

## A review on chemometric techniques

M Venkata Kumari<sup>1\*</sup>, Davulsabgari Asif<sup>2</sup>, Shaik Aasik<sup>2</sup>, Shaik Karishma<sup>2</sup>, Shaik Mehar<sup>2</sup>

<sup>1</sup> Associate Professor, Department of Pharmaceutical Analysis, Nimra College of Pharmacy, Nimra Nagar Jupudi, Ibrahimpatnam, Vijayawada, Andhra Pradesh, India

<sup>2</sup> Department of Pharmacy, Nimra College of Pharmacy, Nimra Nagar Jupudi, Ibrahimpatnam, Vijayawada, Andhra Pradesh, India

### Abstract

Chemometric techniques involve the application of mathematical and statistical methods to analyze chemical data, enhancing the interpretation of complex datasets. These methods are crucial in various fields, including analytical chemistry, environmental monitoring, and pharmaceuticals. By employing multivariate analysis, chemometrics enables the extraction of meaningful information from high-dimensional data, such as spectral or chromatographic data. Common chemometric techniques include principal component analysis (PCA), partial least squares (PLS) regression, and cluster analysis. PCA reduces the dimensionality of data while preserving variance, making it easier to visualize and interpret relationships among variables. PLS regression is particularly effective for predictive modelling, linking multiple dependent and independent variables, and is widely used in calibration tasks. Cluster analysis helps in identifying groups within data, facilitating the classification of samples based on their chemical properties. These techniques also support quality control, method validation, and the development of robust analytical methods. For instance, in the pharmaceutical industry, chemometrics aids in optimizing formulations and ensuring product consistency. Moreover, the integration of chemometrics with machine learning and artificial intelligence is paving the way for advanced predictive analytics and automation in chemical analysis. In summary, chemometric techniques are indispensable for handling complex chemical data, providing powerful tools for data interpretation, quality assurance, and enhanced decision-making across various scientific disciplines. Their ongoing evolution continues to drive innovation in analytical methodologies, contributing significantly to advancements in chemistry and related fields.<sup>[1]</sup>

### Definition

Chemometric is a very wide discipline driven from Mathematics and Statistics methods of applications for interpretation & prediction called as data. Chemometrics Significant role are they can use to build algorithms for identification of analytical data and evaluation of data, it also used in the determination quality process for pharmaceutical drugs. To know more information click on Enter and observe the statistical data & mathematical data.<sup>[2]</sup>

**Keywords:** Chemometric techniques, analytical chemistry, dimensionality of data

### Introduction

Rakesh Shah *et al* 1989 identified the important process parameters and they proposed optimization of coating formulation, drug content uniformity design (mixing time/drum speed) method by chemometrics that widely range from basic statistics to signal processing factorial design calibration techniques curve fitting with factor analysis detection pattern recognition neural network. The term chemometrics, which was first carried out by a Swedish scientist Svante Wold in the year 1971 means

nothing but it is just simply making use of some mathematical and statistical methods to get more relevant information from chromatographic data. The field is then referred to as chemometrics, a term that has been defined in some cases, and not so completely in others i.e. the science used for correlating measurements performed on chemical systems or processes with their state using mathematical methods <sup>[25–30]</sup>. Areas and Principles that Make Help in Chemometrics There are several ways suggested to the <sup>[3]</sup>.

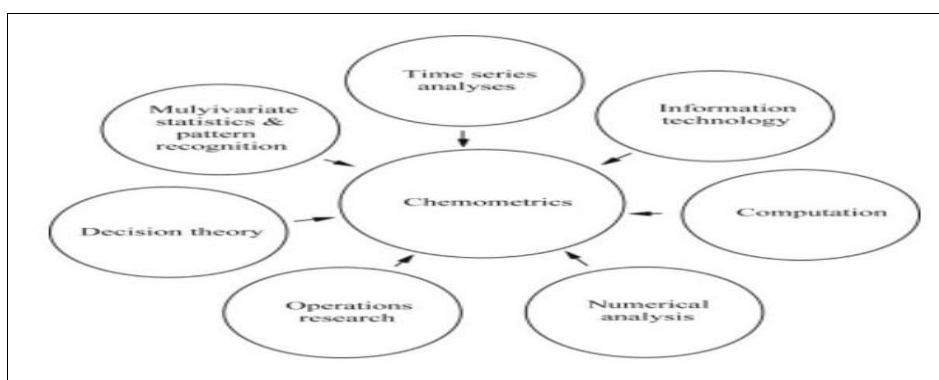


Fig 1

A flow diagram showing the generic steps involved in the HPLC-based analysis of herbal products.

These are involved in the destruction or non invasive online

techniques. Recently, due to the benefits they bring several non-destructive methodologies based on the spectroscopic techniques combined with the chemometric tools.

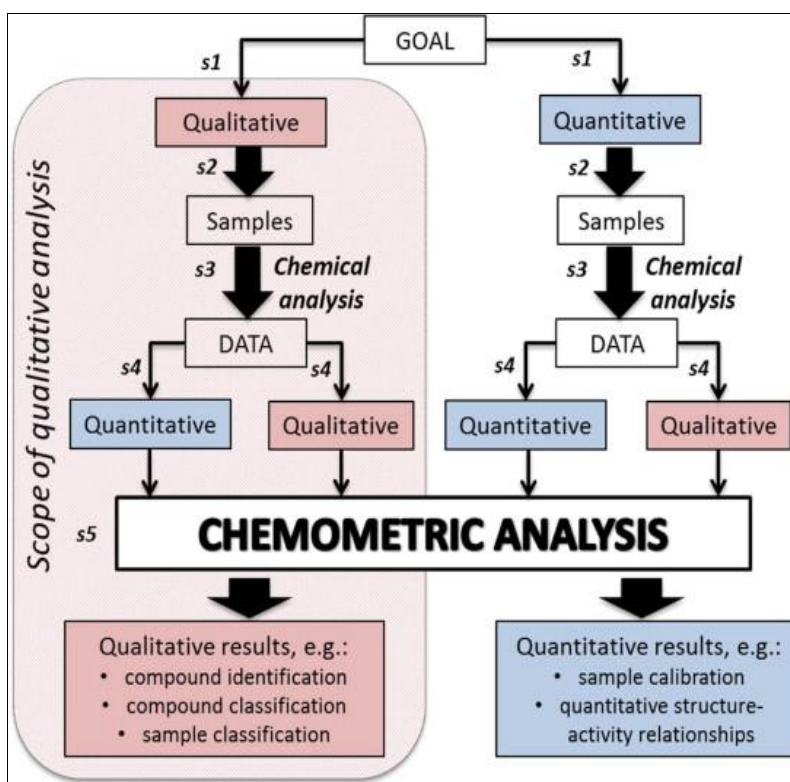


Fig 2

#### Chemometrics as tools for fraud/adulteration detection

Chemometric tools are have the proposed for the pharmaceutical quality check. poor quality of a pharmaceutical can be found in the market mainly due to two reasons

- Low Production of standards
- Fraud Attempts.

Counterfeited drugs are present in different fraud / adulteration. They could contain the no API (active pharmaceutical ingredients ) several methodologies have been proposed to detect the counterfeited / fraud/

Adulteration of the pharmaceutical ingredient. In these a major role played by chemometric.

#### Classification

The chemometric classification methods focus on the possibility of assigning an object (sample) to a class based on the result of a set of variables obtained from chemical measurements Unsupervised pattern recognition (UPR) can be differentiated from EDA in that the purpose of UPR is to detect the similarities among objects, while with EDA there is no specific prejudice as to whether or how many classes will be found<sup>[4]</sup>

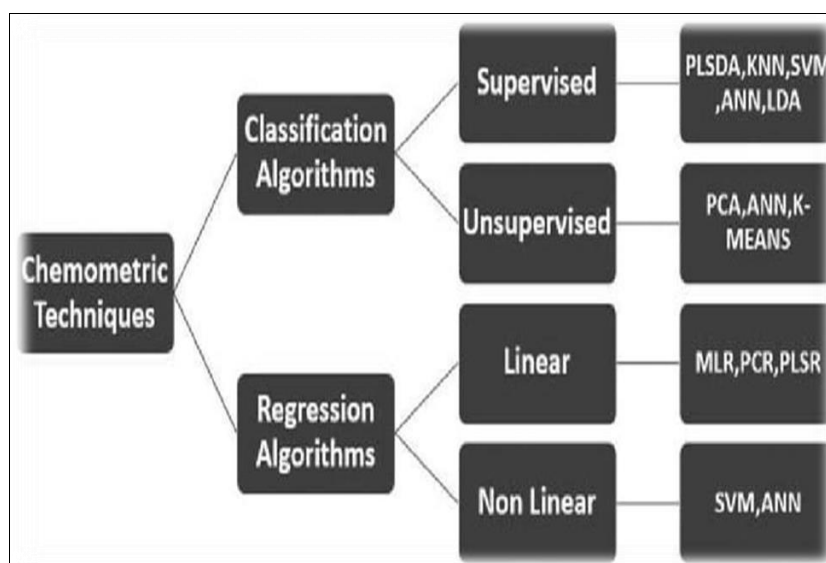


Fig 3

## Classification algorithms

### Supervised method

#### PLS-DA-PLS discriminate analysis

PLS-DA is a PLS-based method that can perform qualitative analysis or discriminate analysis.

Like PLS, the structure of the PLS-DA model is based on variables X and y. However, unlike PLS, the y of PLS-DA

is not fixed. Each column of the Y matrix corresponds to a class and contains 1 if the sample belongs to that class and 0 otherwise (full discriminate coding). On the other hand, if the class is heterogeneous, it will be difficult to model it, since all samples of the class are given the same value (1), but they differ in some way in terms of experience.<sup>[5]</sup>

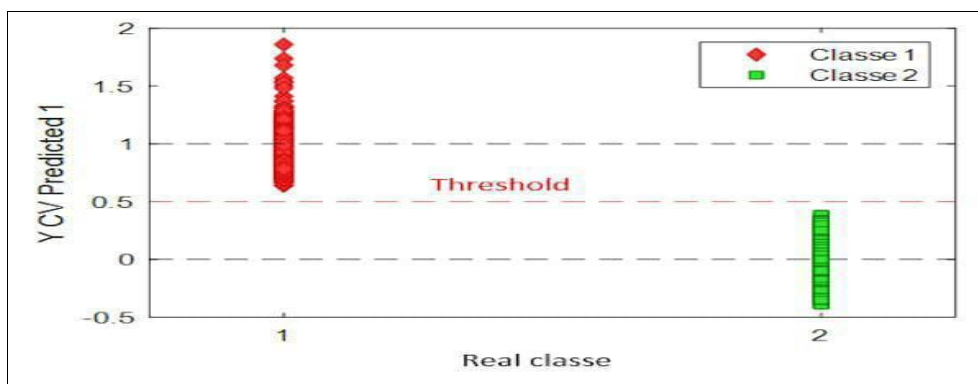


Fig 4

#### SVM (Support vector machine)

This method is often used for nonlinear or complex problems. It depends on finding the boundary separating the two groups. Therefore, only part of the calibration model is actually used: these are the support vectors defining the boundaries. In calibration, the dimensions of this matrix are  $N \times N$  ( $n$  is the number of samples in the calibration). The kernel is usually a Gaussian kernel, and the Gaussian requires a large optimization parameter: sigma, which allows the degree of nonlinearity to be adjusted. The SVM method also needs the optimization of the constant to avoid competition ( $C$  or value). Tuning these two parameters is important to obtain a good and robust model.

#### SIMCA-Soft independent modeling of class analogy

This suitable for high-performance products. Each group  $k$  was modeled separately by PCA. This PCA allows modeling of class differences. A confidence interval is then created for each sample to define the limits of association for that group. This bound may be based on the Euclidean distance of the X residues (called  $Q$ ), the leverage (or equivalent Hotelling  $T^2$  or Mahalanobis distance), or generally a combination of both. If the sample falls into the restricted category, it is divided into  $k$  groups. If multiple classes overlap or are close to each other, the model may be assigned to more than one class or no class at all. In this case, consider using a rejection group. In fact, the PCA model is created by class, independent of other classes. On the other hand, when the signals are very close, the method based on class difference will be better.<sup>[6]</sup>

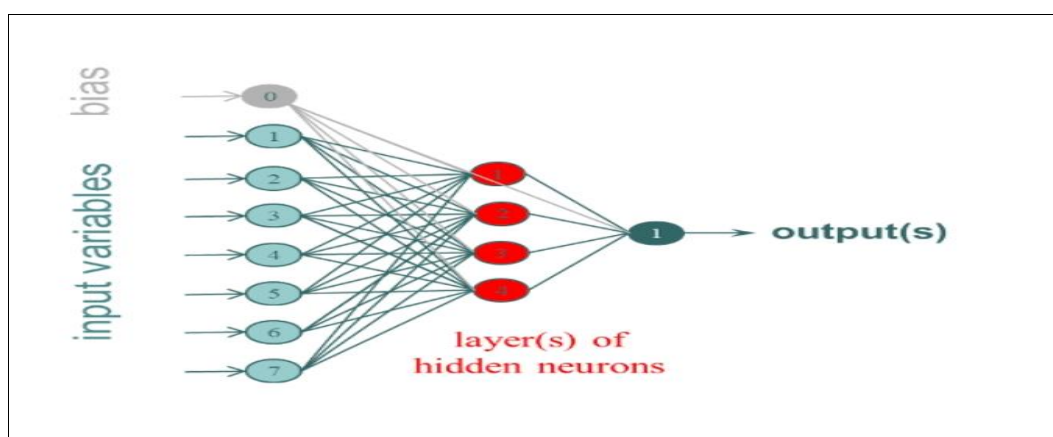


Fig 5

#### Unsupervised classification

The first is the "unsupervised" method (or group), which aims to repeat similar patterns without using prior knowledge. The second is the "supervised classification" approach (or discrimination), which uses class members to build a model. Mining process. They are search tools designed to find "natural" clustering trends from internal data (X) without any prior knowledge of the class structure. Therefore, all these methods only evaluate the similarity of examples according to their X values.<sup>[7]</sup>

#### K-Means

Non-hierarchical clustering methods focus on creating the final partitions of the data. Unlike hierarchical methods, users are required to specify based on their prior knowledge, which can be a strong limit on ideas.

Distribution of  $k$  categories. K-means results depend on the choice of the initial distribution and  $k$  classes.

Do this process until the end of the process is reached (i.e.: no change or maximum return is reached).<sup>[8]</sup>

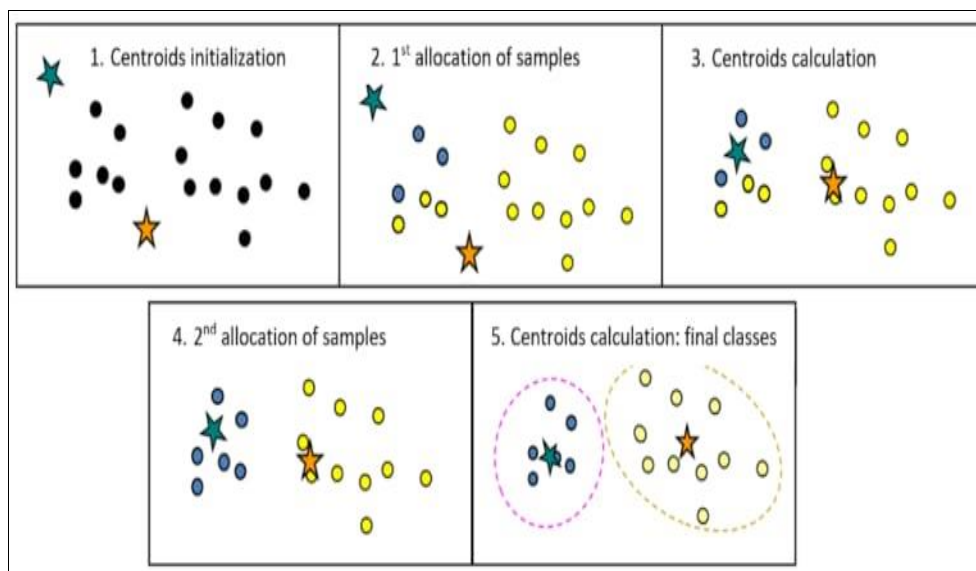


Fig 6

### Principal component analysis

It is the main function of the chemo metrics tool. The PCA method can be used for the following purposes:

- Visualization of X in different positions
- Cluster detection
- Outlier detection
- Impact of anomalies
- Data compression while reducing
- Noise removal<sup>[9]</sup>

PCA can be seen as a better method for visualizing patterns represented by different variables than the preparation of the principles in the new system called principal component (PC). These axes are designed to be different from X to

extract data. The elements of X, while the last component represents the noise. To facilitate, the structure is usually seen in a 2D or 3D plane, which corresponds to the estimation of the structure of the layers in 2 or 3 axes. PCA can form the basis of other different methods such as unsupervised distribution.

### Multi-blocks Analysis

Many blocks of data contain data sets where the same structure has different properties or where many models, each with a different number of items, have different properties.

The patterns in block may be different. Multiple group analysis

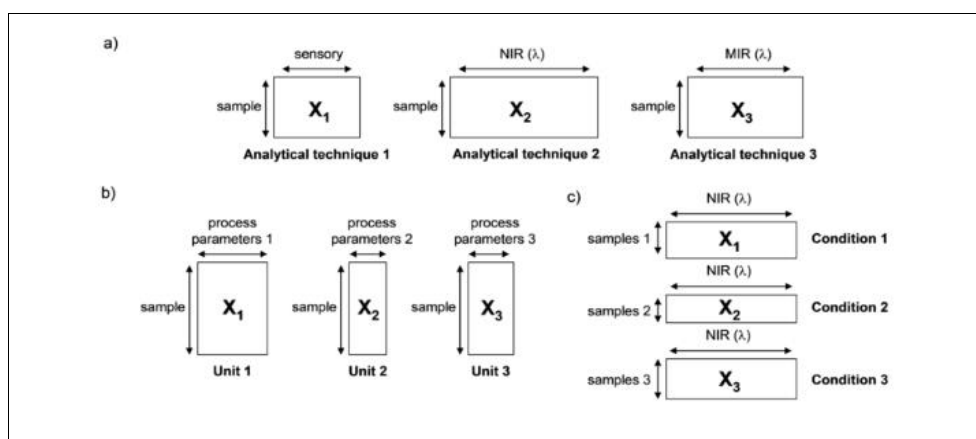


Fig 7

### Independent Component Analysis

ICA is designed to identify the components and conditions present in a mixture or process. The main component of PCA usually identifies a mixture of these pure components and does not necessarily provide a sufficient answer. The weighting coefficient or "ratio" A is proportional to the contribution of the source signal in the mixture. Unlike PCA, the results depend on the number of components to be extracted. So for example, the first component for a 3-component ICA will be different from a 4-component ICA. There are tools that make it easy to choose the number of components, such as "ICA by parts", which makes it

possible to analyze the strength of the model by looking at the correlation between the components of the product, ICA generates the data in parts<sup>[10]</sup>

### Regression algorithms

#### Linear

Multiple linear regression (MLR) MLR is the simplest multivariate regression model. It is a simple linear regression for multivariate factors.<sup>[11-15]</sup>

However, if the explanatory variables are in the same direction, the so-called inverse matrix calculation can make the model unstable.

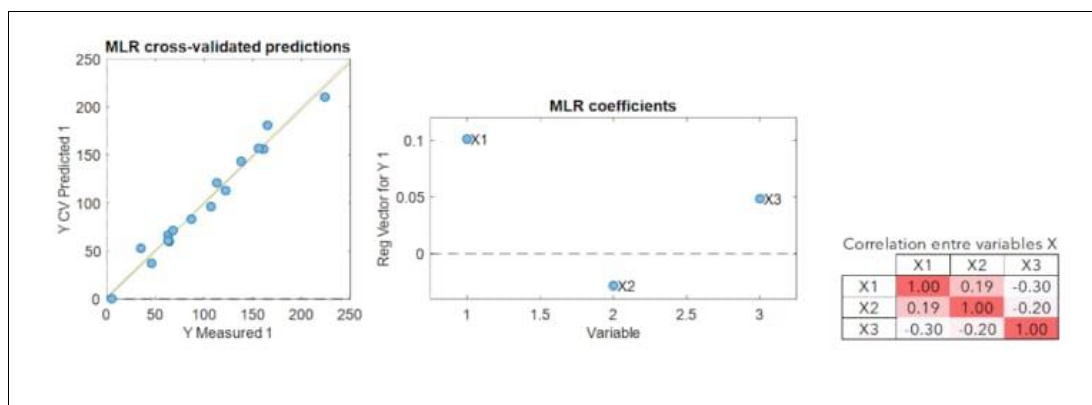


Fig 8

### Non-linear

#### SVM – support vector machines

For nonlinear or complex problems, the SVM (Support Vector Machine) method is used. This depends on finding the boundary separating the two groups. Therefore, in fact, only model is used: these are the support vectors that define the boundaries. In calibration, the dimensions of this matrix

The kernel is usually a optimization parameter of Gaussian diameter: sigma, which allows the degree of nonlinearity to be varied.

Using the method should also optimize a fixed parameter, such that over fitting (C or value) can be avoided. Tuning is important to obtain a good and powerful model. SVM.

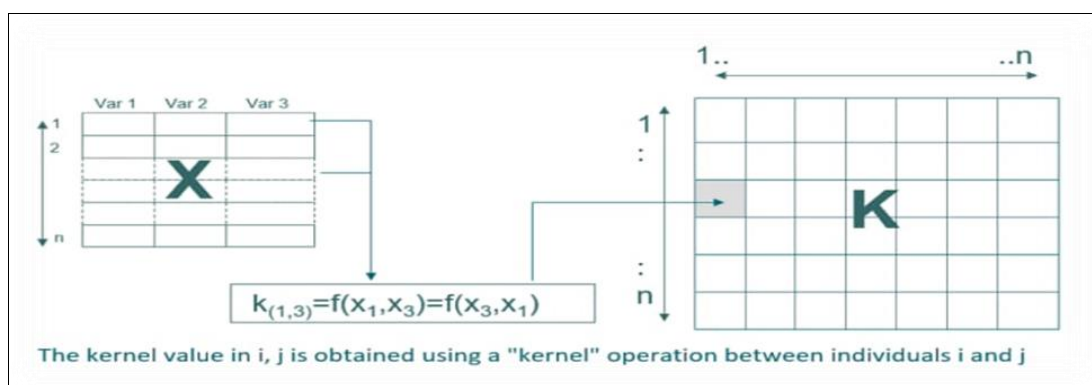


Fig 9

### Random forests

A development of this method, the "Random Forest", makes it possible to overcome the overfitting problem in the CART approach. When a new model is sent to the forest, its final prediction corresponds to:

- In the case of multiple predictions, the average of all prediction trees

### Most rooms

based on various fundamentals The process is designed to connect the two models, to be "integrated", of which the RF is also part. Like CART) is used with a split or random variable.

### Chemometrics in forensic science: approaches and applications

Forensic investigations often rely on physical evidence to reconstruct the circumstances of a crime. However, more objective interpretation of evidence, as well as rigorous management, storage, and evaluation procedures, are still needed.

Chemometrics is recognized as a powerful tool for the interpretation and optimization of analytical procedures in forensic science. However, factors such as sampling, validity, and underlying study design must be carefully considered. This review begins with an overview of selected

chemometric methods and then provides a comprehensive review of studies demonstrating the utility of chemometrics in a variety of legal settings.

The review concludes with a discussion of the problems and consequences of this sudden change.

### Chemometric Approaches in UV-visible Spectroscopy

Chemometrics is a science that involves measuring the effects of chemical systems (including chemical processes) on the state of the system using mathematical or statistical algorithms. It is clear from this definition that chemometrics is data-driven. The goal of many chemometric methods is to construct empirical models from data that allow one or more physical properties to be predicted based on the measurements. The four main aspects of performance that can be improved using chemometric methods are accuracy, precision, robustness, and reproducibility. multidisciplinary analysis. The calculation of compounds with many overlaps in the mixture is always a difficult analytical problem, especially in the non-equilibrium phase of the analyte. With the emergence of the fast scanning PDAs mentioned above and the low-end computers that can process complex data, new horizons have opened up in the mathematical processing of the acquired data. In recent years, two powerful signal processing methods, namely multi-l



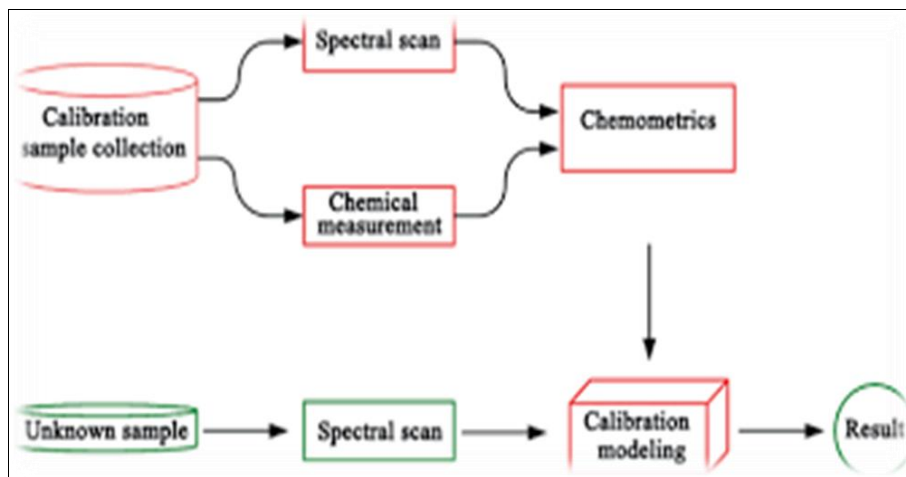


Fig 10

### Chemometric Methods for Spectroscopy-Based Pharmaceutical Analysis

Spectroscopy is widely used to characterize pharmaceutical products or processes, especially due to its desirable characteristics of being rapid, cheap, non-invasive/non-destructive and applicable both off-line and in-/at-/on-line. Spectroscopic techniques produce profiles containing a high amount of information, which can profitably be exploited through the use of multivariate mathematic and statistic (chemometric) techniques.

The present paper aims at providing a brief overview of the different chemometric approaches applicable in the context of spectroscopy-based pharmaceutical analysis, discussing both the unsupervised exploration of the collected data and the possibility of building predictive models for both quantitative (calibration) and qualitative (classification) responses.

### Future of Chemometric Based Spectral Data Analysis

Deep learning methods and machine learning algorithms are expected to take the lead. This might result in the development of “smart” chemometric techniques that have the ability to suggest new analytical approaches or even select the most appropriate analysis method automatically for a given dataset.

Developing explainable AI models in chemometrics will probably be a major focus as AI integration becomes more complex. Creating trust and ensuring the reliability of chemometric analysis, especially for important applications, will rely largely on the transparency with which models arrive at their conclusions.

Emerging spectroscopic techniques: Chemometric analysis will have expanded opportunities as a result of the creation of new spectroscopic techniques with greater sensitivity and resolution. Chemometrics will have to change and grow to deal with the special data structures and complex connections that come with these cutting-edge methods.

Wider uses: Chemometrics is expected to be applied in novel and developing areas outside of traditional chemistry. It could be applied to environmental monitoring, improved material characterization, or even customized treatment using spectrum analysis of biological samples.

To sum up, there can be lots of promise in the investigation of chemometric spectrum data in the future. Chemometrics can be an effective instrument for unlocking insights hidden within spectral data, resulting in innovations and

developments across many scientific and technical areas. This can be achieved by accepting these advancements and developing interdisciplinary collaboration.

### Applications

Chemometrics refers to the application of mathematical and statistical methods to chemical data. It plays a crucial role in various scientific disciplines, particularly in chemistry, biochemistry, pharmacology, and environmental science. Here are some detailed applications of chemometrics.

- **Quality Control and Assurance:** Chemometrics is extensively used in industries such as pharmaceuticals, food production, and manufacturing to ensure the quality and consistency of products. Techniques like multivariate data analysis (MVA) help in monitoring and controlling production processes to maintain desired quality standards.
- **Spectroscopic Analysis:** Spectroscopic techniques like infrared (IR), nuclear magnetic resonance (NMR), and mass spectrometry generate complex data sets. Chemometrics helps in analyzing and interpreting these spectra to identify compounds, quantify substances, and understand molecular structures.
- **Pattern Recognition:** Chemometrics methods such as principal component analysis (PCA), cluster analysis, and discriminant analysis are used for pattern recognition in chemical and biological data. These techniques can classify samples into groups based on similarities or differences in their chemical composition.
- **Quantitative Structure-Activity Relationship (QSAR):** QSAR models in drug discovery and environmental science predict biological activity or properties of chemical compounds based on their molecular structure. Chemometrics assists in building and validating these models using statistical regression and machine learning algorithms.
- **Process Analytical Technology (PAT):** In manufacturing processes, PAT involves real-time monitoring and control to ensure consistent product quality and process efficiency. Chemometrics tools enable the analysis of large datasets from sensors and analytical instruments to optimize manufacturing processes.

- **Environmental Monitoring:** Chemometrics is applied in environmental science to analyze data from water and air quality monitoring systems. It helps in identifying pollutants, assessing their concentrations, and understanding their sources using techniques like chemometric fingerprinting and source apportionment.
- **Metabolomics and Proteomics:** In biological research, chemometrics aids in analyzing large-scale metabolomic and proteomic data. It helps in identifying biomarkers, understanding metabolic pathways, and studying the interactions between genes, proteins, and metabolites.

### Advantages

Chemometrics offers several distinct advantages in the field of chemistry and related disciplines where complex data analysis is crucial. Here are the key advantages of using chemometrics in detail:

- **Multivariate Analysis Capability:** One of the primary strengths of chemometrics lies in its ability to handle multivariate datasets, where multiple variables (such as different wavelengths in spectroscopy or different components in a mixture) are measured simultaneously. Techniques like principal component analysis (PCA), partial least squares regression (PLS), and discriminant analysis can extract relevant information from these datasets, revealing patterns, correlations, and differences that might not be apparent through univariate methods.
- **Quality Improvement and Control:** In industrial applications, chemometrics plays a critical role in ensuring product quality and process efficiency. By analyzing data from production processes (such as chemical reactions or manufacturing steps), chemometrics helps in monitoring and controlling variables to maintain desired product specifications and minimize variability.
- **Enhanced Data Interpretation:** Chemometrics provides tools for robust data interpretation. By reducing data dimensionality and focusing on essential variations, chemometrics helps in understanding complex relationships within data. This capability is particularly useful in fields like spectroscopy, chromatography, and biological assays where numerous variables interact to produce measurable outcomes.
- **Optimization of Experimental Design:** Chemometrics aids in designing experiments more effectively. Techniques such as design of experiments (DOE) optimize the allocation of resources and samples, ensuring that experiments are efficient and yield maximum information. This approach minimizes costs and reduces the number of experiments required to achieve meaningful results.
- **Exploratory Data Analysis:** Chemometrics supports exploratory data analysis, where researchers can visually and statistically explore datasets to uncover trends, outliers, and relationships. This process aids in generating hypotheses, identifying new patterns, and gaining insights into underlying mechanisms within chemical and biological systems.

- **Quality Assurance in Analytical Methods:** Chemometrics helps in validating and improving analytical methods by assessing factors such as accuracy, precision, linearity, and sensitivity. By analyzing calibration curves and evaluating method performance parameters, chemometrics ensures that analytical measurements meet regulatory standards and produce reliable results.
- **Automation and Efficiency:** With advancements in computational methods and software tools, chemometrics has become increasingly automated, making complex data analysis more accessible and efficient. Automated data preprocessing, model development, and validation streamline workflows and reduce the time and effort required for data analysis.
- **Interdisciplinary Applications:** Chemometrics bridges disciplines such as chemistry, biology, pharmacology, environmental science, and engineering. Its methods and principles are widely applicable across diverse fields, fostering interdisciplinary collaboration and enabling researchers to tackle complex scientific challenges from multiple perspectives.
- **Integration of Diverse Data Sources:** Chemometrics facilitates the integration of data from different analytical techniques and instruments. By combining data from spectroscopy, chromatography, mass spectrometry, and other methods, chemometrics provides a comprehensive view of chemical systems, allowing researchers to uncover complex interactions and dependencies.

### Limitations

While chemometric techniques offer numerous advantages, they also come with certain limitations and challenges. Here are some key limitations of chemometric techniques:

- **Data Quality Dependency:** The reliability and accuracy of chemometric analyses heavily depend on the quality of the input data. If the data are noisy, contain outliers, or are biased, it can lead to erroneous conclusions and unreliable models. Preprocessing steps such as data cleaning and outlier detection are critical but can be challenging, especially with complex datasets.
- **Overfitting and Model Complexity:** In some cases, chemometric models can become overly complex, especially when dealing with high-dimensional data or when using powerful modeling techniques. Overfitting occurs when a model captures noise or random fluctuations in the data, leading to poor generalization and unreliable predictions on new data.
- **Interpretability Issues:** Some chemometric models, particularly those based on machine learning algorithms such as neural networks or support vector machines, can be difficult to interpret. This lack of transparency makes it challenging to understand the underlying relationships between variables and to explain the reasoning behind model predictions.

- **Assumption Violations:** Many chemometric techniques are based on statistical assumptions that may not always hold true for real-world data. For example, assumptions of normality, linearity, or homoscedasticity (constant variance) may be violated, affecting the validity of statistical tests and model results.
- **Need for Expertise:** Effective application of chemometric techniques often requires a good understanding of both the underlying statistical principles and the specific domain of application (e.g., chemistry, biology, engineering). Interpretation of results and selection of appropriate models can be challenging without sufficient expertise.
- **Data Preprocessing Complexity:** Preprocessing of data (e.g., normalization, scaling, feature selection) is often necessary before applying chemometric techniques to ensure optimal model performance. However, determining the most suitable preprocessing steps can be time-consuming and requires careful consideration of data characteristics and analytical goals.
- **Ethical and Regulatory Considerations:** In fields such as pharmaceuticals and environmental science, chemometric models may be used to make critical decisions with ethical and regulatory implications. Ensuring transparency, fairness, and accountability in model development and decision-making processes is essential but challenging.

### Conclusion

Chemometrics stands as a pivotal discipline in modern scientific research and industrial applications, offering powerful tools for extracting valuable insights from complex chemical and biological data. Through sophisticated mathematical and statistical techniques, chemometrics addresses a wide range of challenges and facilitates informed decision-making across diverse fields.

- variables simultaneously, chemometrics uncovers hidden patterns, correlations, and dependencies that traditional methods might miss.
- **Enhanced Data Interpretation:** By reducing data dimensionality and focusing on significant
- **Multivariate Data Analysis:** Capable of handling large datasets with multiple variations, chemometrics enables clearer understanding and deeper insights into complex systems.
- **Predictive Modeling:** Through advanced statistical methods and machine learning algorithms, chemometrics develops predictive models for applications ranging from drug discovery to environmental monitoring, improving efficiency and effectiveness.
- **Quality Control and Assurance:** In industrial settings, chemometrics ensures product quality and process optimization by monitoring and controlling variables, minimizing variability, and maintaining consistent standards.

- Despite its strengths, chemometrics faces challenges such as data quality dependencies, model complexity, interpretability issues, and computational intensity. Addressing these challenges requires expertise, careful validation, and consideration of ethical implications, particularly in decision-making contexts.

### References

1. Pierce KM, Hoggard JC, Mohler RE, Synovec RE, Chromatogr J, 2008:A1184:341.
2. Brereton RG, Chemometrics, Data Analysis for the Laboratory and Chemical Plant, Wiley, New York, 2003.
3. Massart DL, Chemometrics: A Textbook, Elsevier Sciences Ltd., New York, 1988.
4. Beebe KR, Pell RJ, Seasholtz MB, Chemometrics: A Practical Guide, WileyInterscience, New York, 1998.
5. Amigo JM, Skov T, Bro R, Chem. Rev, 2010:110:4582.
6. Lavine B, Workman J, Anal. Chem, 2010:82:4699.
7. Pierce KM, Mohler RE, Sep. Purif. Rev, 2012:41:143.
8. Cortes HJ, Winniford B, Luong J, Pursch M, Sep J. Sci, 2009:32:883.
9. Mondello L, Tranchida PQ, Dugo P, Dugo G. Mass Spectrom Rev, 2008:27:101.
10. François I, Sandra K, Sandra P. Anal Chim Acta, 2009:641:14.
11. Horvatovich P, Hoekman B, Govorukhina N, Bischoff R. J Sep Sci, 2010:33:1421.
12. Reichenbach SE, Tian X, Tao Q, Stoll DR, Carr PW. J Sep Sci, 2010:33:1365.
13. Bailey HP, Rutan SC, Carr PW. J Chromatogr A, 2011:1218:8411.
14. Mondello L, Herrero M, Kumm T, Dugo P, Cortes H, Dugo G. Anal Chem, 2008:80:5418.
15. Hoggard JC, Synovec RE. Anal Chem, 2007:79:1611.