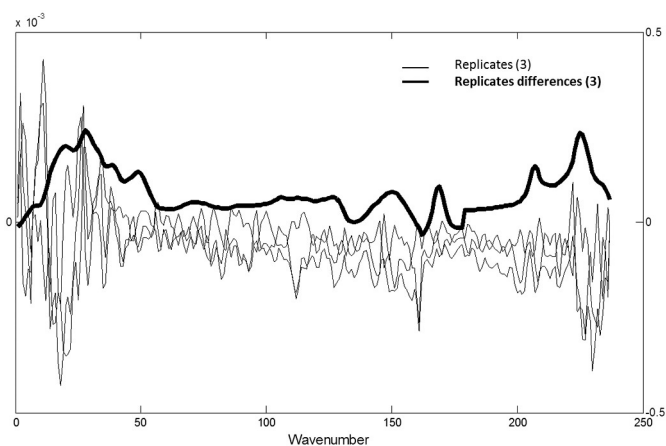


Fig. 3. Three repetitions of the same clean sample (in bold) and the differences two by two (dotted lines).

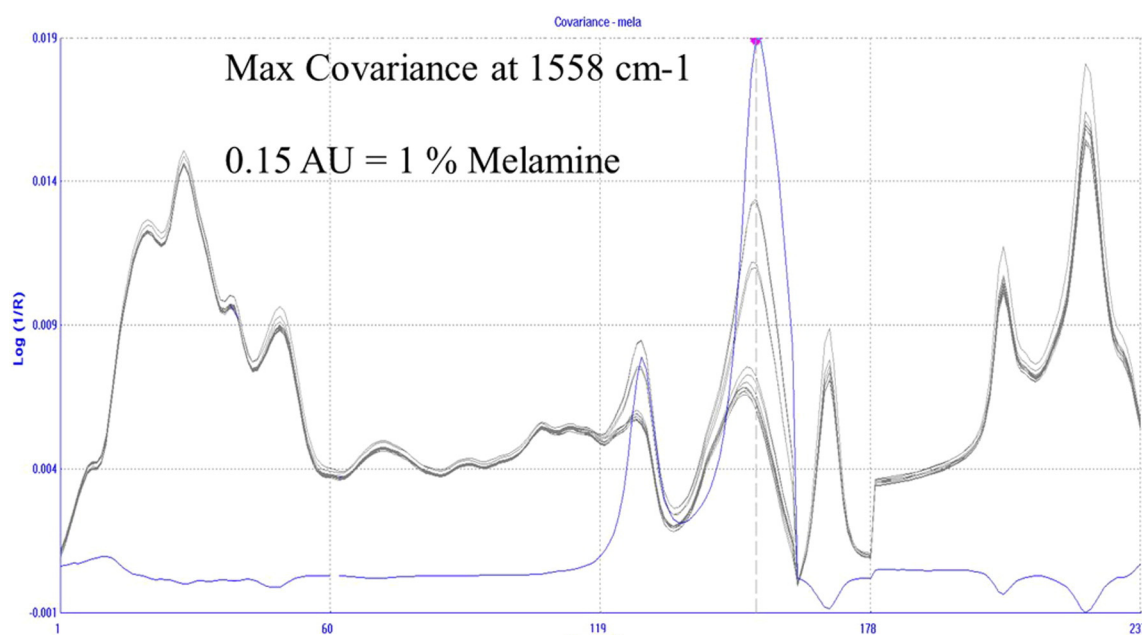


Fig. 4. Maximum covariance at 1558 cm^{-1} .

differences in the spectra larger than 0.000750 μlog . These values were completely masked in the noise.

For the sake of comparison, the outcomes of the proposed study have been compared to those of the use of combined T^2 and Q statistics after a global PCA on the whole set of clean samples. These statistics were able to easily detect contamination at levels up to 0.1% (1000 ppm) and in a less extension at levels of 0.05% (500 ppm), however they were unable to detect contamination at lower values.

5. Conclusion

In this study, a Local MWPCA method was used for the characterization of an important agronomical product and the detection of possible contaminants using vibrational spectroscopy. This application illustrated the potential of this method for detecting abnormal spectra in samples and its capacity to reduce the level of detection obtained with classical techniques. The study was based on using liquid UHT samples contaminated with melamine, making it a targeted study, but the method could be used for detecting abnormalities (inadvertent or deliberate contamination) in the data and as an initial step prior to further analyses. Local MWPCA could also be used in routine lab as a control chart to detect the deviation in time of the instruments, comparing daily a sample from a common batch of milk to previous samples from the same batch. It could be used easily in a reception/production plant for rapid and online quality control. In addition, with the local selection of the most spectroscopically similar samples, the spectral library could be built using different products, which could lead to the development of a unique global model.

Table 1
Optical Density (in Absorbance Units – UA) and S/N ratio for each level of melamine contamination (in percentage and ppm).

% Melamine	ppm	Optical Density (OD)	(SD noise = 0.000200) S/N
1	10,000	0.150000	750
0.1	1,000	0.015000	75
0.01	100	0.001500	7.5
0.005	50	0.000750	3.75
0.001	10	0.000150	0.75
0.0001	1	0.000015	0.075

Conflict of interest

The authors declare that they have no conflict of interest.

Acknowledgments

This work was conducted within the framework of the EU Seventh Framework Programme for research, technological development and demonstration under Grant Agreement No. 613688 FOODINTEGRITY project. The authors are grateful to the technical support of the Valorisation of Products department of the CRA-W for sample preparation and data acquisition.

References

- [1] A. Savitzky, M.J.E. Golay, Smoothing and differentiation of data by simplified least squares procedures, *Anal. Chem.* 36 (8) (1964) 1627–1639.
- [2] H.R. Keller, D.L. Massart, Evolving factor analysis, *Chemom. Intell. Lab. Syst.* 12 (3) (1991) 209–224.
- [3] R. Tauler, Interpretation of environmental data using chemometrics, in: D. Barceló (Ed.), *Sample Handling and Trace Analysis of Pollutants – Techniques, Applications and Quality Assurance*, Elsevier, Amsterdam, The Netherlands, 2000.
- [4] D.S. Broomhead, G.P. King, Extracting qualitative dynamics from experimental data, *Physica D20* (1986) 217–236.
- [5] N. Golyandina, V. Nekrutkin, A. Zhigljavsky, *Analysis of Time Series Structure: SSA and Related Techniques*, Chapman & Hall/CRC, 2001.
- [6] S. Xie, S. Krishnan, Dynamic principal component analysis with nonoverlapping moving window and its applications to epileptic EEG classification, *Sci. World J.* 2014 (2014).
- [7] H. Hassani, Singular spectrum analysis: methodology and comparison, *J. Data Sci.* 5 (2007) 239–257.
- [8] S. Šašić, Y. Katsumoto, H. Sato, Y. Ozaki, Applications of moving window two-dimensional correlation spectroscopy to analysis of phase transitions and spectra classification, *Anal. Chem.* 75 (16) (2003) 4010–4018.
- [9] H. Shinzawa, S. Morita, I. Noda, Y. Ozaki, Effect of the window size in moving-window two-dimensional correlation analysis, *J. Mol. Struct.* 799 (1–3) (2006) 28–33.
- [10] X.L. Chu, Y.P. Xu, S.B. Tian, J. Wang, W.Z. Lu, Rapid identification and assay of crude oils based on moving-window correlation coefficient and near infrared spectral library, *Chemom. Intell. Lab. Syst.* 107 (2011) 44–49.
- [11] X. Liu, U. Kruger, T. Littler, L. Xie, S. Wang, Moving window kernel PCA for adaptive monitoring of nonlinear processes, *Chemom. Intell. Lab. Syst.* 96 (2) (2009) 132–143.
- [12] J.C. Jeng, Adaptive process monitoring using efficient recursive PCA and moving window PCA algorithms, *J. Taiwan Inst. Chem. Eng.* 41 (4) (2010) 475–481.
- [13] S.R. Ryu, I. Noda, Y.M. Jung, Moving window principal component analysis for detecting positional fluctuation of spectral changes, *Bull. Kor. Chem. Soc.* 32 (7) (2011) 2332–2338.

- [14] J.H. Jiang, R.J. Berry, H.W. Siesler, Y. Ozaki, Wavelength interval selection in multi-component spectral analysis by moving window partial least-squares regression with applications to mid-infrared and near-infrared spectroscopic data, *Anal. Chem.* 74 (2002) 3555–3565.
- [15] B. Lu, I. Castillo, L. Chiang, T.E. Edgar, Industrial PLS model variable selection using window variable importance in projection, *Chemom. Intell. Lab. Syst.* 135 (2014) 90–109.
- [16] J.A. McLean, Targeting the untargeted, *Anal. Sci.* 17 (2014) 17–18.
- [17] V. Baeten, P. Vermeulen, J.A. Fernández Pierna, P. Dardenne, From targeted to untargeted detection of contaminants and foreign bodies in food and feed using NIR spectroscopy, *New Food* 17 (3) (2014) 16–23.
- [18] X. Lu, Y. Si-Min, C. Chen-Bo, W. Zhen-Ji, Y. Xiao-Ping, The feasibility of using near-infrared spectroscopy and chemometrics for untargeted detection of protein adulteration in yogurt: removing unwanted variations in pure yogurt, *J. Anal. Methods Chem.* 2013 (2013) 1–9.
- [19] J.C. Moore, A. Ganguly, J. Smeller, L. Botros, M. Mossoba, M.M. Bergana, Standardisation of non-targeted screening tools to detect adulterations in skim milk powder using NIR spectroscopy and chemometrics, *NIR News* 23 (5) (2012) 9–11.
- [20] L. Xu, S.M. Yan, C.B. Cai, X.P. Yu, Untargeted detection and quantitative analysis of poplar balata (PB) in Chinese propolis by FT-NIR spectroscopy and chemometrics, *Food Chem.* 141 (4) (2013) 4132–4137.
- [21] FOSS, Abnormal spectrum screening (ASM), Dedicated Analytical SolutionsA White Paper from FOSS. P/N 1026513, Issue 2, May 20142014.
- [22] O. Abbas, B. Lecler, P. Dardenne, V. Baeten, Detection of melamine and cyanuric acid in feed ingredients by near infrared spectroscopy and chemometrics, *J. Near Infrared Spectrosc.* 21 (2013) 183.
- [23] S.A. Haughey, S.F. Graham, E. Cancouët, C.T. Elliott, The application of near-infrared reflectance spectroscopy (NIRS) to detect melamine adulteration of soyabean meal, *Food Chem.* 136 (2013) 1557.
- [24] J.A. Fernández Pierna, D. Vincke, P. Dardenne, Z. Yang, L. Han, V. Baeten, Line scan hyperspectral imaging spectroscopy for the early detection of melamine and cyanuric acid in feed, *J. Near Infrared Spectrosc.* 22 (2014) 103–112.
- [25] M. Lin, A review of traditional and novel detection techniques for melamine and its analogues in foods and animal feed, *Front. Chem. Eng. Chin.* 3 (2009) 427.
- [26] L. Yuan, E. Ewen, D. Todd, Z. Qiang, S. Jiang-rong, L. Xian-jin, Recent developments in the detection of melamine, *J. Zhejiang Univ. Sci. B* 13 (7) (2012) 525–532.
- [27] T. Davies, *Spectroscopy Europe* 11/4, 1999.
- [28] W.H. Press, S.A. Teukolsky, W.T. Vetterling, B.P. Flannery, *Numerical Recipes in C: The Art of Scientific Computing*, Cambridge University Press, 1992.
- [29] N. Ruiz, 'Near Infrared Spectroscopy: present and future applications', Technical Bulletin of the American Soybean Association, Vol FT52-2001. Available at http://ussec.org/wp-content/uploads/sites/6/2013/06/sm_052010_E.pdf 2001 (last accessed October 2015).
- [30] R.G. Whitfield, M.E. Gerger, R.L. Sharp, Near-infrared spectrum qualification via Mahalanobis distance determination, *Appl. Spectrosc.* 41 (7) (1987) 1204–1213.
- [31] J.E. Jackson, *A User's Guide to Principal Components*, Wiley, 1991.
- [32] C. Lu, B. Xiang, G. Hao, J. Xu, Z. Wang, C. Chen, Rapid detection of melamine in milk powder by near infrared spectroscopy, *J. Near Infrared Spectrosc.* 17 (2009) 59–67.
- [33] L.J. Mauer, A.A. Chernyshova, A. Hiatt, A. Deering, R. Davis, Melamine detection in infant formula powder using near- and mid-infrared spectroscopy, *J. Agric. Food Chem.* 57 (10) (2009) 3974–3980.
- [34] K. Norris, Letter to the Editor: hazards with near-infrared spectroscopy in detecting contamination, *J. Near Infrared Spectrosc.* 17 (4) (2009) 165–166.