

APPLICATION OF CHEMOMETRICS

Bernard G.M. Vandeginste

Department of Analytical Chemistry, University of Nijmegen, 6525 ED Nijmegen,
The Netherlands

Abstract - Chemical analysis on a micro scale and at trace concentrations demand the utmost of today's technical achievements. An essential condition hereby is to design and select optimal measurement systems. Likewise, it is necessary to extract maximum information from the analytical data. Chemometrics uses and develops mathematical and statistical methods in order to meet the requirements for the solution of today's analytical problems. An overview will be given, first of the chemometric "tools", whereafter three examples will illustrate the potentials of the methods for (i) the optimization of an analytical procedure, (ii) the obtaining of maximum information from the analytical data and (iii) for the combination and classification of analytical results.

INTRODUCTION

In spite of the incredible diversity of analytical methods and procedures it is not difficult to discover a certain similarity in the way analytical problems are solved. The first action is to translate the analytical problem into its analytical and economical terms (Fig. 1).

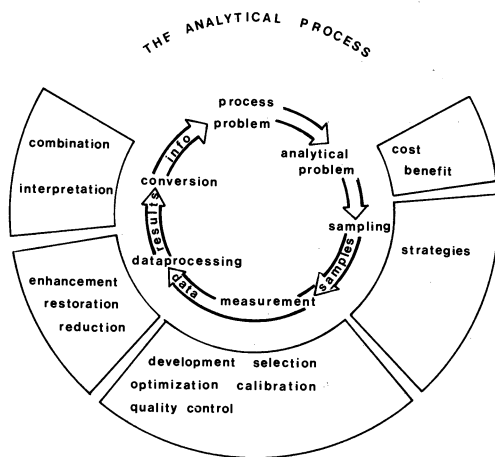


Fig. 1. The analytical process.

The latter may sound a bit surprising, but one should realize that anyhow it is senseless to produce analytical information which is less worth than the analysis costs. The analysis costs are determined by the sampling programme in conjunction with the quality of the analytical information. The analytical procedure is the tool for the production of the analytical information.

Many operational conditions of the analytical procedure determine the quality of the produced data. Under bad conditions the method will probably give inaccurate or imprecise results, or it may take too long before obtaining the result. Classically, the best or optimal conditions are determined by varying one parameter at a time. However, such univariate methods require many experiments and often find a wrong optimum (Fig. 2a). Instead, a multivariate approach is necessary (Fig. 2b), which is discussed below.

With the ongoing revolutionary developments in micro electronics the era of 'self optimizing' analytical instruments, controlled by built-in optimization algorithms is at the doorstep. As calculation capability is growing sheaper one has to consider the costs of extracting more

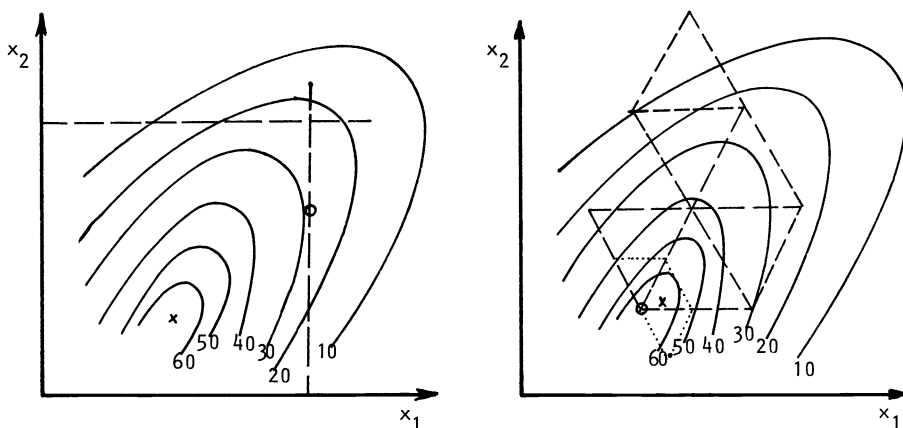


Fig. 2a. Univariate optimization

Fig. 2b. Simplex optimization

information from the data against the production of more data by carrying out additional measurements (often requiring an additional and thus a more expensive sampling). For instance, the experimental design of the measurements influences the amount of analytical information which is obtained on the analyzed system. A good design may yield more information from even a smaller number of measurements in comparison to a bad design. In the field of calibration for instance, parameters of the calibration function may be calculated in a recursive way: i.e. after each measurement the parameters of the calibration function are re-estimated and their values are fed back to the experimental design.

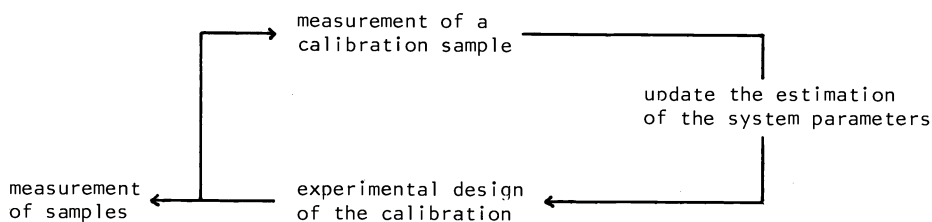


Fig. 3. Recursive on-line calibration

After the measurement process, the enhancement and restoration of the unrefined data can be started (e.g. a smoothing or deconvolution), which is followed by a proper data analysis. Good data processing methods are essential in order to utilize the full capabilities of analytical methods. In HPLC for instance, the potentials of a diode array detector are only fully exploited when the proper multivariate methods are applied like a principal components analysis followed by a curve resolution step. This is illustrated on an example given below. The product of an analytical procedure is a set of analytical results (e.g. concentrations) which are obtained for a number of samples. Such a set of results has to be combined into chemical information (or knowledge) on the sampled system (the system under investigation). For instance, the results of a water analysis of a drinking water production plant have to be combined into one or more quality parameters, which allow proper control actions. In this respect, the chemometrician plays an important role. Data structures are often so complex, that advanced mathematical and statistical methods are necessary to recover the hidden information. An illustrative example is the detection of a periodicity in data (or time) series, by a relatively simple autocorrelation analysis (Fig. 4). An analysis of variance is another example.

Many analytical problems are put in terms whether an object on which several parameters have been measured, belongs to one or another category. An example is the classification (or identification) of minerals on the basis of their X-ray spectra. Classification methods are unified under the term Pattern Recognition. Many applications of pattern recognition in analytical chemistry have been reported. A major operation in pattern recognition is to provide displays for the analyst in a reduced 2-D space, with the conservation of as much as possible of the original structure in the data (e.g. clusters). Developments are ongoing now to incorporate pattern recognizers in the analytical instruments for an on-line interpretation and/or classification of the results. At our laboratory, the possibilities are investigated to analyze surface images obtained with a raster-electron-microscope (R.E.M.), by digital image

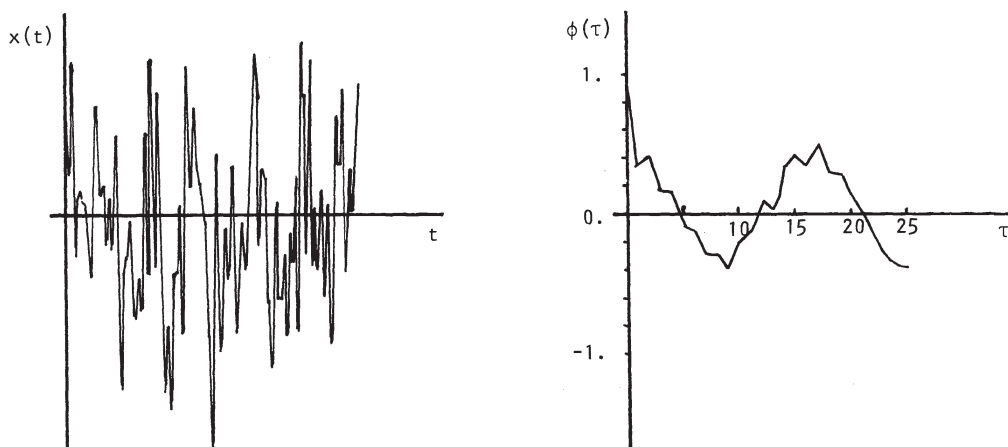


Fig. 4. Time series and its autocorrelation function.

processing. Two-dimensional data matrices ($n \times m$) can be displayed under the form of an image of $n \times m$ pixels, having a given grey level (or colour). Because the image is available in a digitized form, many operations can be carried out on the image: e.g. an image enhancement or an edge detection (to be compared with a peak detection in spectra).

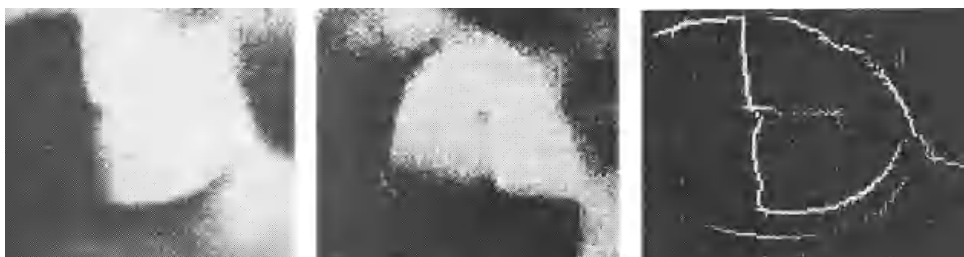


Fig. 5a. SIMS Fe image ($128 \times 128 \times 8$) of a 5×5 mm Fe blade.

Fig. 5b. Image after a high emphasis spatial filtering.

Fig. 5c. Edge detection.

After the image processor has located the coordinates of interest (e.g. particles of a given size when analyzing minerals in coal), a micro processor controlled sample table is moved to that location, and an energy-dispersed X-ray (EDAX) spectrum is taken, whereafter the pattern recognizer gives the mineral constitution.

Much work has still to be done in that respect, but it illustrates the track, chemometricians are moving on: that is, to turn instruments into intelligent analyzers.

THREE APPLICATIONS

The on-line optimization and calibration of an analytical procedure

Analytical chemists are right to be concerned whether analytical procedures operated on a routine basis are optimally performing within the calibrated conditions. For the finding of the optimal conditions one uses an optimization procedure and defines a criterion to be optimized. An optimization method ought to be fast (small number of measurements) and to lead to the right optimum. Figure 2a shows that the classical method of optimizing one parameter at a time may fail in finding the optimum. The reason is that the optimum for one parameter may depend on the level of all other parameters. Multi variate sequential methods (Ref. 1) vary all parameters at a time in a search of the steepest (and thus shortest) route to the optimum (Fig. 2b). The measurements are arranged in geometrical figures (Simplices which are triangles in the 2-d case). The coordinates of the corners of the figure

represent a set of values of the parameters. By dropping the corner with the worst response and by moving into the opposite direction, a new Simplex is constructed. The coordinates of the new corner represent the values of the parameters for the next measurement. In the neighbourhood of the optimum the Simplices start to turn about the optimum, but by taking some special precautions the optimum can be further approached.

In practice, modified Simplex procedures are used which perform better than the above mentioned symmetrical one (Ref. 2).

Implicitly it has been assumed that optimization was wanted for one criterion only. For example, maximal sensitivity, or minimal analysis time, or maximal peak resolution per unit of time etc. Although many analytical problems are of a multicriteria nature, optimization in analytical chemistry has been mainly limited to unicriteria methods. Because of its sequential nature, the Simplex algorithm is very efficient for on-line applications in so-called self optimizing instruments. The principle is relatively simple. The analytical instrument is hooked to a microprocessor loaded with a Simplex algorithm. The coordinates of the corners of the starting Simplex are entered by the analyst. From that point the self-optimization proceeds until the optimal settings of the instrumental parameters are found, without any further human interaction. The microcomputer has a triple task in that respect: the control of the Simplex, the calculation of the response by a processing of the unrefined data (e.g. peak height) and the transfer of the new values of the instrument settings to the instrument. As far as I know computer controlled self optimizing instrumentation has been realized in flow injection analysis and furnace atomic absorption spectrometry. In FIA (Ref. 3) the instrument sets the optimal flows of the reagents and carrier stream, in order to obtain maximal sensitivity. In AAS (Ref. 4) the instrument optimizes the temperatures and times of the drying, ashing, atomizing and burning stages, using standards. However, in AAS, optimum values found for a standard may be off-optimum for an unknown sample. Hence, more specific intelligence has to be incorporated in advance in the instrument or should be acquired by the instrument itself. A form of artificial intelligence may also be realized in calibration. A calibration function which takes into account interferences and matrix effects is given by: $\underline{R} = \underline{K} \cdot \underline{C}$, where \underline{R} is the response vector of the responses measured at NS sensors, \underline{C} is the concentration vector of the NA analytes and \underline{K} is a NS x NA matrix of the sensitivities of the NA analytes at NS sensors. In the case of matrix effects every element k_{ij} of \underline{K} is a function of all concentrations $c_{i \neq j}$.

Saxberg and Kowalski (5) developed the Generalized Standard Addition Method (GSAM) for the calibration of such systems. All elements of \underline{K} are determined by adding standard solutions of all analytes to the sample, according to a given experimental design. After each addition the responses are measured at all sensors. After completion of the additions \underline{K} and \underline{C}_0 (concentrations of the NA analytes in the unknown sample) are calculated. During the additions, however, the measurements contain already information on the system state (the \underline{K} matrix). In classical calibration one measures first all standard solutions, without using that implicit information for an eventual adaptation of the calibration design. Recently (Ref. 6), a recursive calibration method has been developed whereafter each measurement of a standard, calibration constants are updated (Fig. 6a).

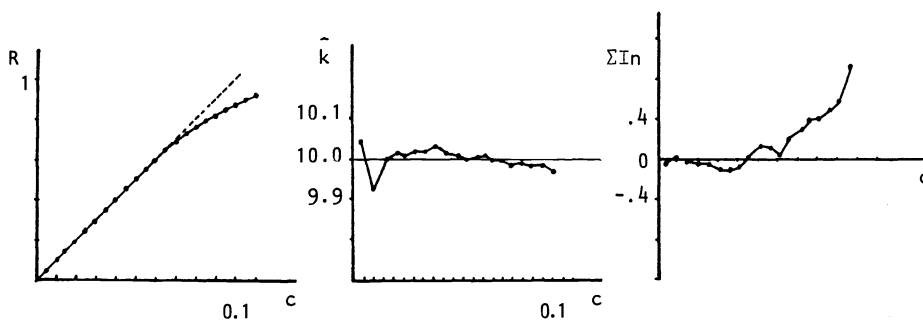


Fig. 6a Calibration curve with non-linearity.

Fig. 6b Estimated k values.

Fig. 6c Cumulative innovation for detection of non-linearity.

The general expression of the recursive algorithm is (Ref. 6)

$$\text{New estimate of } K = \text{old estimate of } K + \text{correction} \quad (1)$$

The correction term is a function of the value of the last measurement.

The application of a recursive algorithm allows to calibrate according to the scheme given in Fig. 3. Important features of this design are:

(i) an on-line control of the validity of the underlying model by monitoring the innovation which is the difference between the measured response, R_n and the expected response \hat{R}_n , based on the old estimate \hat{K}_n . When the model deviates from linearity, the innovations will keep equal signs, instead of fluctuating about zero (Fig. 6c);

(ii) the availability of real time information whether the estimates of \hat{K} are known within the desired limits.

In my opinion, future instrumentation will become the more and more intelligent. First steps in that direction are the realisation of self-optimizing and self-calibrating instruments.

The acquisition of maximal information from analytical data

The upcoming of hyphenated methods, like GC-MS and recently, HPLC-UV/VIS, has enormously augmented the potentials of chemical analysis. Gas chromatography, for instance, is a poor qualitative method but is a powerful quantitative method for complex mixtures. On the other hand mass spectrometry is a powerful qualitative method, when dealing with non-mixtures. The combination of both methods unifies outstanding qualitative and quantitative capabilities. As long as all analytes of interest are well separated, all desired analytical information is easily acquired. In many instances, however, good resolution may be only achieved after a time consuming optimisation, with a risk that even then the analytes remain poorly resolved. In that situation, the limits of traditional qualitative and quantitative analysis have been reached. However, by a multivariate approach, which is explained below, the number of components in a poorly resolved elution profile can be estimated. More important, however, is the fact that the pure spectra of the analytes can be estimated also, provided that the number of analytes in the profile is less than four.

In this paper attention will be focussed on the HPLC/UV-VIS method. Contrary to the single or dual wavelength HPLC, where one or two chromatograms are recorded, HPLC equipped with a diode array detector may produce full spectra (e.g. 20 or more wavelengths) every second. The resulting analytical data, obtained for a strongly overlapped system of 3 analytes (Fig. 7a) consists of a data matrix of, for instance, NS spectra recorded over NW wavelengths (Fig. 7b).

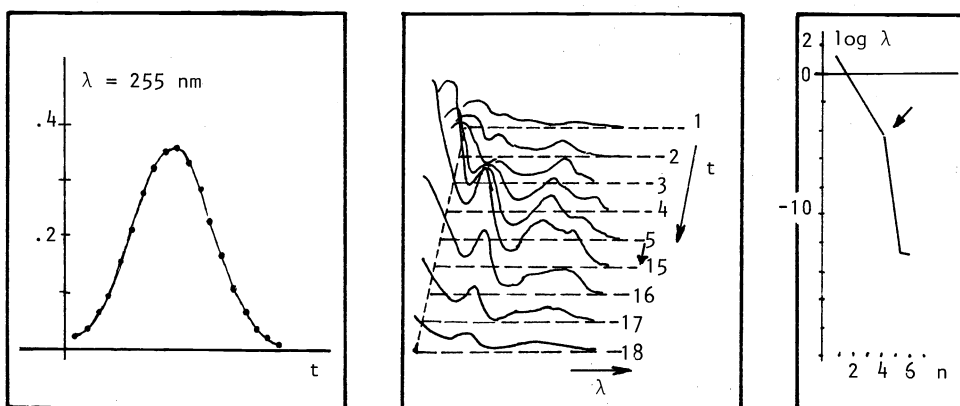


Fig. 7a 3-compound elution profile (Sampling interval: 1 sec.).

Fig. 7b Spectra recorded over elution profile.

Fig. 7c Determination of the number of components.

The absorbances at each wavelength during the elution are linear combinations of the absorbances of every component present in the elution profile. This property allows to find the number of independent factors (or compounds). This can be illustrated on a simplified example, where spectra are recorded at 3 wavelengths only ($NW = 3$). Every spectrum (i) consists then of three absorbances (a_{i1}, a_{i2}, a_{i3}), which can be represented as a point in a three dimensional space (Fig. 8a). Let us consider two cases:

(i) All spectra ($i = 1, NS$) are combinations of one component only. Then all points (spectra) will lie on a straight line. Because of noise superposed on the data some random deviations from that line will be observed in practice. Looking to the same data but in a little different way, one can say that the line is pointed in the direction of maximal variation (or variance) in the data (Fig. 8b).

(ii) All spectra ($i = 1, NS$) are combinations of spectra of two compounds. It can be easily checked that all points (spectra) lie now about a plane surface. In terms of variations, the plane defines two orthogonal directions of variation (Fig. 8c).

In principle, when noise is present on the data, one will find three orthogonal directions

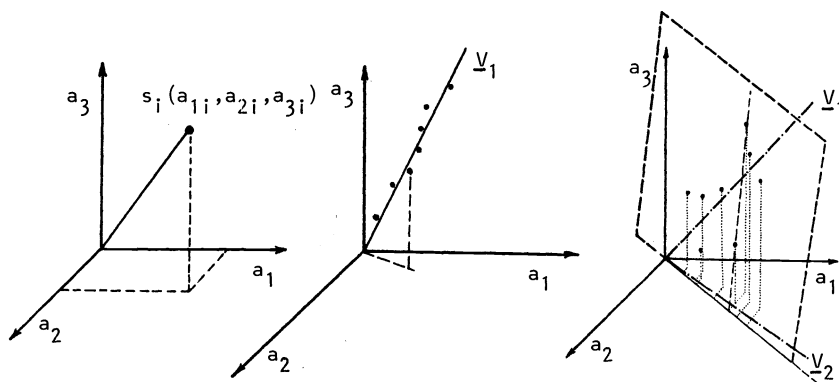


Fig. 8a Representation of a spectrum in the NW-dimensional space (NW = 3).

Fig. 8b One-component spectra in the NW-dimensional space (NW = 3).

Fig. 8c Two-component spectra in the NW-dimensional space (NW = 3).

of variation in a three dimensional space. These directions are called the eigenvectors of the variance-covariance matrix of the datamatrix. The eigenvalue of this matrix is a measure for the variance spanned by its respective eigenvector. From the Fig. 8b and 8c it is readily seen that the number of significant directions of variance (the other directions describe only the noise which is present in the data) is also the number of factors (here compounds) in the data. The mathematical method is called principal components analysis (Ref. 7). Figure 7c gives a plot of the eigenvalues obtained for the unresolved elution profile of three diphenyl amines (Fig. 7a). The figure indicates that there are 3 components with a small fourth one (is a slight overlap of a peak at the left wing).

The next problem is to calculate the spectra of the pure compounds and to reconstruct the elution profiles for the determination of the sequence of the retention times.

Lawton and Sylvestre (Ref. 8) gave a solution for a two component system. Chen (Ref. 9) developed an algorithm for a three component system of GC-MS data.

At our laboratory the algorithm of Chen has been adapted to handle HPLC/UV-VIS data up to three unresolved components.

If three compounds are present, three significant eigenvalues are found, associated with three eigenvectors v_1 , v_2 and v_3 . The coordinates of the spectra, first represented in a NW-th dimensional space, (NW is the number of wavelengths) can be calculated in the v_1, v_2, v_3 space

$$\underline{M}_i = a_{i1}v_1 + a_{i2}v_2 + a_{i3}v_3 \quad (2)$$

For an easy 2-d representation, one can use the orthogonal coordinates (θ, φ) , where

$$\underline{M}_i = \cos \theta v_1 + \cos \varphi \sin \theta v_2 + \sin \theta \sin \varphi v_3 \quad (3)$$

All spectra (after being normalized on a norm = 1) can be represented as a vector (angle φ and length θ) in a θ - φ plot. The next step is to extrapolate the measured spectra in order to find estimations of the pure spectra S_1, S_2 and S_3 . For $\varphi = 0$ to 2π , the value of θ is calculated for which one of the elements of \underline{M}_i becomes zero (and all other elements are positive). This gives a plot as shown in Fig. 9.

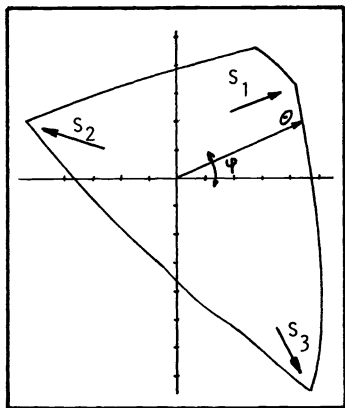


Fig. 9 θ - φ plot of candidate pure spectra

An algorithm, which will be published elsewhere, has been developed at our laboratory in order to select the three pure spectra from such a Θ - \varnothing plot. Simulations carried out at our laboratory (Ref. 10) revealed that good qualitative and quantitative (after proper calibration) information is obtainable also under conditions of poor chromatographical resolution of compounds with very similar spectra. The strength of the method is well demonstrated when comparing the estimated pure spectra and the real pure spectra (Fig. 10a, 10b, 10c) of the three compound system being eluted with the elution profile of Fig. 7a. The spectra are perfectly recognizable allowing the identification of the compounds.

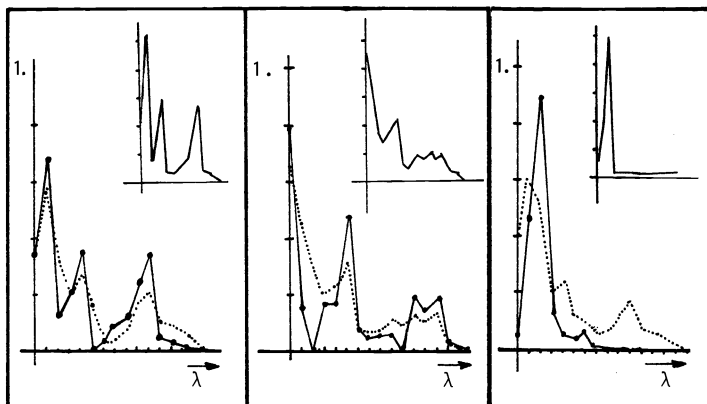


Fig. 10a, b, c Estimated pure spectra (—); purest spectra (.....) in the mixture (spectra in the corner are the true spectra).

The elution profiles are calculated by solving the overdetermined set of equations at every elution time t_i :

$$\underline{M}_i = \sum_{j=1}^3 c_{i,j} \underline{S}_j$$

In matrix notation: $\underline{M}_i = \underline{S} \cdot \underline{c}_i$, where \underline{M}_i is the spectrum recorded at time t_i , \underline{c}_i is the vector of the three unknown profiles at time t_i , \underline{S} is a 3 column matrix of the three estimated pure spectra. The solution at t_i is given as: $\underline{c}_i = [\underline{S}^T \cdot \underline{S}]^{-1} \cdot \underline{S}^T \cdot \underline{M}_i$.

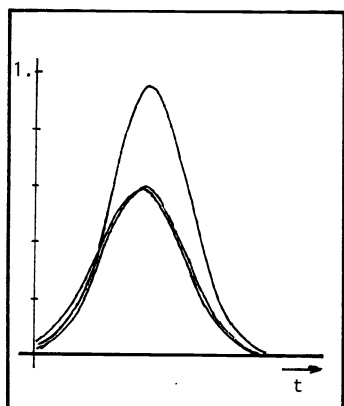


Fig. 11 Calculated elution profiles of the compounds.

Confidence intervals of the elution profile are found by repeating the calculations with the purest measured spectra (dots in Fig. 10a, b, c).

A comparison of the retention times found by the program, with the true retention times of the model system, and a comparison of the estimated pure spectra with the true spectra illustrate the promising capabilities of the method in view of the poor chromatographical resolution.

The combination and classification of analytical results

The formulation of the analytical information in many cases requires the combination of many analytical results. A direct assignment of a property to an object is often very difficult. An empirical approach instead is to measure certain variables or features on objects with a known property. Thereafter one tries to derive a classification rule which attributes the right property to certain regions in the multi dimensional space, spanned by the features. The classification rule is thereafter used to classify the unknowns. The basic principle of such pattern recognition method is that the features form the axes of a multi dimensional feature space. Feature values of an object locate that object in that space. The basic assumption is that the closer objects are located, the more similar they are. When the common property of the objects is unknown, one tries to locate clusters in the feature space. Analysis of the clusters may lead to the discovering of the common property. Pattern recognition consists of a collection of algorithms for many types of operations: display of the multi dimensional space; transformation of the axes (features) in order to enhance the separation of the categories; modelling of the clusters, etc.

Pattern recognition has been successfully applied on a wide variety of analytical problems for more than 10 years (Ref. 11, 12). It has been and still is one of the main subjects of research of chemometricians. Patterns are somehow associated with images. An image consists of a raster of points, called pixels, having a given grey level. Every 2-dimensional matrix can be represented as an image. Chemical images have been recently obtained in Secondary Ion Mass Spectrometry (SIMS) (Ref. 13). A narrow ion beam is moved over the sample surface, causing the reflection of secondary ions. The mass of these ions is measured with a quadrupole mass spectrometer. When the mass spectrometer is tuned on a fixed mass unit, a digitized matrix is obtained of intensities as a function of the spatial coordinates. Such a matrix can be displayed as an image, whereafter various image processing methods are applicable.

It is clear that the combination of image processing and pattern recognition may become a powerful tool for the surface analyst. The image processing (IP) is the 2-dimensional equivalent of the data processing of spectra, chromatograms etc. Peak-find procedures in spectra become edge detectors in IP, deconvolution becomes image restoration and the 1-d smoothing procedures are image enhancement techniques.

At our laboratory a system is under development for the automatic analysis of images of a Raster Electron Microscope (R.E.M.) in order to find the spots of interest for a local X-Ray analysis (RMA). Ideally, the analysis should be concluded by an automatical interpretation of the X-Ray spectra. The aim is thus to carry out a surface analysis according to the scheme shown in Fig. 12.

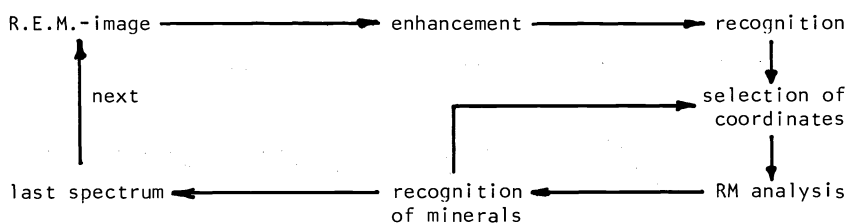


Fig. 12 Scheme of an intelligent R.E.M./R.M.A. instrument

The X-Ray spectra are recorded at a number of selected windows of energy. In our particular case where the mineral constitution in coal is determined, these windows are for K, Ca, S, Fe, Si, and Al. In the terminology of multivariate statistics, every spectrum is a point in the 6-th dimensional feature space. In order to investigate the possibilities of an automatical analysis of the spectra, spectra were recorded at 400 spots arranged in a regular 20 x 20 grid over the sample surface. 160 of them were essentially null spectra (no mineral present at that spot). In order to find the structure in the data, one needs a plot of all spectra in one or more 2-dimensional plots. Referring to the previous paragraph, a projection of all points on a plane in the direction of the two principal orthogonal axes of variation (eigenvectors v_1, v_2 of the two largest eigenvalues) will conserve the most of the variance in the data. Figure 13 shows a plot of the data projected on the (v_1, v_2) plane after normalization. It is obvious that some structure becomes visible in the data.

Spectra are lined in triangles, to be compared with diagrams obtained for ternary mixtures. The corners are the simplest spectra and lines connecting the corners represent spectra which are linear combinations of the spectra in the corners. The fact that spectra are nicely lined up indicates that most spectra are pure spectra or are spectra of mixtures of 2 to 3 minerals only.

An obvious continuation is to categorize the spectra according to the windows, which give a signal. In this way a total of 17 categories were obtained: K; S; Si; K-Ca; Fe-S; Si-Al; Ca-S; Ca-Si; S-Si, Ca-Si-S; K-Ca-S; S-Si-Fe; Si-S-Al; K-S-Si-Al; Ca-S-Si-Fe; S-Si-Fe-Al; K-Si-S-Fe-Al; K-S-Si-Al. By performing a disjoint principal components analysis on each of the categories, information is obtained about the number of components present in each group (Table 1) of spectra and about the location, size and direction of the 17 categories of spectra in

the 6-d space. This information in fact is a model for each group, known as SIMCA (Ref. 14).

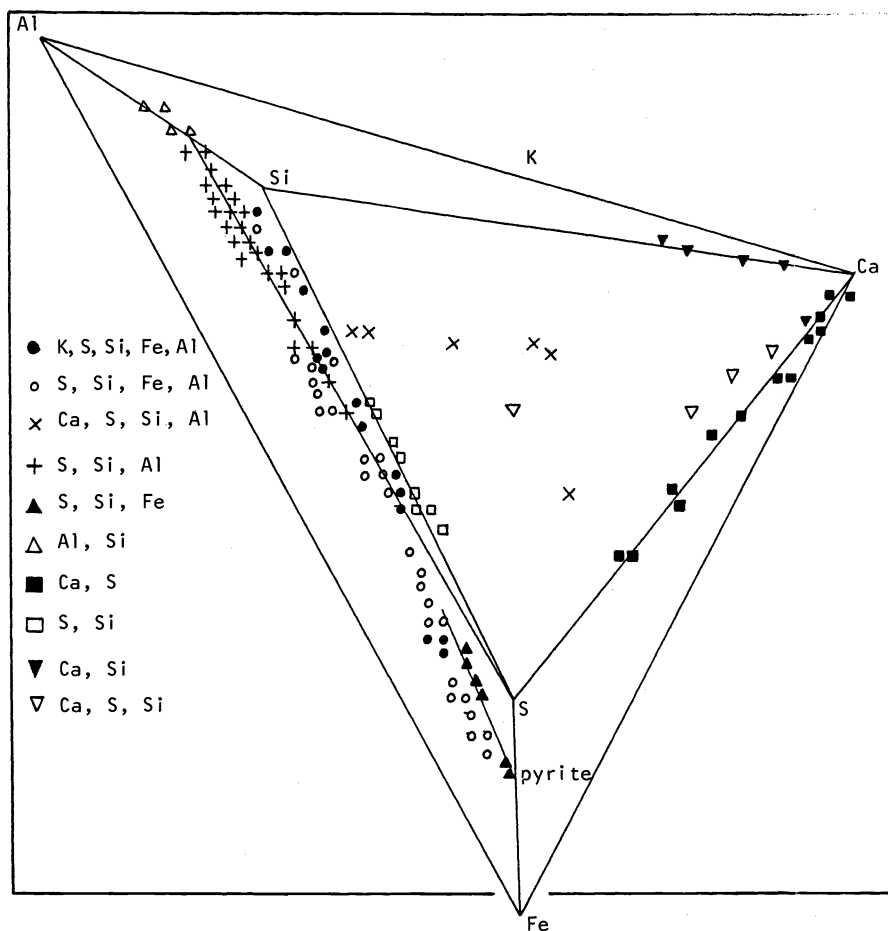


Fig. 13 Karhunen-Loeve projection (v_1, v_2) of X-ray spectra of minerals in coal (spectra are categorized according to the combination of spectral windows with a signal).

TABLE 1. Results of a disjoint principal components analysis on 17 categories of line combinations

category	information preserved with						vectors
	1	2	3	4	5	6	
K, Ca	93.6	100					
Ca	100						
Ca, S	93.2	100					
Ca, Si	99.9	100					
Si	100						
S, Si, Fe	99.5	99.9	100				
S, Si, Fe, Al	70.6	97.3	99.8	100			
K, Ca, S	100						
S, Fe	100						
Ca, S, Si	82.5	98.3	100				
Ca, S, Si, Fe	68.2	100					
S, Si, Al	99.0	100					
S, Si	100.0						
K, S, Si, Fe, Al	77.4	91.4	96.8	99.9	100		
Si, Al	98.5	100					
Ca, S, Si, Al	97.1	99.3	100.0				
K, S, Si, Al	92.3	100					

The SIMCA model is capable to classify unknown spectra in one of the categories and is capable to recognize whether the unknown is a member of a new category. The classification, however, is in this example a trivial matter as it classifies the spectra according to the windows showing a signal. The data listed in Table 1, however, show that only a very limited number of minerals is present in the system. One should be aware, however, that the maximal number of compounds which can be detected in a category is equal to the number of windows having a signal. Thus for a group of spectra with 3 signals a maximum of 3 components can be found.

From Table 1 it follows that most 2 line spectra are mainly 1 compound. Three groups of 3 line spectra (K-Ca-Si; Si-S-Fe and S-Si-Al) are seen as one mineral (with a fairly constant contamination of S). Two four line spectra are also seen as one mineral. All other categories are mainly combinations of two minerals. At this point the study may proceed to curve resolution as explained in the previous paragraph. In a further research, spectra of the pure compounds will be estimated for each group, whereafter the composition of the mixtures can be estimated, completing the interpretation of the spectra.

Acknowledgement - The author wishes to thank B.R. Kowalski, Laboratory for Chemometrics, University of Washington, Seattle, for making available to us the computer package 'ARTHUR' for pattern recognition. L. de Galan, Technical University, Delft, The Netherlands, is gratefully acknowledged for making available the HPLC/UV-VIS data.

REFERENCES

1. D.L. Massart, A. Dijkstra and L. Kaufman, Evaluation and Optimization of Laboratory Methods and Analytical Procedures, 2nd. edn. Elsevier, Amsterdam (1980).
2. P.F.A. van der Wiel, Anal. Chim. Acta, accepted for publication (1983).
3. T.J. Sly, D. Betteridge, D. Wibberley and D.G. Porter, J. Automatic Chem., **4**, 186 (1982).
4. P.F.A. van der Wiel, L.G. van Dongen, B.G.M. Vandeginste and G. Kateman, Laboratory Microcomputer, accepted for publication (1983).
5. B.E.H. Saxberg and B.R. Kowalski, Anal. Chem., **51**, 1031 (1979).
6. B.G.M. Vandeginste, J. Klaessens and G. Kateman, Anal. Chim. Acta, **150**, (1), 71-86 (1983).
7. E.R. Malinowski and D.G. Howery, Factor Analysis in Chemistry, J. Wiley & Sons, New York (1980).
8. W.H. Lawton and E.A. Sylvestre, Technometrics, **13**, 617 (1971).
9. Jie-Hsung Chen and Lian-Pin Huang, Anal. Chim. Acta, **133**, 271-281 (1981).
10. R. Essers, B.G.M. Vandeginste and G. Kateman, to be published (1983).
11. D.L. Massart and L. Kaufman, The Interpretation of Analytical Chemical Data by the use of Cluster Analysis, J. Wiley & Sons, New York (1983).
12. B.R. Kowalski (Ed.), Chemometrics, Theory and Application, Amer. Chem. Soc. Symposium, Ser. **52** (1977).
13. B.G.M. Vandeginste and B.R. Kowalski, Anal. Chem., **55**, 557-564 (1983).
14. S. Wold, Pattern Recognition, **8**, 127-139 (1976).