**Shiraz University**
**Faculty of Sciences**

**M.Sc. Thesis in Analytical Chemistry**

**APPLICATION OF CHEMOMETRICS METHODS IN STUDYING THE EFFECTS OF SOLUTE STRUCTURES ON THE ELECTROCHEMISTRY OF STEROIDS AND ELECTRONIC ABSORPTION SPECTRA OF SUBSTITUTED SCHIFF BASES**

**By**

**MAHDIEH YAZDANI**

**Supervised by:**

**B. HEMMATEENEJAD Ph.D.**

**September 2009**

In the Name of GOD

Declaration

I, Mahdieh Yazdani, a chemistry student, majored in analytical chemistry from the faculty of science declare that this thesis is the result of my research and I had written the exact reference and full indication wherever I used other's sources. I also declare that the research and the topic of my thesis are not reduplicative and guarantee that I will not disseminate its accomplishments and not make them accessible to others without the permission of the university. According to the regulations of the mental and spiritual ownership, all rights of this belong to Shiraz University.

Name: Mahdieh Yazadni

Date: 2009.9.5

IN THE NAME OF GOD

# APPLICATION OF CHEMOMETRICS METHODS IN STUDYING THE EFFECTS OF SOLUTE STRUCTURES ON THE ELECTROCHEMISTRY OF STEROIDS AND ELECTRONIC ABSORPTION SPECTRA OF SUBSTITUTED SCHIFF BASES

## BY

MAHDIEH YAZDANI

## THESIS

SUBMITTED TO THE SCHOOL OF GRADUATE STUDIES IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE DEGREE OF MASTER OF SCIENCE (M.Sc.)

## IN

ANALYTICAL CHEMISTRY
SHIRAZ UNIVERSITY
SHIRAZ
ISLAMIC REPUBLIC OF IRAN

EVALUATED AND APPROVED BY THE THESIS COMMITTEE AS: EXCELLENT

B. HEMMATEENEJAD, Ph.D., ASSOC. PROF. OF CHEMISTRY
(CHAIRMAN)

A. SAFAVI, Ph.D., PROF. OF CHEMISTRY

H. SHARGHI, Ph.D., PROF. OF CHEMISTRY

G. ABSALAN, Ph.D., ASSOC. PROF OF CHEMISTRY

SEPTEMBER 2009

*Dedicated to:*

*My Parents*

*&*

*My Siblings*

*&*

*In memory of:*

*My Grandfather*

# ACKNOWLEDGMENTS

# ABSTRACT

Application of chemometrics methods in studying the effects of solute structures on the electrochemistry of steroids and electronic absorption spectra of substituted schiff bases

BY

MAHDIEH YAZDANI

A data set consists of a diverse set of steroids have been used to develop QSER models for their $E_{1/2}$ by means of MLR and PCR analyses. In MLR analysis, the QSER models were constructed by first grouping descriptors and then stepwise selection of variables from each group (MLR1) and stepwise selection of predictor variables from the pool of all calculated descriptors (MLR2). Similar procedure was used in PCR analysis. In overall, the $R^2$ of CV and prediction of the QSER models resulted from MLR1, MLR2 and PCR1 approaches were higher than 90%, which show the high ability of the models to predict $E_{1/2}$ of the steroids.

A series of Schiff bases were studied for their delicate changes in absorption electronic spectra by changing substituents and solvents. Schiff base derivatives of different substituents in the aromatic ring were used. At the same time, different solvents having different solvatochromic parameters were used. Linear relationships were established to show how the variations of $\nu_{max}$ were related to molecular descriptors and solvatochromic parameters by changing substituents and solvent, respectively. MLR and factor analysis (FA) were applied to find the meaningful chemical factors and provide the regression models. It was found that the $\nu_{max}$ was mainly controlled by the solvent's polarity and the descriptors of Mor30p, Mor22v and E2v were selected for *para*, *meta* and *ortho* positions, respectively.

**TABLE OF CONTENT**

**CONTENT**                                                                                    **PAGE**

# LIST OF FIGURES

VIII

# GLOSSARY OF ABBREVIATIONS:

| | |
|---|---|
| AET | Apparent Error in Target |
| CV | Cross Validation |
| $E_{1/2}$ | Half-wave reduction potential |
| FA | Factor Analysis |
| FE | Feature Extraction |
| FS | Feature Selection |
| GA | Genetic Algorithm |
| HPLC | High Performance Liquid Chromatography |
| HPLC-DAD | HPLC-Diode Array Detector |
| LFER | Linear Free Energy Relationships |
| LMO | Leave Many Out |
| LOO | Leave One Out |
| MLR | Multiple Linear Regression |
| NIPALS | Nonlinear Iterative Partial Least Squares |
| NIR | Near Infrared |
| PCA | Principal Component Analysis |
| PCR | Principal Component Regression |
| PCs | Principal Components |
| PLS | Partial Least Squares |
| QSAR | Quantitative Structure Activity Relationship |
| QSER | Quantitative Structure Electrochemistry Relationship |
| QSPR | Quantitative Structure Property Relationship |
| $RMSE_{cv}$ | Root Mean Squared Error of Cross Validation |
| $RMSE_{p}$ | Root Mean Squared Error of Prediction |
| SARs | Structure Activity Relationships |
| SSPE | Sum of Squares of Prediction Errors |
| Std | Standard Deviation |
| SVD | Singular Value Decomposition |
| TFA | Target Factor Analysis |

# CHAPTER ONE

# INTRODUCTION

## 1.1 Introduction

The field of chemistry is currently facing major changes. Optical, mechanical, and microelectronic technologies have advanced rapidly in recent years. Computer power has increased dramatically as well. All these developments, together with other factors, provide a new opportunity but also challenge to chemists in research and development [1].

Sound chemical information that forms the basis of many of humanity's important decision-making processes depends on three critical properties of the measurement process, including its (1) chemical properties, (2) physical properties, and (3) statistical properties. The conditions that support sound chemical measurements are like a platform supported by three legs. Credible information can be provided only in an environment that permits a *thorough understanding and control* of these three critical properties of a chemical measurement:

1. Chemical properties, including stoichiometry, mass balance, chemical equilibria, kinetics, etc.

2. Physical properties, including temperature, energy transfer, phase transitions, etc.

3. Statistical properties, including sources of errors in the measurement process, control of interfering factors, calibration of response signals, modeling of complex multivariate signals, etc.

If any one of these three legs is missing or absent, the platform will be unstable and the measurement system will fail to provide reliable results, sometimes with catastrophic consequences. It is the role of statistics and chemometrics to address the third critical property. It is this fundamental role that provides the primary motivation for developments in the field of chemometrics. Sound chemometric methods and a well-trained work force are necessary for providing reliable chemical information for humanity's decision-making activities [2, 3].

2

## 1.2 Chemometrics

The term chemometrics was first coined in 1971 to describe the growing use of mathematical models, statistical principles, and other logic-based methods in the field of chemistry and, in particular, the field of analytical chemistry. Chemometrics is an interdisciplinary field that involves multivariate statistics, mathematical modeling, computer science, and analytical chemistry. Some major application areas of chemometrics include (1) calibration, validation, and significance testing; (2) optimization of chemical measurements and experimental procedures; and (3) the extraction of the maximum of chemical information from analytical data [2, 4]. A reasonable definition of chemometrics remains as how do we get chemical relevant information out of measured chemical data, how do we represent and display this information, and how do we get such information into data? As mentioned by Wold [5]. Chemometrics is considered by some chemists to be a subdiscipline that provides the basic theory and methodology for modern analytical chemistry [1].

In many respects, the field of chemometrics is the child of statistics, computers, and the "information age." Rapid technological advances, especially in the area of computerized instruments for analytical chemistry, have enabled and necessitated phenomenal growth in the field of chemometrics over the past 30 years. The relationship of chemometrics to different disciplines is indicated in Figure 1.1. On the left are the enabling sciences, mainly quite mathematical and not laboratory based. Statistics, of course, plays a major role in chemometrics. Statistical approaches are based on mathematical theory, so statistics falls between mathematics and chemometrics. Computing is important as much of chemometrics relies on software. Engineers, especially chemical and process engineers, have an important need for chemometric methods in many areas of their work, and

3

**Figure 1.1:** How chemometrics relates to other disciplines

have a quite different perspective from the mainstream chemist. On the right are the main disciplines of chemistry that benefit from chemometrics. Analytical chemistry is probably the most significant area. Environmental chemists, biologists, food chemists as well as geochemists, chemical archaeologists, forensic scientists and so on depend on good analytical chemistry measurements and many routinely use multivariate approaches especially for pattern recognition, and so need chemometrics to help interpret their data. These scientists tend to identify with analytical chemists. The organic chemist has a somewhat different need for chemometrics, primarily in the areas of experimental design (e.g. optimising reaction conditions) and QSAR (quantitative structure–analysis relationships) for drug design. Finally, physical chemists such as spectroscopists, kineticists and materials scientists often come across methods for signal deconvolution and multivariate data analysis [6].

To finalize the point, it should be said that chemometrics is the interface between chemistry and mathematics. As Kowalski [7-10] clearly stated "Chemometric tools are vehicles that can aid chemists to move more

4

PDF created with pdfFactory Pro trial version www.pdffactory.com

efficiently on the path from measurements to information to knowledge"[11].

## 1.3  Historical Development of Chemometrics

The term *chemometrics* was introduced by Svante Wold [5] and Bruce R. Kowalski in the early 1970s. Terms like *biometrics* and *econometrics* were also introduced into the fields of biological science and economics. Afterward, the International Chemometrics Society was established. Since then, chemometrics has been developing and is now widely applied to different fields of chemistry, especially analytical chemistry.

### 1.3.1  Development of Chemometrics

Chemometrics has developed over the past decade from a fairly theoretical subject to one that is applied in a wide range of sciences. Early development went hand in hand with the development of scientific computing, and primarily involved using multivariate statistical methods for the analysis of analytical chemistry data. Early chemometricians were likely to be FORTRAN programmers using mainframes and statistical libraries of subroutines. The earliest applications were primarily to analytical chemical datasets, often fairly simple in nature, for example a high performance liquid chromatography (HPLC) cluster of two or three peaks or a set of UV/visible mixture spectra.

Chemometrics as a discipline became organized in the 1980s, with the first journals, meetings, books, societies, packages and courses dedicated to the subject. Historically HPLC and near infrared (NIR) spectroscopy provided particularly important growth points in the 1980s, partly due to the economic driving forces at the time and partly due to the research interest and contacts of the pioneers. Industrial applications were

5

particularly important in the developing phase. In the 1990s the application of chemometrics started expanding, with special emphasis on certain industries especially the pharmaceutical industry, where large companies provided significant funding for developments. A new and exciting phase has emerged as from the late 1990s, involving very complex datasets. This new and exciting phase is possible due to the capacity of analytical instruments to acquire large amounts of data rapidly, for example via autosamplers, and chemometrics becomes a type of data mining. All the building blocks are available, for example signal processing, chromatographic alignment, data scaling and pattern recognition but the problems are now much more complex. In addition many of the newer applications are biologically driven and so there is emerging a new interface between chemometrics and bioinformatics, especially the relationship between analytical data as obtained by spectroscopy or chromatography and genetics, requiring some appreciation of gene or protein sequence and similarities.

Chemometrics became much more widespread today. Individual researchers in academic positions, some with prior training in some of the more established groups, and others who read papers and books that had been written over this early period, dabbled with chemometrics, and so the applications spread not just to dedicated research groups but to investigators working in different environments [12].

## 1.3.2 Chemometrics in Iran

Recently, Hemmateenejad has published a paper [13] to represent the activity of the Iranian chemometrics community and to increase the awareness about the studies of chemometrics in Iran. In this article a list of the publication of the Iranian scientists in the chemometrics was collected from the ISI web of science database. A rapid growth is observed in the

6