# Relationship between Duality Concept and Complementarity-Coupling Theorems in Case of Reducing the Rotational Ambiguity

Ph.D. Thesis

**Elnaz Tavakkoli**

**Supervisors: Prof. H. Abdollahi**

**Prof. R. Rajkó**

# Abstract

Self Modeling Curve Resolution (SMCR) is a class of techniques concerned with estimating pure profiles underlying a set of measurements on chemical systems. In general, the estimated profiles are ambiguous (non-unique) except if some special conditions fulfilled. Implementing the adequate information can reduce the so-called rotational ambiguity effectively, and in the most desirable cases lead to the unique solution. Therefore, studies on circumstances resulting in unique solution are of particular importance. In bilinear chemical (e.g., spectroscopic) data matrix, there is a natural duality between its row and column vector spaces using minimal constraints (non-negativity of concentrations and absorbances). The correspondence between the points in the considered space and the hyperplanes in the dual space can be applied to extract information in SMCR analysis. Infact, the conditions for achieving a unique solution can be investigated based on the duality concept as a general principle

## I. Duality based investigations in SMCR

## A. Clarification and visualization of the duality concept of curve resolution using three-component non-negative bilinear chemical data

In this work, it has been intended to clarify the duality concept and implement it in data structure visualization. Accordingly, duality principle is used to visualize the relationship between data points in each space and non-negativity borders in its dual space. Duality concept state that one point in one data space is corresponding with one hyperplane in the other space. Therefore, by implementing this concept and using data points, it has been visualize that the intersection of the all defined hyperplanes generates the outer boundaries in the dual mode. In order to reach a better understanding of this concept and its application in SMCR, hypothetical chromatographic system of three components is visualized in PC spaces.

## B. Duality based direct resolution of unique profiles using zero concentration region information

The duality relation can also be the proposed approach to exploit more information of the chemical data matrix if additional knowledge of the system is available. In this section, by applying the reliable information about the concentration window of components, the conditions of unique resolution has been explored. Therefore, the conditions of the unique solution according to duality concept and using zero concentration region information is intended to show using simulated and experimental datasets. It gives an easy way for finding how to use zero concentration region information based on duality in special cases to obtain unique solutions in SMCR methods. Additionally, the knowledge of a pure profile (e.g. spectrum) is used to discovery of information in SMCR method. In this way, graphical tools of PCA are utilized for illustrating the data structure.

## II .Development on generalized rank annihilation problem by implementing duality

The analytical chemist is frequently confronted with the problem of analyzing complex mixtures in the presence of any component in the sample that is not included in the calibration model. In these cases, it is desirable to be able to obtain quantitative information for a particular component without concern for the rest of the components in the sample. The property of quantitation of an analyte in the presence of unknown constituents is called second-order advantage. Additionally, in recent studies it has been shown that when the solution is unique or information is available for obtaining the unique solution, the duality concept,is a useful approach to extract information in SMCR. In a bilinear data matrix, there is a natural duality between its row and column spaces and it is possible to transform the coordinates in row space to the coordinates in column space and vice versa without using any constraint. In this section, a novel algorithm to achieve "second order advantage" is introduced based on duality. Moreover, informative geometrical visualizations beside mathematical formulas are presented for the proposed

method. In this way, one of the most common approaches for visualization of data, Principle Component Analysis (PCA), is used. It is shown in order to determine $\lambda$, the value related to the contribution of the constituent of interest in an unknown mixture, the subspace of the all interfering compounds should be defined properly. Therefore, a matrix of $\lambda$ values will be obtained in a systematic way and it is possible to calculate more precise estimation of $\lambda$.

## III. A rank reduction based normalization and visualization of rank annihilation

## A. A rank reduction based normalization

In this work a special normalization for a particular component in the system by using pure component spectra is derived. By implementing the general rank annihilation Wedderburn formula, it has been shown that the introduced normalization specifies the intensity of the same component in the dual mode. The interesting point is that the derived normalization is particularly applied for the analyte of interest without considering the remain part of the abstract space. By the help of the figures and plots, the relation between "known profile" and "normalization" is visualized. The introduced normalization result in the rank reduction for the known component.

## B. Visualization of rank annihilation

Although rank annihilation is a crucial concept in chemometrics, but it has not been investigated considering the dimensionality of the data. In this work, illustrative visualization of rank annihilation procedure has been provided. It has been illustrated that the reduction of the dimension of the residual, when standard matrix is subtracted from the mixture, can be monitored in rank annihilation problem. Therefore, in this section graphical visualizations are provided to depict the relation of "rank" and "dimensionality" when the standard component is annihilated from the data.

## IV. Soft-trilinear constraints for improved quantitation in Multivariate Curve Resolution

When a set of samples is to be analyzed with one data matrix per sample, the data is often presumed to have "trilinear" structure if the profile for each compound does not change shape or position from one sample to the other. By applying this information as a trilinearity constraint in SMCR methods, overlapping peaks related to the pure compounds of interest can be resolved in a unique way. In practice, many systems have non-trilinear behavior due to deviation from ideal response, for example a sample matrix effect, or changes in instrumental response (e.g., shifts and/or changes in the shape of chromatographic peaks). In such cases, the trilinear model is not valid because every analyte does not have the same peak shape or position in every sample. In such cases, the unique profiles obtained by strictly enforced trilinearity constraints will not necessarily produce true profiles because the data set does not follow the assumed trilinear behavior.

In this section, we introduce "soft-trilinearity constraints" and a new MATLAB program to permit peak profiles of the components of interest to have small deviations in their shape and position from sample to sample. In order to visualize the results, soft-trilinearity constraints were incorporated into a systematic grid search algorithm for the case of a three-component system [1]. This algorithm is general and can be applied to any MCR method. Results are provided for simulated noisy data with non-trilinear behavior and one experimental data set. The results show that implementing soft-trilinearity constraints reduces the range of feasible solutions considerably compared to the application of simpler non-negativity constraints. The results of this approach are compared to other methods including PARAFAC2 and MCR-ALS with hard-trilinearity constraints. It is shown that the methods employing hard-trilinearity constraints lead to incorrect solutions, or produce solutions outside the range feasible solutions.

# Contents

VI

# Table of Figures:

# List of Schematics

# List of tables