



Chemometrics applied to unravel multicomponent processes and mixtures

Revisiting latest trends in multivariate resolution

A. de Juan*, R. Tauler

Chemometrics Group, Universitat de Barcelona, Diagonal 647, 08028 Barcelona, Spain

Received 31 March 2003; received in revised form 2 June 2003; accepted 4 June 2003

Abstract

Progress in the analysis of multicomponent processes and mixtures relies on the combination of sophisticated instrumental techniques and suitable data analysis tools focused on the interpretation of the multivariate responses obtained. Despite the differences in compositional variation, complexity and origin, the raw measurements recorded in a multicomponent chemical system can be very often described with a simple model consisting of the composition-weighted sum of the signals of their pure compounds.

Multivariate resolution methods have been the tools designed to unravel this pure compound information from the non-selective mixed original experimental output. The evolution of these chemometric approaches through the improvement of exploratory tools, the adaptation to work with complex data structures, the ability to introduce chemical and mathematical information in the algorithms and the better quality assessment of the results obtained is revisited. The active research on these chemometric area has allowed the successful application of these methodologies to chemical problems as complex and diverse as the interpretation of protein folding processes or the resolution of spectroscopic images.

© 2003 Elsevier B.V. All rights reserved.

Keywords: Chemometrics; Multivariate resolution; Multicomponent system; Curve resolution

1. Introduction

The characterization of a multicomponent system requires as a first step recording a relevant and sufficiently discriminating experimental output. This can be done by analysing a sample with a hyphenated technique or monitoring a process in a multivariate fashion. In these and similar examples, all the data collected form a table or data matrix where one direction (the elution or the process direction) is related

to the compositional variation of the system and the other direction refers to the variation in the response collected. The existence of these two directions of variation helps to differentiate between components.

In this context, it is important to note the wide concept of component, defined as any entity giving a distinct and real pure response. This includes examples as diverse as a chemical compound, a conformational state or a pollution source, whose response could be a profile including the relative apportionment of its different pollutants. Within the variety of multicomponent systems, processes and mixtures can be placed at the two extremes. The term process holds for reaction data, where the compositional changes respond

* Corresponding author. Tel.: +34-934034445;

fax: +34-934021233.

E-mail address: annaj@apolo.ubi.es (A. de Juan).

to a known physicochemical model, or for any evolving chemical system (e.g. a chromatographic elution), whose sequential compositional variation can be caused by physical or chemical changes and whose underlying physicochemical model, if any, is too complex or simply unknown. Mixtures would have a complete random variation along the compositional direction of the data set. An example could be a series of spectra collected in independent multicomponent samples. Other data sets lie between these two groups because they lack the global continuous compositional evolution of a process, but they can present it locally. It is the case of spectroscopic images that can have a smooth compositional variation in neighbouring pixels or of environmental data sets, where close geographical sampling points can be compositionally related.

In the analysis of any multicomponent system, the main goal is transforming the raw experimental output into useful information. We aim at passing from a plain visual sketch of the overall measured variation in our chemical data to a clear description of the contribution of each of the components present in the mixture or the process. Despite the diverse nature of multicomponent systems, the variation in their related experimental outputs can, in many cases, be expressed as a simple composition-weighted linear additive model of pure responses, with a single term per component contribution. Although such a model is known to be often followed because of the nature of the instrumental responses measured, the information related to the individual contributions involved cannot be drawn in a straightforward way from the raw measurements. The common purpose of all multivariate resolution methods is filling in this gap and providing the linear model of individual component contributions using solely the raw experimental measurements. Resolution methods are powerful approaches with low demanding needs; thus, none of the pure components in a system should be known beforehand and any information available about the system may be used, but is not essential either. Actually, the only indispensable condition is the inner linear structure of the data set. The mild requirements needed have promoted the use of resolution methods to tackle many chemical problems that could not be solved otherwise.

The evolution of resolution approaches through the improvement of exploratory tools, the adaptation to work with complex data structures, the ability to in-

roduce chemical and mathematical information in the algorithms and the better quality assessment of the results obtained is revisited here.

2. Resolution methods: general definition and limitations

All resolution methods decompose mathematically a global mixed instrumental response into the pure contributions due to each of the components in the system [1–9]. This global response is organized in a matrix D containing raw information about all the components present in the data set. Resolution methods allow for the decomposition of the initial mixture data matrix D into the product of two data matrices C and S^T , each of them including the pure response profiles of the n mixture or process components associated with the row and the column direction of the initial data matrix, respectively (see Fig. 1). In matrix notation, the expression valid for all resolution methods is:

$$D = CS^T + E \quad (1)$$

where $D(r \times c)$ is the original data matrix, $C(r \times n)$ and $S^T(n \times c)$ are the matrices containing the pure response profiles related to the data variation in the row direction and in the column direction, respectively, and $E(r \times c)$ is the error matrix, i.e. the residual variation of the data set that is not related to any chemical contribution. Parameters r and c are the number of rows and the number of columns of the original data matrix, respectively, and n is the number of chemical components in the mixture or process. C and S^T often refer to concentration profiles and spectra (hence their names), although resolution methods are proven to work in many other diverse problems [10–18].

From the early days in resolution research, the mathematical decomposition of a single data matrix, no matter the method used, is known to be subject to ambiguities [1]. This means that many sets of paired C - and S^T -type matrices can reproduce the original data set with the same fit quality. In plain words, the correct reproduction of the original data matrix can be achieved by using response profiles differing in shape (rotational ambiguity) or in magnitude (intensity ambiguity) from the sought (true) ones [19].

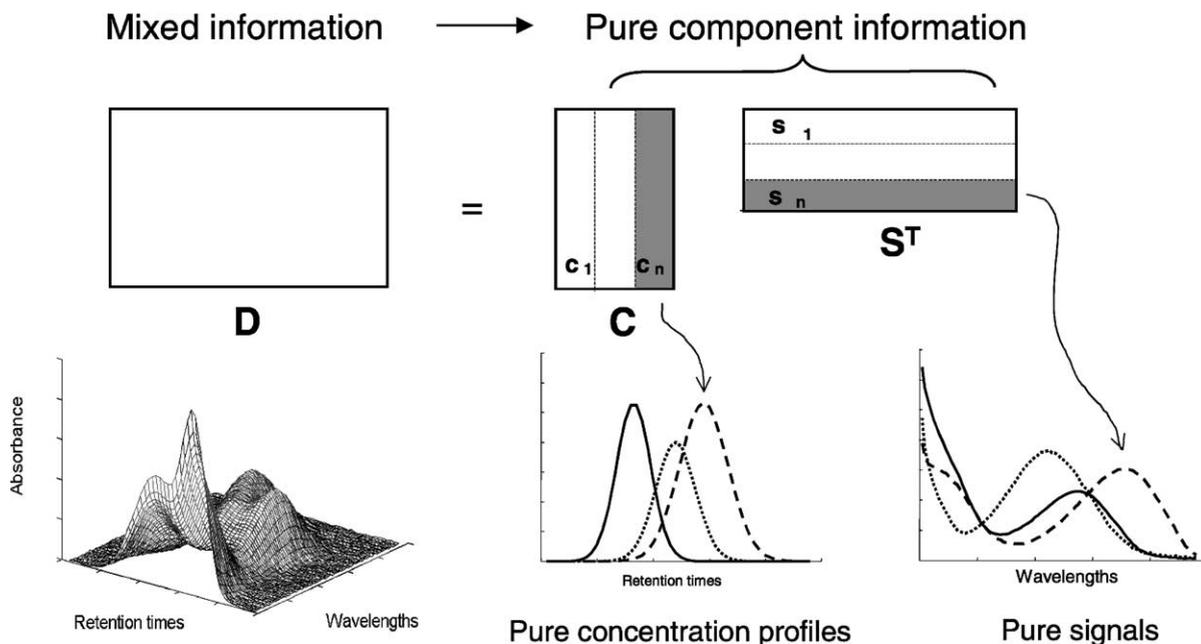


Fig. 1. Data set decomposition provided by resolution methods and graphical representation of the information obtained for a HPLC-DAD example.

These two kinds of ambiguities can be easily explained. The basic equation associated with resolution methods, $D = CS^T$, can be transformed as follows:

$$D = C(TT^{-1})S^T \quad (2)$$

$$D = (CT)(T^{-1}S^T) \quad (3)$$

$$D = C'S'^T \quad (4)$$

where $C' = CT$ and $S'^T = (T^{-1}S^T)$ describe the D matrix as correctly as the true C and S^T matrices do, though C' and S'^T are not the sought solutions. The rotational ambiguity problem indicates that a resolution method can potentially provide as many solutions as T matrices can exist, i.e. infinite, unless C and S are forced to obey certain conditions. In a hypothetical case with no rotational ambiguity, the basic resolution model could still be rewritten as shown below:

$$D = \sum_{i=1}^n \left(\frac{1}{k_i} c_i \right) (k_i s_i^T) \quad (5)$$

$$D = C'S'^T \quad (6)$$

where k_i are scalars. The concentration profiles of the new C' matrix would have the same shape as the real

ones, but being k_i times smaller, whereas the spectra of the new S' matrix would be shaped like the S spectra, though k_i times more intense.

How seriously ambiguity problems can affect resolution results has a different answer in each data set. In different manners, the main purpose of past and current investigation is finding ways to suppress, when feasible, or minimize this source of uncertainty in resolution results and to assess the effect on the quality of the profiles recovered.

3. Resolution-oriented exploratory analysis

One of the attractive features of the resolution methods is that there is no need of any previous knowledge, either chemical or mathematical, to analyse the data set of interest. However, the preliminary information that can be obtained from an exploratory analysis of the measurements can influence positively the resolution of the system. Thus, exploratory information can be used to build good initial estimates of concentration profiles and responses. It can also be included in the optimisation process in the form of constraints

and can help to validate the results obtained. When a good exploratory analysis is carried out, the resolved profiles cannot disagree with the previous information obtained.

Many exploratory procedures are often methods derived from Principal Component Analysis (PCA), one of the most basic and widely used chemometric tools devoted to find the number and direction of the relevant sources of variation in a bilinear data set [7,20,21]. The information provided by the global and local application of Principal Component Analysis to the data set can be essential in resolution. Thus, the first step in many resolution methods is the determination of the total number of chemical components in the data set and the ambiguity of the final solutions depends basically on the distribution and overlap of these compounds along the data set, i.e. on the local rank information obtained.

Much recognition should be given to the pioneering methods Evolving Factor Analysis (EFA) [22–24] and Fixed Size Moving Window-Evolving Factor Analysis (FSMW-EFA) [25,26] which are still widely used. EFA is specially applicable to sequential processes and performs subsequent PCA runs on windows gradually enlarged by addition of a row in the process direction. This analysis is performed building the windows from the beginning to the end of the process (forward EFA) and in the opposite direction (backward EFA). The proper display of the evolution of the eigenvalues obtained from these PCA runs allows for the detection and location of the emergence and decay of the compounds in a process and, as a consequence, the concentration window and the zero-concentration region for each component in the system are easily set (see Fig. 2b). The diverse uses of this latter information have given rise to most of the non-iterative resolution methods [27–36]. The information obtained from EFA allows for the derivation of estimates of the concentration profiles for these compounds, that are often used as starting point in iterative resolution methods [37,38]. FSMW-EFA conducts a series of PCA runs moving a window of a fixed size one row downwards every time from top to bottom of the data set. The visualisation of these results gives a local rank map of the data set, i.e. a representation of how many components are simultaneously present in the different zones of the data set (see Fig. 2b). The local rank map helps to locate easily selective zones in the data set and to know

the degree of compound overlap along the data set. It is especially sensitive to the detection of minor species. New algorithms based on EFA and FSMW-EFA have refined the performance of the parent methods [39,40] and have widened their applicability to the study of systems with concurrent processes [41] or complex spatial structure, such as spectroscopic images [42].

Other exploratory methods in resolution work finding the purest variables (rows or columns) in a data matrix, i.e. those marking the most dissimilar row and column profiles in the raw measurements [43]. Some of these procedures work on the abstract space of principal components [44,45] and some others directly on the space of the real variables of the measurements [46–48]. As additional benefits, some of them can also be used to determine the total number of components in the data set and are very sensitive to detect minor compounds. The purest variables in the concentration direction provide the purest responses in the original measurements and the purest response variables give the purest concentration profiles. These purest profiles are used as starting point in iterative resolution methods and are either the most dissimilar in the data set or those spanning the data space most efficiently according to the criterion of the selection method. This kind of initial estimates are particularly useful for data sets where the evolution of the concentration profiles of the different compounds is not sequential, as it is the case of spectroscopic images [42,49] or environmental data [14,15].

4. Introducing knowledge in the optimisation: the role of constraints

Although resolution does not require previous knowledge about the chemical system under study, additional knowledge, when existing, can be used to improve the results obtained and to tailor the sought pure profiles according to certain known features.

The working procedure of the different resolution methods is very diverse but many of them start with initial estimates of C or S and work by optimising iteratively the concentration and/or response profiles using the available information about the system [19,50–55]. The introduction of this information is carried out through the implementation of constraints. A constraint can be defined as any mathematical

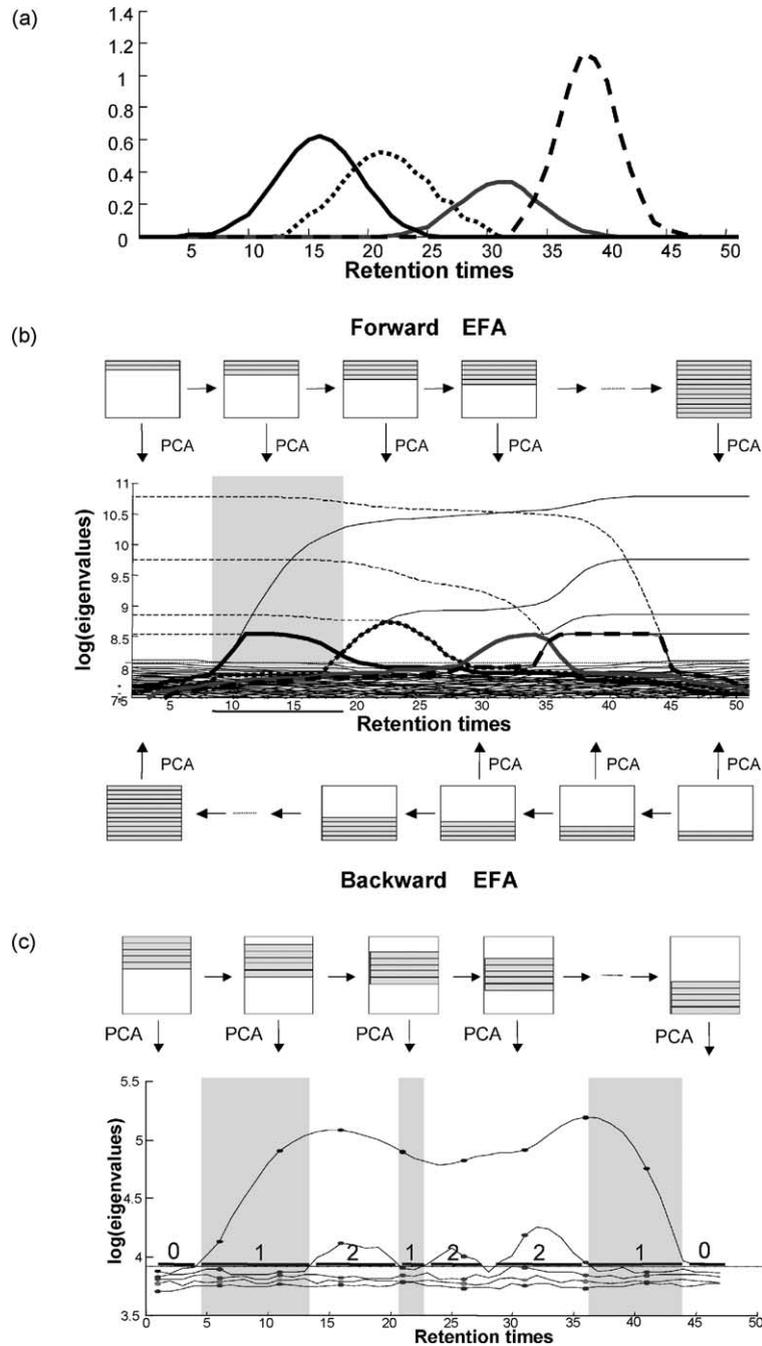


Fig. 2. (a) Concentration profiles of an HPLC-DAD data set. (b) Information derived from the data set in (a) by Evolving Factor Analysis (EFA): scheme of PCA runs performed. Combined forward EFA (solid black lines) and backward EFA (dashed black lines) plot. The thick lines with different linestyles are the derived concentration estimates. The shaded zone marks the concentration window for the first eluting compound. The rest of the elution range is the zero-concentration window. (c) Information derived from the data set in (a) by Fixed Size Moving Window-Evolving Factor Analysis (FSMW-EFA): scheme of the PCA runs performed. The straight lines and associated numbers mark the different windows along the data set as a function of their local rank (number). The shaded zones mark the selective concentration windows (rank 1).

or chemical property systematically fulfilled by the whole system or by some of its pure contributions [56]. Constraints are translated into mathematical language and force the iterative optimisation process to model the profiles respecting the conditions desired.

Although generally accepted nowadays, there is still a certain controversy related to the application of constraints, due presumably to misuses and rudimentary implementations of these conditions in the past. The application of constraints should be always prudent and soundly grounded and they should only be set when there is an absolute certainty about the validity of the constraint. Even a potentially applicable constraint can play a negative role in the resolution process when factors like experimental noise or instrumental problems distort the related profile so as it is no longer obeyed or when the profile is modified so roughly that the convergence of the optimisation process is seriously damaged. Resolution methods have progressed through the new formulation and better implementation of constraints.

Most of the constraints firstly implemented were directly linked to chemical properties fulfilled by the pure concentration or response profiles. Non-negativity applied to concentrations and to many instrumental signals [1,3,19,50,51,53,57]. Unimodality (i.e. the presence of only one maximum per profile) was useful in many concentration profiles related to processes, like reaction profiles or peaks in a chromatographic elution process [19,51,53,58,59]. It also applied to some particular signals, like some voltammetric responses [60]. Closure (or mass balance) was valid in many reaction systems as well [19,22,28,38,53]. The most recent progress in chemical-related constraints refers to the implementation of a physicochemical model into the resolution process [61–69]. Such a strategy has reconciled the separate worlds of hard- and soft-modelling and has allowed for the resolution of chemical systems unable to be successfully tackled by any of these two pure methodologies. In resolution methods where hard- and soft-modelling coexist, advantages from both aspects can be taken. Thus, the strictness of the hard model constraint decreases dramatically the ambiguity of the related profiles and provides fitted parameters of physicochemical [61–68] and analytical interest [69], such as equilibrium and kinetic constants or total analyte concentrations (see Fig. 3). The soft-part of

the algorithms allows for the modelling of complex systems, where the central reaction system evolves in the presence of absorbing interferents [63,69].

Chemical information associated with the knowledge of pure spectra or pure concentration profiles can also be introduced in the optimisation as an additional equality constraint [5,19,28,53,62]. The known profiles may be set to be invariant along the iterative process. Following this concept, the knowledge of a profile does not need to be complete to be used. When some elements are known, they can also be fixed. This has opened the possibility to use two-way resolution methods for quantitative purposes. Thus, some data sets analogous to those used in multivariate calibration, formed by the signals recorded in a series of calibration and unknown samples, can be analysed. The quantitative information is obtained by resolving the system using an additional constraint that fixes the known concentration values of the analyte(s) in the calibration samples in the related concentration profile(s) [70]. This approach provides also the qualitative information related to resolution methods in the form of pure signal profiles associated with the analyte and interferents and has proven to be as powerful as classical multivariate methods in examples where the net signal of the analyte is not very minor.

Other constraints are related to mathematical features and can be applied to all data sets, regardless of their chemical nature. These constraints are associated with the concept of local rank, i.e. how the number and distribution of components vary locally along the data set. The key constraint within this family is selectivity. It holds for concentration and spectral windows where only one component is present and suppresses completely the ambiguity linked to some of the profiles in the system. The strong effect of this constraint and the easy link with the analogous chemical concept explain its early and wide application in resolution problems [1,5,19,53]. Not so common, but equally recommended is the use of other local rank constraints in iterative resolution methods [19,53,71]. Setting which components are absent in data set windows with a number of components smaller than the total rank always contribute to the correct resolution of profiles and can be particularly helpful in multicomponent systems like spectroscopic images or mixtures, where more process-related constraints (unimodality, closure, etc.) are not applicable.

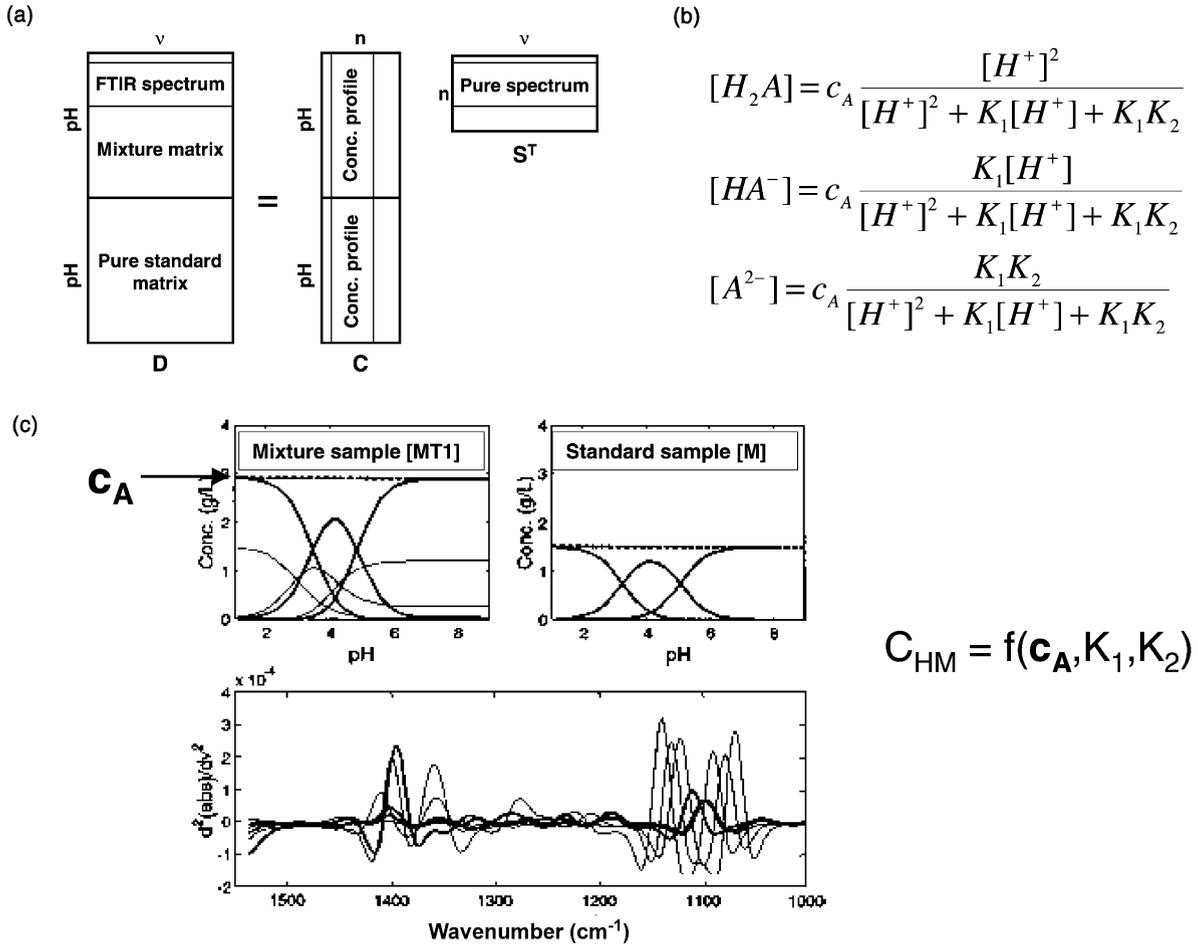


Fig. 3. Application of a hard-model constraint based on the acid–base equilibrium model to an FTIR titration of an analyte in the presence of interferences. (a) Structure of the column-wise data set with a sample and a standard matrix. (b) Physicochemical model used for the diprotic acid analyte in the hard-model constraint. (c) Pure concentration profiles and 2nd derivative spectra obtained. The bold lines refer to the analyte system. Thin solid lines come from an evolving interferent. C_A : total concentration of the analyte in the sample is obtained as a parameter obtained from the fitted hard-model (see [69] for more details).

Improvements in the implementation have significantly increased the efficiency and reliability of constraints. Flexibility and smoothness could be the two key words to summarize the main objectives in this area. In this context, the word flexibility encompasses a double meaning. On one hand, it concerns the complete freedom of combinations in the way to constrain the profiles linked to the different directions and components in a data set. On the other hand, it affects the degree of tolerance in the application of the constraints, i.e. how strict we are going to be to consider

that a profile does not obey a certain constraint. To cope with noise-related problems in real data, the implementation of constraints often allows for small deviations from the ideal behaviour before correcting a profile [19,53,56].

Once a constraint should enter in action, there is a big variability in the way a profile can be corrected. When well implemented and fulfilled by the data set, constraints can be seen as the driving forces of the iterative process to the right solution and, often, they are found not to be active in the last part of

the optimisation process. The goal is modifying the wrong-behaved profile upsetting as little as possible the convergence of the iterative optimisation. The first methods to be used were substitution methods, where the wrong elements in a profile were updated by others that fulfilled the sought constraint (e.g. negative values were replaced by zeroes in non-negativity). These straight replacements could eventually bring the system far from the convergence pathway when done too abruptly in very complex systems. In these approaches, there has been a significant gain in smoothness of constraint formulation. The most recent alternative tends to implement the constraints in a least-squares sense [72–74], keeping unmodified the monotonic convergence associated with the main least-squares optimisation step, characteristic of many iterative resolution methods. Constraints like non-negativity, unimodality and other equality constraints have been successfully implemented following this approach [75–77]. Despite their optimal mathematical properties, there is still a need for a

larger flexibility and speed in the use of these more rigorously implemented constraints. Given the iterative nature of resolution methods, the power of substitution strategies to solve systems that need a very flexible application of constraints and the accepted fact that the final results obtained are not significantly different between rigorous and approximate constraint implementation methods in most cases, both alternatives should be kept in use and under active research to widen the span of chemical problems to be resolved.

5. Resolution of complex data arrangements

New kinds of data sets, such as those produced by the most recent instrumentation, like spectroscopic images, can be analysed with resolution methods that can adapt easily to their huge size and complex spatial structure [42,49,78]. A spectroscopic image can be displayed as a data cube with two dimensions related to the x , y coordinates of the surface scanned and a

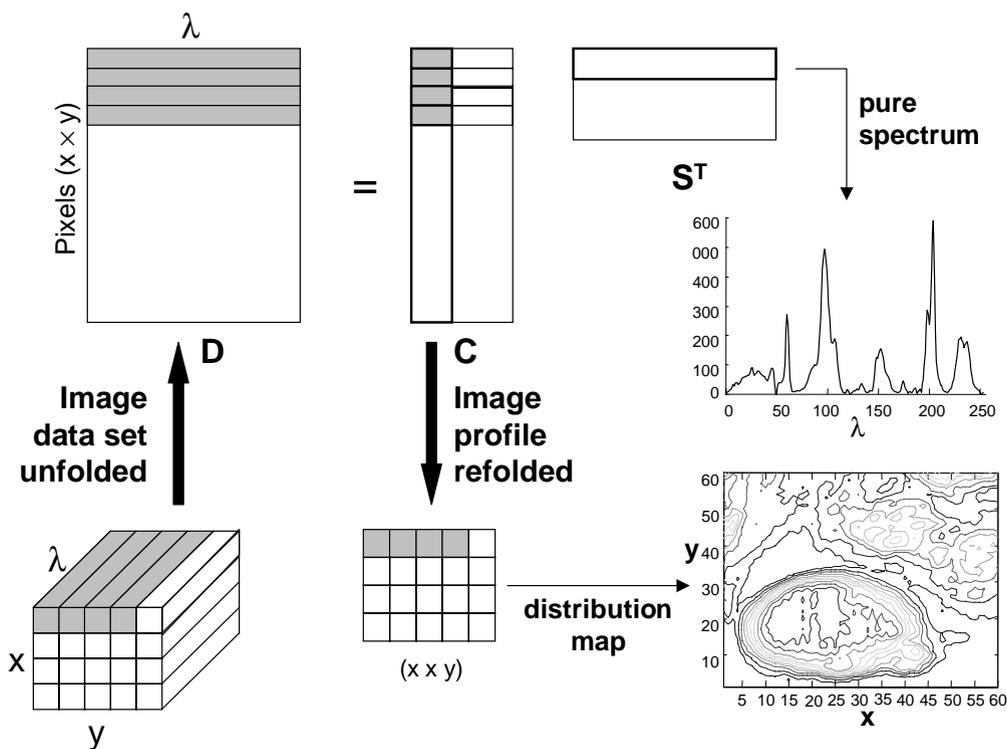


Fig. 4. General procedure followed for the resolution of spectroscopic images.

third spectral dimension. However, the measurement variation in an image data set follows a bilinear model, where the mixed signal recorded in each pixel is described by the concentration-weighted sum of the pure signals of the chemical compounds present. The pure spectra and distribution maps for the different compounds are resolved as in any other two-way data set. The only additional operations required are purely formal and consist of unfolding the original image cube into a matrix of pixel spectra and refolding the elements in the obtained concentration profiles according to the original spatial structure of the image to get the distribution map (see Fig. 4).

More important than the treatment of formally complex data sets, like the images mentioned above, advances in resolution have helped to expand the applicability of resolution methods from the treatment of one experiment (a single data matrix or a two-way data set) to the analysis of more complex data sets, organised forming a data cube or as row-, column- or row- and column-wise augmented matrices. These complex arrangements are known under the name of three-way data sets (see Fig. 5).

The resolution of three-way data sets provides pure profiles for the compounds that relate to the three informative directions of the data set. These triads can be directly provided by the algorithms [79–92], e.g. Parallel Factor Analysis (PARAFAC) or Tucker models, or the profiles related to the third direction can be easily derived from the results of the analysis of augmented matrices when they have a relevant chemical meaning, like in Multivariate Curve Resolution-Alternating Least Squares (MCR-ALS) (see Fig. 6) [53]. Working with three-way data sets has overcome typical problems associated with two-way data analysis. As a general gain, there is a significant decrease in the ambiguity of the resolution results because of the use of more and richer information about the multicomponent system under analysis.

The implementation of constraints finds in the three-way data sets an area for further development. Thus, the concept of flexibility extends to the possibility of constraining the profiles in the three directions of a data cube differently [81,93]. In augmented matrices, the profiles constrained and the constraints applied can also vary in the related C and S submatrices [7,19,53].

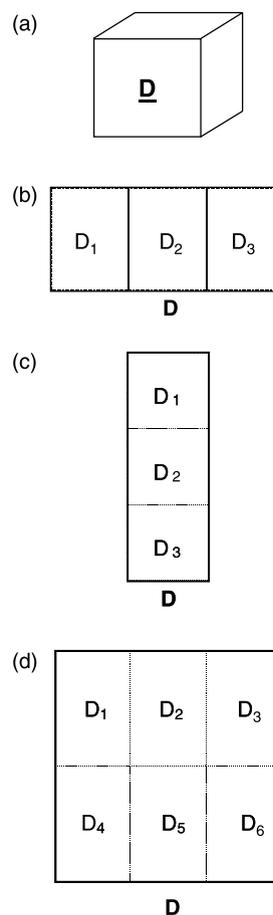


Fig. 5. Different arrangements of three-way data sets. (a) Data cube; (b) row-wise augmented matrix; (c) column-wise augmented matrix; (d) row-and column-wise augmented matrix.

Other kinds of constraints have arisen because of the special structure of three-way data sets. In mathematical terms, three-way data sets divide into trilinear and non-trilinear. A data set is said to be trilinear when each compound in all experiments treated together can be described by a triad of invariant pure profiles. This would be the case of a three-way data set formed by several samples monitored by fluorescence. The three kinds of profiles to be recovered (excitation spectrum, emission spectrum and sample-to-sample quantitative variation) are invariant for one compound, i.e. the pure spectral shape does not change from sample to sample. A non-trilinear data set could come from several HPLC-DAD runs, with the derived elution

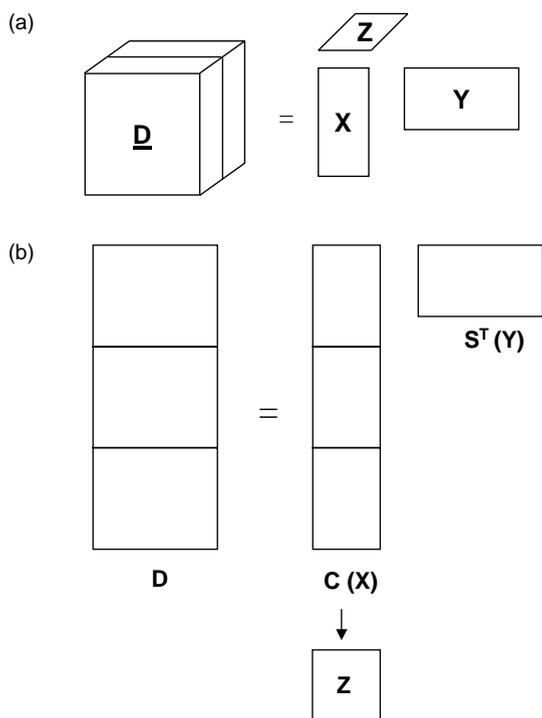


Fig. 6. Examples of data set decomposition derived from three-way resolution methods. (a) Resolution of trilinear data by PARAFAC. (b) Resolution of a column-wise data matrix by MCR-ALS.

profiles, spectra and quantitative profiles. In this case, the elution profile of a compound can suffer shifts or shape changes among runs and a common invariant elution profile per compound for all runs would not describe the system properly. Trilinearity is the most essential three-way related constraint because, when applicable, it ensures the uniqueness of the resolution results, i.e. the resolved profiles do not have ambiguity [79–84]. Some resolution methods, like PARAFAC or Direct Trilinear Decomposition, have an inner trilinear structure, whereas others, like Multivariate Curve Resolution-Alternating Least Squares, can optionally apply this constraint [7,53]. It is crucial to know the real inner structure of the three-way data set to be analysed (trilinear or non-trilinear) to select the three-way resolution method that can provide the best results [95–97]. Other constraints based on the chemical knowledge, like the correspondence of species, allow the introduction of information related to the presence or absence of

compounds in the different matrices treated together [53,94].

Working with three-way data sets improves the description of complex processes. In these situations, the same process can be monitored with different instrumental techniques and the different sets of measurements collected can be simultaneously analysed in a row-wise augmented data matrix [98,99]. Using this strategy, intermediates in protein folding processes have been successfully detected and modelled. In this concrete example, only the simultaneous analysis of measurements collected with circular dichroism in the far-UV (sensitive to changes in the protein secondary structure) and in the near-UV (sensitive to changes in the protein tertiary structure) provided enough information to resolve the concentration profile and spectrum of an intermediate conformation (with a denatured tertiary structure and a native-like secondary structure) in the thermal-induced protein folding of α -apolactalbumin (see Fig. 7). In other cases, the analysis of column-wise augmented matrices formed by experiments on the same system in slightly different conditions (e.g. different metal-to-ligand ratio in complexation processes or different temperatures) can help in the resolution of all species in all experiments, even in those where they have a minor contribution or an incomplete evolution [59,71,100–102]. Analysis of column- and row-wise augmented matrices is feasible when the two circumstances mentioned above concur [103,104].

If the two-way resolution analysis has mainly focused on the qualitative information derived from the shapes of the pure signal and concentration profiles, the resolution of three-way data sets often provides additional quantitative information. Indeed, the third direction in many three-way data sets is devoted to explain the sample-to-sample differences in concentration for the different compounds (e.g. in the example of series of chromatographic runs or the series of excitation–emission fluorescent samples). In contrast to multivariate calibration methods, the information related to a sample is not a vector (e.g. a spectrum) but a whole matrix. The quantitative information is derived from the varying scale of the pure concentration profile of a compound in the different matrices (samples) [94,105–111]. The richer information about the sample and the possibility to use the powerful three-way resolution methods has two obvious advantages with

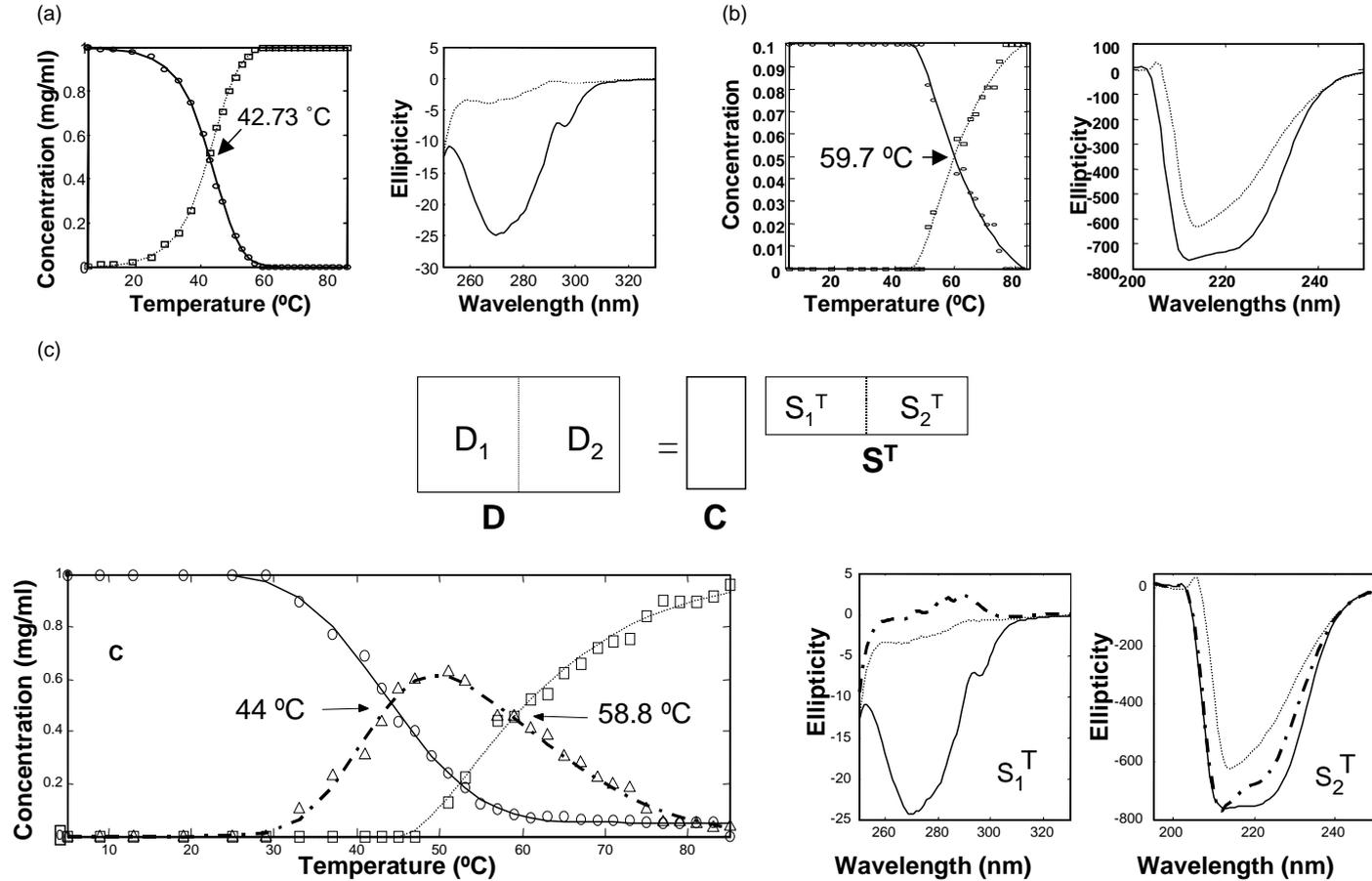


Fig. 7. Resolution of the protein folding of α -apolactalbumin. (a) Detection of changes in protein secondary structure (far-UV circular dichroism measurements). (b) Detection of changes in protein tertiary structure (near-UV circular dichroism measurements). (c) Complete description of protein folding. Resolution of the row-wise data set formed by near-UV (D_1) and far-UV (D_2) circular dichroism measurements. Solid line: native conformation; dash-dotted line: intermediate conformation; dotted line: unfolded conformation.

respect to traditional calibration approaches; strictly speaking, only one standard is needed because the concentration of the analyte in any other sample can be derived from the scale ratio with the profile in the standard and, equally important, the quantitation can be performed in the presence of unknown interferences.

Three-way resolution methods have allowed the analysis of rank-deficient systems, that could not be tackled with any two-way resolution approach. The name of rank-deficient systems originates because the number of mathematically distinguishable pure contributions (rank) is lower than the total number of chemical compounds in the data set. This happens when the profiles of some compounds (signals or concentration profiles) can be described as linear combinations of the profiles of others. This would happen when two compounds have identical signals, identical concentration profiles (e.g. B and C in an $A \rightarrow B + C$ reaction) or when several closed systems coexist [112,113]. In this situation, the failure of the two-way resolution methods does not come from the limitations of these approaches, but from the inner structure of the data matrix. Rank-deficiency can only be solved by matrix augmentation in the direction (signal or concentration) where this problem is present, i.e. adding one or more matrices with new and appropriate information so that the linear dependence of the pure profiles of the original rank-deficient data matrix does not hold anymore (e.g. column-wise augmenting a matrix with the reaction $A \rightarrow B + C$ with a matrix with the process $A \rightarrow B$). The analysis of the three-way full-rank augmented matrix will provide the correct resolved profiles for the compounds present in the original two-way rank-deficient matrix. [69,109,114,115].

6. Quality assessment of results

Until recently, research in resolution was mainly focused on the design and improvement of the algorithms used and not much in providing procedures to assess the quality of the results obtained. The main sources of uncertainty associated with the resolution results are the ambiguity of the recovered profiles and the experimental noise of the data. Providing methodologies to quantify this uncertainty is not only a topic

of interest in the current research, but a necessary step to introduce the use of resolution methods in standard analytical procedures.

The possible existence of ambiguity in resolution is known since the earliest research in this area [1]. After years of experience, it has been possible to set resolution theorems that indicate clearly the conditions needed to recover uniquely the pure concentration and signal profiles of a compound in a data set. These conditions depend mainly on the degree of overlap among the region of occurrence of the compound and the rest of constituents and on the general distribution of the different compound windows along the data set [116]. Therefore, in the same system, some profiles can be recovered uniquely and some others will be necessarily affected by a certain ambiguity. When ambiguity exists, a compound is represented by a band of feasible solutions and not by a unique profile. Calculating the boundaries of these bands is not straightforward and the first attempts proposed were valid only for systems with two or three components [3]. More recent approaches extended their applicability to systems with no limit in the number of contributions [4,117]. The latest tendency uses optimisation strategies to find the minimum and maximum solution band boundaries by minimising and maximising objective functions subject to selected constraints that, using different parameters, represent the ratio between the signal contribution of a certain compound and the total signal from all compounds in the data set. These strategies are more powerful than previous ones and allow for an accurate study of the effect of the different constraints in the magnitude of the bands of feasible solutions (see Fig. 8) [118,119].

Even in the absence of ambiguity, the experimental error contained in real data can propagate in the resolution results. This source of uncertainty affects the results of all kinds of data analysis methods and, in simpler approaches, like multivariate or univariate calibration, is easily quantified with the use of well established and generally accepted figures of merit. Although some figures of merit have been proposed for higher order calibration methods, [120] finding analytical expressions to assess the error associated with resolution results is an extremely complex problem because of the huge number of non-linear parameters that are calculated, as many as elements in

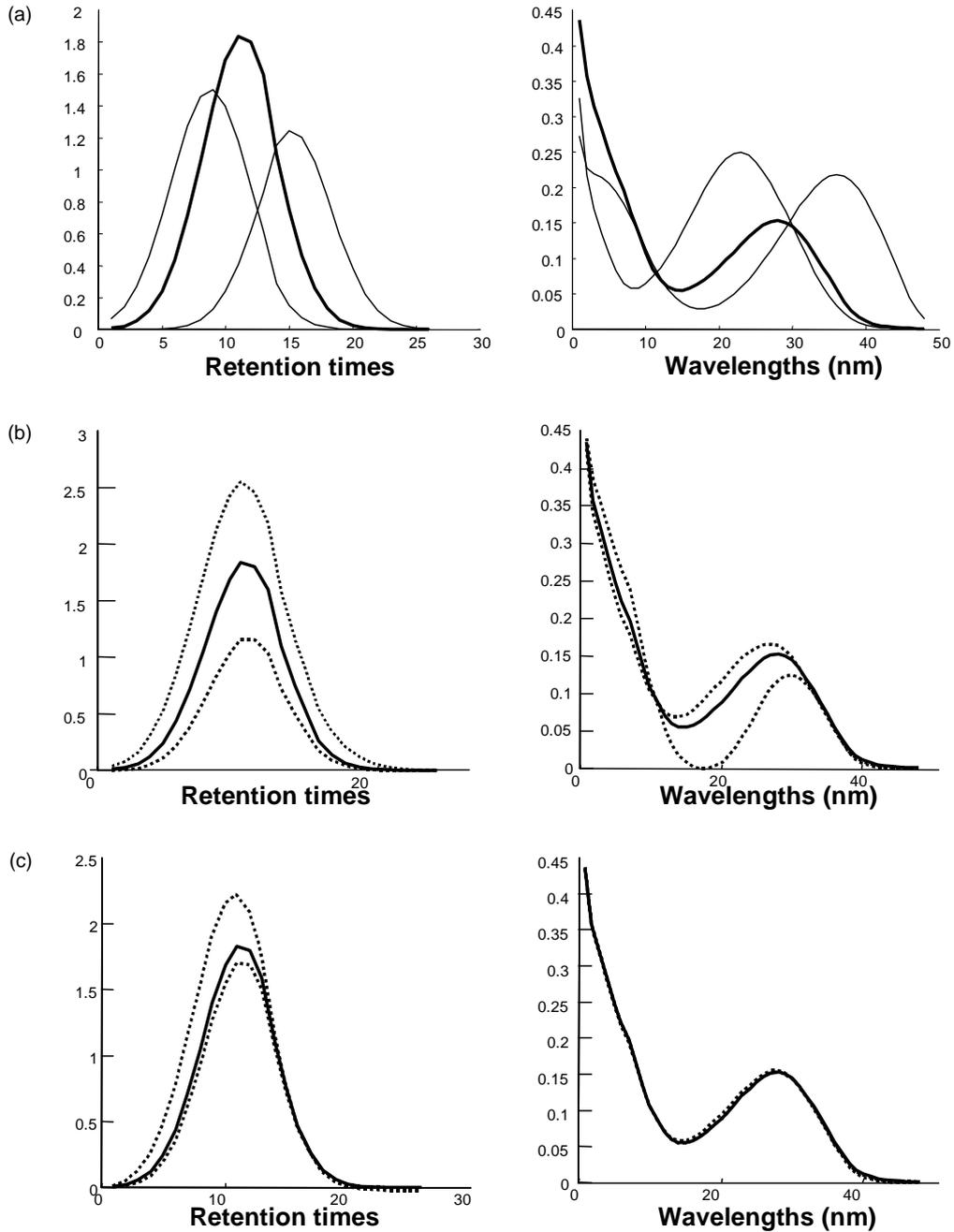


Fig. 8. (a) True concentration profiles and spectra of an HPLC-DAD data set. (b) Resolved concentration profile and spectrum for the bold component in (a) (solid line: result from MCR-ALS analysis, dashed lines: upper and lower solution boundaries subject to non-negativity, unimodality and spectra normalisation). (c) Resolved concentration profile and spectrum for the bold component in (a) (solid line: result from MCR-ALS analysis, dashed lines: upper and lower solution boundaries subject to local rank, non-negativity and spectra normalisation) (for more detail, see [119]).

all the pure profiles resolved. To overcome this problem and still give a reliable approximate estimation of the error propagation in resolution, other strategies known under the general name of resampling are used [121–123]. In these strategies, an estimation of the dispersion in the resolution results that would be obtained after the resolution of a huge number of replicates is obtained. To simulate these replicates, a data matrix can be resolved repeatedly removing some rows or columns every time or, in the case of a three-way data set removing complete data matrices [62,124] (jackknife), the complete data set can be resolved multiple times after adding a certain amount of noise on top of the experimental measurements (noise-added method) [121] or the replicates of the data set can be constructed by addition of a certain amount of noise to a noise-free simulated or reproduced data set (Monte Carlo simulations). These strategies that provide an enormous number of results coming from the different resolution runs allow for an estimation of the uncertainty coming from the noise in the resolved profiles and, eventually, for the computation of the accuracy linked to some parameters, like rate constants or equilibrium constants, that can be indirectly derived from the profiles obtained [62,125].

Although ambiguity and noise are two distinct sources of uncertainty in resolution, their effect on the resolution results cannot be considered completely independent. Thus, the boundaries of the compound windows can be clearly blurred due to the noise effect and this can give rise to ambiguities that would be absent in noise-free data sets. A definite advance would be the proposal of approaches that may consider this combined effect in the estimation of resolution uncertainty.

7. Conclusions

The evolution of multivariate resolution methods is continuous, fast and efficient. With the improvements in exploratory tools, introduction of chemical information, adaptation for the analysis of complex data structures and quality assessment of the results, it can be envisioned an important increase in the range of application of these methods and a more generalised and standardised use of multivariate resolution by the analytical community.

References

- [1] W.H. Lawton, E.A. Sylvestre, *Technometrics* 13 (1971) 617.
- [2] E.A. Sylvestre, W.H. Lawton, M.S. Maggio, *Technometrics* 16 (1974) 353.
- [3] O.S. Borgen, B.R. Kowalski, *Anal. Chim. Acta* 174 (1985) 1.
- [4] O.S. Borgen, N. Davidsen, Z. Myngyang, O. Oyen, *Mikrochim. Acta II* (1986) 1.
- [5] J. Craig Hamilton, P.J. Gemperline, *J. Chemometr.* 4 (1990) 1.
- [6] W. Windig, *Chemom. Intell. Lab. Sys.* 16 (1992) 1.
- [7] A. de Juan, E. Casassas, R. Tauler, *Soft modeling of analytical data*, *Encyclopedia of Analytical Chemistry: Instrumentation and Applications*, Wiley, New York, 2000, p. 9800.
- [8] Y. Liang, O.M. Kvalheim, *Fresenius J. Anal. Chem.* 370 (2001) 694.
- [9] J.H. Jiang, Y. Ozaki, *Appl. Spectr. Rev.* 37 (2002) 321.
- [10] M. Esteban, C. Ariño, J.M. Díaz-Cruz, M.S. Díaz-Cruz, R. Tauler, *Trends Anal. Chem.* 19 (2000) 49.
- [11] C.G. Zampronio, L.A.B. Moraes, M.N. Eberlin, R.J. Poppi, *Anal. Chim. Acta* 446 (2001) 495.
- [12] B.H. Cruz, J.M. Díaz-Cruz, I. Sestakova, J. Velek, C. Ariño, M. Esteban, *J. Electroanal. Chem.* 520 (2002) 111.
- [13] M.C. Antunes, J.E. Simão, A.C. Duarte, M. Esteban, R. Tauler, *Anal. Chim. Acta* 459 (2002) 291.
- [14] P.K. Hopke, *Receptor Modeling in Environmental Chemistry*, Wiley, New York, 1985.
- [15] 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, 2000, p. 689.
- [16] J.S. Salau, R. Tauler, J.M. Bayona, I. Tolosa, *Environ. Sci. Tech.* 31 (1997) 3482.
- [17] Y.L. Xie, P.K. Hopke, P. Paatero, *J. Chemometr.* 12 (1998) 357.
- [18] X.H. Song, A.V. Polissar, P.K. Hopke, *Atmos. Environ.* 35 (2001) 5277.
- [19] R. Tauler, A.K. Smilde, B. R. Kowalski, *J. Chemometr.* 9 (1995) 31.
- [20] E.R. Malinowski, *Factor Analysis in Chemistry*, third ed., Wiley-VCH, New York, 2002.
- [21] D.L. Massart, B.G.M. Vandeginste, L.M.C. Buydens, S. de Jong, P.J. Lewi, J. Smeyers-Verbeke, *Handbook of Chemometrics and Qualimetrics*, in: *Data Handling in Science and Technology*, vol. 20. Elsevier, Amsterdam, 1997.
- [22] H. Gampp, M. Maeder, C.J. Meyer, A.D. Zuberbühler, *Talanta* 32 (1985) 1133.
- [23] M. Maeder, A.D. Zuberbühler, *Anal. Chim. Acta* 181 (1986) 287.
- [24] M. Maeder, *Anal. Chem.* 59 (1987) 527.
- [25] H.R. Keller, D.L. Massart, *Anal. Chim. Acta* 246 (1991) 379.
- [26] H.R. Keller, D.L. Massart, Y.Z. Liang, O.M. Kvalheim, *Anal. Chim. Acta* 263 (1992) 29.
- [27] M. McCue, E.R. Malinowski, *J. Chromatog. Sci.* 21 (1983) 229.

- [28] H. Gampp, M. Maeder, C.J. Meyer, A.D. Zuberbühler, *Anal. Chim. Acta* 193 (1987) 287.
- [29] E. R. Malinowski, *J. Chemometr.* 6 (1992) 29.
- [30] E.R. Malinowski, *J. Chemometr.* 10 (1996) 273.
- [31] O.M. Kvalheim, Y.Z. Liang, *Anal. Chem.* 64 (1992) 936.
- [32] Y.Z. Liang, O.M. Kvalheim, H.R. Keller, D.L. Massart, P. Kiechle, F. Erni, *Anal. Chem.* 64 (1992) 946.
- [33] R. Manne, H. Shen, Y. Liang, *Chemom. Intell. Lab. Syst.* 45 (1999) 171.
- [34] R. Manne, B.V. Grande, *Chemom. Intell. Lab. Syst.* 50 (2000) 35.
- [35] C. Mason, M. Maeder, A.C. Whitson, *Anal. Chem.* 73 (2001) 1587.
- [36] J.H. Jiang, S. Šašić, R. Yu, Y. Ozaki, *J. Chemometr.* 17 (2003) 186.
- [37] H. Gampp, M. Maeder, Ch.J. Meyer, A. Zuberbühler, *Talanta* 33 (1986) 943.
- [38] R. Tauler, E. Casassas, *J. Chemometr.* 3 (1988) 151.
- [39] J. Toft, O.M. Kvalheim, *Chemom. Intell. Lab. Syst.* 19 (1993) 65.
- [40] A.C. Whitson, M. Maeder, *J. Chemom.* 15 (2001) 475.
- [41] A. de Juan, S. Navea, J. Diewok, R. Tauler, *Chemom. Intell. Lab. Syst.*, submitted for publication.
- [42] A. de Juan, R. Dyson, C. Marcolli, M. Rault, R. Tauler, M. Maeder, *Trends Anal. Chem.* (2003), in press.
- [43] F. Cuesta-Sánchez, B. van den Bogaert, S.C. Rutan, D.L. Massart, *Chemom. Intell. Lab. Syst.* 34 (1996) 139.
- [44] E.R. Malinowski, *Anal. Chim. Acta* 134 (1982) 129.
- [45] A. de Juan, B. van den Bogaert, F. Cuesta Sánchez, D.L. Massart, *Chemom. Intell. Lab. Syst.* 33 (1996) 133.
- [46] W. Windig, J. Guilment, *Anal. Chem.* 63 (1991) 1425.
- [47] F.C. Sánchez, J. Toft, B. van den Bogaert, D.L. Massart, *Anal. Chem.* 68 (1996) 79.
- [48] B.V. Grande, R. Manne, *Chemom. Intell. Lab. Syst.* 50 (2000) 19.
- [49] J.J. Andrew, T.M. Hancewicz, *Appl. Spectrosc.* 52 (1998) 797.
- [50] P.J. Gemperline, *J. Chem. Inf. Comput. Sci.* 24 (1984) 206.
- [51] B.G.M. Vandeginste, W. Derks, G. Kateman, *Anal. Chim. Acta.* 173 (1985) 253.
- [52] R. Tauler, E. Casassas, *Anal. Chim. Acta* 223 (1989) 257.
- [53] R. Tauler, *Chemom. Intell. Lab. Syst.* 30 (1995) 133.
- [54] E.J. Karjalainen, *Chemom. Intell. Lab. Syst.* 7 (1989) 31.
- [55] M. Leger, P.D. Wentzell, *Chemom. Intell. Lab. Syst.* 62 (2002) 171.
- [56] A. de Juan, Y. Vander Heyden, R. Tauler, D.L. Massart, *Anal. Chim. Acta* 346 (1997) 307.
- [57] E. Spjøtvoll, H. Martens, R. Volden, *Technometrics* 24 (1982) 173.
- [58] P.J. Gemperline, *Anal. Chem.* 58 (1986) 2656.
- [59] R. Tauler, A. Izquierdo-Ridorsa, E. Casassas, *Chemom. Intell. Lab. Syst.* 18 (1993) 293.
- [60] J.M. Díaz-Cruz, R. Tauler, B. Grabaric, M. Esteban, E. Casassas, *J. Electroanal. Chem.* 393 (1995) 7.
- [61] S. Bijlsma, A.K. Smilde, *Anal. Chim. Acta* 396 (1999) 231.
- [62] S. Bijlsma, A.K. Smilde, *J. Chemometr.* 14 (2000) 541.
- [63] A. de Juan, M. Maeder, M. Martínez, R. Tauler, *Chemom. Intell. Lab. Syst.* 54 (2000) 123.
- [64] J.M. Díaz-Cruz, J. Agulló, M.S. Díaz-Cruz, C. Ariño, M. Esteban, R. Tauler, *Analyst* 126 (2001) 371.
- [65] E. Bezemer, S.C. Rutan, *Anal. Chem.* 73 (2001) 4403.
- [66] E. Bezemer, S.C. Rutan, *Chemom. Intell. Lab. Syst.* 59 (2001) 19.
- [67] P. Jandanklang, M. Maeder, A.C. Whitson, *J. Chemometr.* 15 (2001) 511.
- [68] W. Windig, J.P. Hornak, B. Antalek, *J. Magnet. Res.* 133 (1998) 298.
- [69] J. Diewok, A. de Juan, M. Maeder, R. Tauler, B. Lendl, *Anal. Chem.* 75 (2003) 641.
- [70] M.C. Antunes, J.E.J. Simao, A.C. Duarte, R. Tauler, *Analyst* 127 (2002) 809.
- [71] A. de Juan, M. Maeder, M. Martínez, R. Tauler, *Anal. Chim. Acta* 442 (2001) 337.
- [72] C.L. Lawson, R.J. Hanson. *Solving Least-Squares Problems*, Prentice-Hall, Englewood Cliffs, NJ, 1974.
- [73] R.J. Hanson, K.H. Haskell, *ACM Trans. Math. Softw.* 8 (1982) 323.
- [74] K.H. Haskell, R.J. Hanson, *Math. Program.* 21 (1981) 98.
- [75] R. Bro, S. de Jong, *J. Chemometr.* 11 (1997) 393.
- [76] R. Bro, N.D. Sidiropoulos, *J. Chemometr.* 12 (1998) 223.
- [77] M.H. Van Benthem, M.R. Keenan, D.M. Haaland, *J. Chemom.* 16 (2002) 613.
- [78] J.H. Wang, P.H. Hopke, T.M. Hancewicz, S.L. Zhang, *Anal. Chim. Acta* 476 (2003) 93.
- [79] R.A. Harshman, *UCLA Working Papers in Phonetics*, 1970, p. 1.
- [80] J.D. Carroll, J. Chang, *Psychometrika* (1970) 283.
- [81] R. Bro, *Chemom. Intell. Lab. Syst.* 38 (1997) 149.
- [82] N.M. Faber, R. Bro, P.K. Hopke, *Chemom. Intell. Lab. Syst.* 65 (2003) 119.
- [83] E. Sánchez, B.R. Kowalski, *Anal. Chem.* 58 (1986) 496.
- [84] E. Sánchez, B.R. Kowalski, *J. Chemometr.* 4 (1990) 29.
- [85] L.R. Tucker, in: C.W. Harris (Ed.), *Problems in Measuring Change*, vol. 122, University of Wisconsin Press, Madison, 1963.
- [86] L.R. Tucker, *Psychometrika* 31 (1966) 279.
- [87] P.M. Kroonenberg, J. de Leeuw, *Psychometrika* 45 (1980) 69.
- [88] A.K. Smilde, R. Tauler, J.M. Henshaw, L.W. Burgess, B.R. Kowalski, *Anal. Chem.* 66 (1994) 3345.
- [89] A.K. Smilde, Y.D. Wang, B.R. Kowalski, *J. Chemometr.* 8 (1994) 21.
- [90] C.A. Andersson, R. Bro, *Chemom. Intell. Lab. Syst.* 42 (1998) 93.
- [91] H.A.L. Kiers, J.M.F. Ten Berge, R. Bro, *J. Chemometr.* 13 (1999) 275.
- [92] R. Bro, J.J. Workman Jr., P.R. Mobley, B.R. Kowalski, *Appl. Spec. Rev.* 32 (1997) 237.
- [93] C.A. Andersson, R. Bro, *Chemom. Intell. Lab. Syst.* 52 (2000) 1.
- [94] R. Tauler, D. Barceló, *Trends Anal. Chem.* 12 (1993) 319.
- [95] A. de Juan, S.C. Rutan, R. Tauler, D.L. Massart, *Chemom. Intell. Lab. Syst.* 19 (1998) 19.

- [96] A.K. Smilde, R. Tauler, J. Saurina, R. Bro, *Anal. Chim. Acta* 398 (1999) 237.
- [97] A. de Juan, R. Tauler, *J. Chemometr.* 15 (2001) 749.
- [98] J. Mendieta, M.S. Díaz-Cruz, M. Esteban, R. Tauler, *Biophys. J.* 74 (1998) 2876.
- [99] S. Navea, A. de Juan, R. Tauler, *Anal. Chem.* 64 (2002) 6031.
- [100] R. Tauler, B.R. Kowalski, S. Fleming, *Anal. Chem.* 65 (1993) 2040.
- [101] K. de Braekeleer, A. de Juan, D.L. Massart, *J. Chromatogr. A* 832 (1999) 67.
- [102] R. Gargallo, R. Tauler, A. Izquierdo-Ridorsa, *Anal. Chem.* 69 (1997) 1785.
- [103] S. Navea, A. de Juan, R. Tauler, *Anal. Chim. Acta* 446 (2001) 187.
- [104] J. Jaumot, N. Escaja, R. Gargallo, C. González, E. Pedroso, R. Tauler, *Nucl. Acid Res.* 30 (2002) e92/1.
- [105] R. Bro, H. Heimdahl, *Chemom. Intell. Lab. Syst.* 34 (1996) 85.
- [106] J. Saurina, S. Hernández-Cassou, R. Tauler, *Anal. Chem.* 69 (1997) 2329.
- [107] R.D. Jiji, G.A. Cooper, K.S. Booksh, *Anal. Chim. Acta* 397 (1999) 61.
- [108] J. Saurina, R. Tauler, *Analyst* 125 (2000) 2038.
- [109] J. Diewok, A. de Juan, R. Tauler, B. Lendl, *Appl. Spec.* 56 (2002) 40.
- [110] P.D. Wentzell, S.S. Nair, R.D. Guy, *Anal. Chem.* 73 (2001) 1408.
- [111] R.P.H. Nikolajsen, K. Booksh, A. Hansen, R. Bro, *Anal. Chim. Acta* 475 (2003) 137.
- [112] M. Amrhein, B. Srinivasan, D. Bonvin, M.M. Schumacher, *Chemom. Intell. Lab. Syst.* 33 (1996) 17.
- [113] A. Izquierdo-Ridorsa, J. Saurina, S. Hernández-Cassou, R. Tauler, *Chemom. Intell. Lab. Syst.* 38 (1997) 183.
- [114] J. Saurina, S. Hernández-Cassou, R. Tauler, A. Izquierdo-Ridorsa, *J. Chemometr.* 12 (1998) 183.
- [115] M.M. Reis, S.P. Gurden, A.K. Smilde, M.M.C. Ferreira, *Anal. Chim. Acta* 422 (2000) 21.
- [116] R. Manne, *Chemom. Intell. Lab. Syst.* 27 (1995) 89.
- [117] R.C. Henry, B.M. Kim, *Chemom. Intell. Lab. Syst.* 8 (1990) 205.
- [118] P.J. Gemperline, *Anal. Chem.* 71 (1999) 5398.
- [119] R. Tauler, *J. Chemometr.* 15 (2001) 627.
- [120] K. Faber, A. Lorber, B.R. Kowalski, *J. Chemometr.* 11 (1997) 419.
- [121] N.M. Faber, N. Gouda, *J. Chemometr.* 15 (2001) 169.
- [122] R. Wehrens, H. Putter, L.M.C. Buydens, *Chemom. Intell. Lab. Syst.* 54 (2000) 35.
- [123] G. Meinrath, *Chemom. Intell. Lab. Syst.* 51 (2000) 175.
- [124] J. Riu, R. Bro, *Chemom. Intell. Lab. Syst.* 65 (2003) 35.
- [125] J. Jaumot, R. Gargallo, R. Tauler, in preparation.