

Spectroscopy Sensitivity Study by Linear Programmin

by

Fei Chen

B.Sc., University of Victoria, 2017

A Dissertation Submitted in Partial Fulfillment of the  
Requirements for the Degree of

MASTER OF SCIENCE

in the Department of Computer Science

© Graduate Advisor, 2017  
University of Victoria

All rights reserved. This dissertation may not be reproduced in whole or in part, by  
photocopying or other means, without the permission of the author.

Spectroscopy Sensitivity Study by Linear Programmin

by

Fei Chen

B.Sc., University of Victoria, 2017

Supervisory Committee

---

Dr. Ulrike Stege, Co-Supervisor  
(Department of Computer Science)

---

Dr. Dennis Hore, Co-Supervisor  
(Department of Chemistry)

## Supervisory Committee

---

Dr. Ulrike Stege, Co-Supervisor  
(Department of Computer Science)

---

Dr. Dennis Hore, Co-Supervisor  
(Department of Chemistry)

## ABSTRACT

This document is a possible Latex framework for a thesis or dissertation at UVic. It should work in the Windows, Mac and Unix environments. The content is based on the experience of one supervisor and graduate advisor. It explains the organization that can help write a thesis, especially in a scientific environment where the research contains experimental results as well. There is no claim that this is the *best* or *only* way to structure such a document. Yet in the majority of cases it serves extremely well as a sound basis which can be customized according to the requirements of the members of the supervisory committee and the topic of research. Additionally some examples on using L<sup>A</sup>T<sub>E</sub>X are included as a bonus for beginners.

# List of Tables

Table 1.1	Sample input of the diet problem . . . . .	5
Table 3.1	Experiment 1 and 2 setting using toy model . . . . .	21
Table 3.2	Experiment 3 setting of toy model . . . . .	23
Table 3.3	Experiment 4 and 5 setting of toy model . . . . .	26
Table 3.4	Constraint study based on Experiment 4 . . . . .	27
Table 3.5	Constraint study based on Experiment 5 of toy model . . . . .	28
Table 4.1	Experiment 1 to Experiment 4 setting for methionine candidates	31
Table 4.2	Experiment 5 to Experiment 9 setting for methionine candidates	33
Table 4.3	Experiment 10 to Experiment 16 setting for methionine candidates	35
Table 4.4	Experiment 17 and 18 to explain the limitation of our LP model for methionine molecule . . . . .	36
Table 5.1	Detailed experiment set setting for the mixture of amino acids .	45

# List of Figures

Figure 2.1	The Euler angles represented as the spherical polar angles $\theta$ , $\phi$ and $\psi$ , and the illustration of the three successive rotations that transform the lab $x$ , $y$ , $z$ coordinate system into the molecular $a$ , $b$ , $c$ frame intrinsically and extrinsically [8]. . . . .	10
Figure 2.2	IR $x$ -polarized spectra of methionine's four candidates and target. The candidates are with $\theta$ of $0^\circ$ , $20^\circ$ , $40^\circ$ and $60^\circ$ . The target is produced by combining 10% of <i>candidate_ir_x_0</i> , 50% <i>candidate_ir_x_20</i> and 40% <i>candidate_ir_x_40</i> . . . . .	13
Figure 2.3	IR $z$ -polarized spectra of methionine's four candidates and target. The candidates are with $\theta$ of $0^\circ$ , $20^\circ$ , $40^\circ$ and $60^\circ$ . The target is produced by combining 10% of <i>candidate_ir_x_0</i> , 50% <i>candidate_ir_x_20</i> and 40% <i>candidate_ir_x_40</i> . . . . .	14
Figure 2.4	Raman $xx$ -polarized spectra of methionine's four candidates and target. The candidates are with $\theta$ of $0^\circ$ , $20^\circ$ , $40^\circ$ and $60^\circ$ . The target is produced by combining 10% of <i>candidate_ir_x_0</i> , 50% <i>candidate_ir_x_20</i> and 40% <i>candidate_ir_x_40</i> . . . . .	14
Figure 2.5	SFG $yyz$ -polarized spectra of methionine's four candidates and target. The candidates are with $\theta$ of $0^\circ$ , $20^\circ$ , $40^\circ$ and $60^\circ$ . The target is produced by combining 10% of <i>candidate_ir_x_0</i> , 50% <i>candidate_ir_x_20</i> and 40% <i>candidate_ir_x_40</i> . . . . .	15
Figure 3.1	<i>cosine</i> - polarize IR spectra of toy model candidates . . . . .	17
Figure 3.2	Toy model Experiment 2 resulting <i>cosine</i> -polarized IR spectrum plotted with the target spectrum; and the residual plot between the spectra. . . . .	22
Figure 3.3	Toy model Experiment 3 resulting <i>cosine</i> -polarized IR spectrum plotted with target spectrum; and the residual plot between the two spectra . . . . .	24

Figure 3.4 <i>sine</i> -polarized IR spectra of toy model candidates with $\theta$ value expanded from $0^\circ$ to $90^\circ$ . . . . .	25
Figure 3.5 IR spectra plotted by the return compositions from the constraint study based on Experiment 4 of toy model . . . . .	28
Figure 3.6 IR spectra plotted by the return compositions from the constraint study based on Experiment 5 of toy model . . . . .	29
Figure 4.1 Compare target IR spectra with the ones generated by the return composition of Experiment 1, 2 and 3 . . . . .	32
Figure 4.2 IR spectra plotted by using target composition and return composition of Experiment 17 . . . . .	37
Figure 4.3 Raman spectra plotted by using the target composition and the return composition of Experiment 17 . . . . .	38
Figure 4.4 SFG spectra plotted by using the target composition and the return composition of Experiment 17 . . . . .	39
Figure 4.5 IR spectra plotted by using the target composition and the return composition of Experiment 18 . . . . .	40
Figure 4.6 Raman spectra plotted by using the target composition and the return composition of Experiment 18 . . . . .	41
Figure 4.7 SFG spectra plotted by using the target composition and the return composition of Experiment 18 . . . . .	42
Figure 5.1 Accuracy analysis for experiments considering a mixture of amino acids with candidates from $0^\circ$ to $80^\circ$ on $\theta$ for each amino acid. Accuracy indicates how many times each experiment in the set return a composition matches the target one. . . . .	47
Figure 5.2 IR Spectra Plotted by Result Composition and Target Composition. . . . .	49
Figure 5.3 Accuracy analysis for experiments considering a mixture of amino acids with candidates from $0^\circ$ to $180^\circ$ on $\theta$ for each amino acid. Accuracy indicates how many times each experiment in the set return a composition matches the target one. . . . .	50
Figure 5.4 Target composition of one random run of six mixed amino acids with candidates expanded from $0^\circ$ to $180^\circ$ on $\theta$ for each amino acid. More detailed data of this target composition can be found in A.1 in the Appendix. . . . .	51

Figure 5.5	Return composition of experiment 2 for one random run of six mixed amino acids with candidates expanded from $0^\circ$ to $180^\circ$ on $\theta$ . More detailed data of this target composition can be found in A.2 in the Appendix. . . . .	52
Figure 5.6	Return composition of experiment 6 for one random run of six mixed amino acids with candidates expanded from $0^\circ$ to $180^\circ$ on $\theta$ . More detailed data of this target composition can be found in A.3 in the Appendix. . . . .	52
Figure 5.7	IR $z$ projection spectrum for alanine candidate with $\theta$ of $0^\circ$ is identical to alanine candidate with $\theta$ of $180^\circ$ . . . . .	53
Figure 5.8	Raman $zz$ projection spectrum for alanine candidate with $\theta$ of $0^\circ$ is identical to alanine candidate with $\theta$ of $180^\circ$ . . . . .	54
Figure 5.9	SFG $zzz$ projection spectrum for alanine candidate with $\theta$ of $0^\circ$ is not identical to alanine candidate with $\theta$ of $180^\circ$ , but symmetric along wavelength . . . . .	55
Figure 6.1	Target composition for one random run of experiment set with scaling factor for mixed amino acids with $\theta$ expended from $0^\circ$ to $80^\circ$ . . . . .	58
Figure 6.2	Return composition of Experiment 2 for one random run of experiment set with scaling factor for mixed amino acids with $\theta$ expended from $0^\circ$ to $80^\circ$ . . . . .	59
Figure 6.3	Experiment accuracy analysis for experiments using experimental spectra data that contains scaling factor that is smaller than 1 and candidates with $\theta$ from $0^\circ$ to $80^\circ$ . . . . .	60
Figure 6.4	Target composition of one random run of experiments containing scaling factor and the mixed amino acids' candidates with $\theta$ expended from $0^\circ$ to $180^\circ$ . . . . .	61
Figure 6.5	Return composition of Experiment 2 for one random run of experiments containing scaling factor and the mixed amino acids' candidates with $\theta$ expended from $0^\circ$ to $180^\circ$ . . . . .	61
Figure 6.6	Return composition of Experiment 6 for one random run of experiments containing scaling factor and the mixed amino acids' candidates with $\theta$ expended from $0^\circ$ to $180^\circ$ . . . . .	62

## ACKNOWLEDGEMENTS

I would like to thank:

**My husband**, for supporting me in the low moments.

**Dr. Ulrike Stege**, for all the support, encouragement, inspiration and patience. I can only finish my thesis with her all help and courage.

**Dr. Dennis Hore**, for always giving me new ideas and wonderful discusses.

**Kuo Kai Hung**, for previous working and information sharing.

**PITA and Dennis groups**, for all the fun and knowledge we share in our weekly meeting.

*I believe I know the only cure, which is to make one's centre of life inside of one's self, not selfishly or excludingly, but with a kind of unassailable serenity-to decorate one's inner house so richly that one is content there, glad to welcome any one who wants to come and stay, but happy all the same in the hours when one is inevitably alone.*

Edith Wharton



## DEDICATION

Just hoping this is useful!

# Chapter 1

## Introduction

### 1.1 Background and Motivation

An interface is what forms a common boundary between two phases of matter. The phases of matter can be of any forms, i.e, solid, liquid, and gas. The behavior of a surface greatly affects the properties of a material, such as oxidation, corrosion, chemical activity, deformation and fracture, surface energy and tension, adhesion, bonding, friction, lubrication, wear and contamination. Therefore, surface characterization identification remains an active area of research in the physics, chemistry, and biotechnology communities as well as in modern electronic technology. It also plays a crucial role in surface science. Among various surface properties, molecular orientation is a key factor of all, because molecular orientation greatly affects molecules' surface properties in aspects such as: adhesion, lubrication, catalysis, bio-membrane functions and so on. [9]

Many experimental techniques have been applied in the study of molecular orientation at interfaces. Among them the optical methods are more preferable. Such methods include infrared (IR) absorption, Raman scattering and visible-infrared sum-frequency generation (SFG) spectroscopies. All these vibrational spectra carry quantitative structural information of molecules at interfaces. Although each of them has its own strengths and shortcomings, they all share the following advantages when compared with other non-optical methods. First of all, they all can be applied to any interfaces accessible by light. Second, they are non-destructive. Third, they are highly sensitive to good spatial, temporal and spectral resolutions [1], [9]. An important ad-

vantage of SFG techniques is that it can discriminate against bulk contributions. This means that its result will not take the effect from the bulk. In order to extract the quantitative structural information that molecules carry at interfaces, different spectroscopy techniques and analyse are required. Combining different spectroscopy techniques is a very effective way to achieve the goal of molecular coordination study at interfaces. However, finding the most effective ways to combine these techniques may not be clear most of time.

In order to analyze these vibrational spectra, various factors need to be considered. For example, a molecule’s vibrational mode in the molecular frame, the orientation average of the molecules adsorbed onto the interface based on the mathematical distribution function and projecting the vibrational mode properties from molecular frame to laboratory frame. The main focus of our study is to combining Linear Programming (LP) with different spectral information to obtain molecular coordination distribution at interfaces. In the following study, we will explore how LP can facilitate extracting quantitative structural information of molecules at interfaces.

Our approach is to first study our LP model’s properties by applying it to a toy model of a small molecule. After that, the LP model is applied to the real molecules to further explore the possibilities of our LP model. The real molecules that we are focusing are six amino acids: methionine, leucine, isoleucine, alanine, threonine and valine.

Before introducing the LP model and the molecule coordination studies, the basic theory of the IR, Raman and SFG spectra is introduced.

## 1.2 Experimental Probes: IR, Raman, SFG

Vibrational spectra (IR, Raman and SFG) are produced by the changes of a molecule’s dipole moment and polarizability. The dipole moment and polarizability are changing as the molecule’s conformation is changing.

IR is the absorption of the absorption-transmission-reflection mode (resonant). The physical principle is the variation of the static dipole moment  $\mu$  (the first rank

tensor) along the normal coordinates  $Q$ :  $\partial\mu/\partial Q$ .

$$I_{IR} \approx \left| \frac{1}{\sqrt{2m_Q w_Q}} \frac{\partial\mu}{\partial Q} \right|^2 \quad (1.1)$$

where  $m_Q$  is the reduced mass of the normal mode, and  $w_Q$  is the resonance frequency. The dipole moment  $\mu$  is a vector of  $x$ ,  $y$  and  $z$ . The dipole moment derivatives can be expressed as Equation 1.2. The IR spectra can be obtained from 3 polarizations:  $x$ ,  $y$ ,  $z$ .

$$\frac{\partial\mu}{\partial Q} = \begin{bmatrix} \partial\mu_x/\partial Q \\ \partial\mu_y/\partial Q \\ \partial\mu_z/\partial Q \end{bmatrix} \quad (1.2)$$

Raman is scattered from a molecule sample. Unlike IR, Raman spectra relate to the variation of the molecular polarizability  $\alpha$  (the second rank tensor) along the normal coordinates  $Q$ :  $\partial\alpha/\partial Q$ .

$$I_{Raman} \approx \left| \frac{1}{\sqrt{2m_Q w_Q}} \frac{\partial\alpha^{(1)}}{\partial Q} \right|^2 \quad (1.3)$$

where  $m_Q$  and  $w_Q$  are the same as defined in Equation 1.1. The polarizability is coupled with  $(x, y, z)$  components of the driving field and  $x, y, z$  components of the induced polarization. Therefore, there are 9 elements in the polarizability, which can be expressed as Equation 1.4. It results in 9 polarizations of Raman spectra:  $xx$ ,  $yy$ ,  $zz$ ,  $xy$ ,  $xz$ ,  $yx$ ,  $yz$ ,  $zy$  and  $zx$ .

$$\frac{\partial\alpha^{(1)}}{\partial Q} = \begin{bmatrix} \frac{\partial\alpha_{xx}^{(1)}}{\partial Q} & \frac{\partial\alpha_{xy}^{(1)}}{\partial Q} & \frac{\partial\alpha_{xz}^{(1)}}{\partial Q} \\ \frac{\partial\mu_{yx}}{\partial Q} & \frac{\partial\alpha_{yy}^{(1)}}{\partial Q} & \frac{\partial\alpha_{yz}^{(1)}}{\partial Q} \\ \frac{\partial\mu_{zx}}{\partial Q} & \frac{\partial\alpha_{zy}^{(1)}}{\partial Q} & \frac{\partial\alpha_{zz}^{(1)}}{\partial Q} \end{bmatrix} \quad (1.4)$$

SFG stands for sum frequency generation vibrational spectroscopy. SFG is a surface-specific technique. It is a non-linear optical process. It is sensitive to the molecular orientation in odd orders. Comparing to linear optical spectroscopy, the biggest advantage of SFG is that it is surface specific. The spectroscopy signal only

comes from the surface, not the bulk. SFG is the variation of the outer product of dipole moment and polarizability,  $\chi^{(2)}$  (the third rank tensor):  $\partial\mu/\partial Q \otimes \partial\alpha/\partial Q$ . Therefore, there are 27 elements for SFG spectra, which result in 27 polarizations of SFG spectra.

$$I_{SFG} \approx \left| \frac{1}{2m_Q w_Q} \left( \frac{\partial\alpha^{(1)}}{\partial Q} \otimes \frac{\partial\mu}{\partial Q} \right) \right|^2 \quad (1.5)$$

### 1.3 Linear programming

LP problems are an optimization ones of a specific form. The standard form of LP is a minimization problem that has an objective function and a number of constraints as shown in Equation 1.6 [4]:

$$\begin{aligned} & \text{minimize} && c_1x_1 + c_2x_2 + \dots + c_nx_n \\ & \text{subject to} && a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ & && a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2 \\ & && \cdot && \cdot \\ & && \cdot && \cdot \\ & && \cdot && \cdot \\ & && a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n = b_m \\ & && x_1 \geq 0, x_2 \geq 0, \dots, x_n \geq 0, \end{aligned} \quad (1.6)$$

where  $x_i$  are the decision variables,  $a_{ij}$  is a matrix of know coefficients,  $b_i$  and  $c_i$  are vectors of known coefficients. The expression to be minimized is called objective function. The equalities and the inequalities are the constraints that all the decision variables need to subject to. These constraints specify a convex polytope that the objective function need to optimize over.

The diet problem is a popular example to illustrate the concept of LP. It is described as follows: a restaurant would like to achieve the minimal nutrition requirements with the lowest price over some the food selections as shown in Table 1.1. For each meal, the minimum requirements for vitamin A, vitamin C and dietary fiber are

Food	Carrot	Cabbage	Cucumber	Required per dish
Vitamin A [mg/kg]	35	0.5	0.5	0.5mg
Vitamin C [mg/kg]	60	300	10	15mg
Dietary Fiber [g/kg]	30	20	10	4g
price[\$/kg]	0.75	0.5	0.15	-

Table 1.1: Sample input of the diet problem

0.5 mg, 15 mg and 4 g. The restaurant has three food options: raw carrot, raw white cabbage and pickled cucumber. The table also displays the nutrition content and the price of each ingredient. With all the information, we want to know how much carrot, cabbage and cucumber is needed in each meal, so that the minimal nutrition requirements can be met with the lowest price. In summary, the goal is to minimize the price, and the constraints are the nutrition requirements. Therefore, the following LP model is come up as shown in Equations 1.7 – 1.13.

$$\text{minimize} \quad 0.75x_1 + 0.5x_2 + 0.15x_3 \quad (1.7)$$

$$\text{subject to} \quad 35x_1 + 0.5x_2 + 0.5x_3 \geq 0.5 \quad (1.8)$$

$$60x_1 + 300x_2 + 10x_3 \geq 15 \quad (1.9)$$

$$30x_1 + 20x_2 + 10x_3 \geq 4 \quad (1.10)$$

$$x_1 \geq 0 \quad (1.11)$$

$$x_2 \geq 0 \quad (1.12)$$

$$x_3 \geq 0 \quad (1.13)$$

In this LP model,  $x_1$ ,  $x_2$  and  $x_3$  are the decision variables. Each decision variable presents the amount of each ingredient. Equation 1.7 is the objective function to minimize. Equation 1.8 to Equation 1.10 describe the nutrition requirements. Equation 1.11 to Equation 1.13 ensure the amount of each ingredient to be greater than 0. With the existing LP solvers that implemented Simplex Method, the optimal solution can be obtained within a second.

For a LP problem, there are only three kinds of solutions: feasible and bounded solutions, feasible and unbounded solutions, and infeasible solutions. If the solution space is feasible and bounded, then there is one optimum solution. If it is feasible

but unbounded, then there is a solution space with an infinite number of optimal solutions [2].

A general LP problem can be a minimization or maximization problem. Its constraints can be equalities or inequalities. For each non-standard LP problem, there are ways to convert it into its standard form. Furthermore, for a LP problem that contains  $n$  decision variables, its solution would be in a  $n$ -dimensional space called  $R^n$ . Each constraint is a hyperplane. It divides the  $R^n$  space into two half-spaces. Therefore, all the constraints together cut this  $R^n$  space into a convex polyhedron when there are feasible solutions. This makes LP a convex problem. The benefit of a convex problem is that the local optimal solution is the global optimum. LP solvers return the optimal solution. If a LP problem has a unique optimal solution, this solution is a vertex of the convex polyhedron. In another word, LP is a convex, deterministic process. It is guaranteed to converge to a single global optimum if there is a solution space.

Another advantage of LP is it can deal with thousands of variables, which makes it suitable for the study of a molecule’s coordination composition at interfaces. Furthermore, LP problems are intrinsically easier to solve than many non-linear problems.

Various algorithms are available in solving LP problems, such as: Simplex algorithm, Interior point, and Path-following algorithms. Both Interior Point and Simplex are common and mature algorithms that work well in practice. Simplex is comparatively easier to understand and implement than Interior Point. Simplex method takes the advantage of the geometric concept that it visits the vertices of the feasible set (convex polyhedron), and check the optimal solution among each visited vertex. The converging approach is also different for these two methods. If there are  $n$  decision variables, usually Simplex will converge in  $O(n)$  operations with  $O(n)$  pivots. Interior point traverses the edges between vertices on a polyhedral set. Generally speaking, Interior point method is faster for larger problems with sparse matrix. However, when experimenting with these two methods, the speed of them is not much different from each other for our study. For our study, Simplex method has proved to be efficient and effective, and it is used for all the experiments.

Last but not the least advantage of LP is its speed. For any LP problem, if it

has an optimal solution, this solution is always a vertex. Simplex method is based on this insight, namely that it starts at a vertex, then pivot from vertex to vertex, until it reaches the optimum. Although it has been shown that Simplex method is not a polynomial algorithm, in practice it usually takes  $2n - 3n$  steps to solve a problem ( $n$  is the number of decision variables).

The LP solver we use is called “GNU linear programming tool kit”(GLPK). It has implemented both Simplex and Interior Point methods in ACNSI C. It is open-source and intended to solve large scale LP problems.

## 1.4 Aims and scope

Given some target experimental spectra and a set of candidates spectra, then figuring out the right combination of candidates for the target spectra is the goal in this study. The approach is to build a LP model, and check if the optimal solutions returned by the solver match the target composition pre-generated. Spectral information of different spectroscopy techniques is applied to the LP model, then we analyze which spectral information helps to obtain the target composition with the highest accuracy. Furthermore, various types of candidate situations are considered, such as: candidates coming from one type of molecule; candidates coming from a mixture of molecules. At last, the experimental spectral information is brought into consideration.

## 1.5 Overview of The Thesis

Chapter 1 briefly introduce the aim and scope of the current study. Chapter 2 explains the current approaches to extract the molecular structure at interfaces, as well as how produce IR, Raman and SFG spectra. Chapter 3 aims to use a simplified molecule model to study the properties of our LP model. Chapter 4 applies the LP model to one type of molecule at interfaces. Chapter 5 applies the LP model to a mixture of different molecules at interfaces. Chapter 6 applies the LP model to experimental spectral data. Chapter 7 is the conclusion and future work.



# Chapter 2

## Methods

### 2.1 Current Approaches to Molecular Structure Elucidation

Currently, there are two main approaches in studying the orientation distribution of molecules at interface. One is comparing the experimental spectra with few predicted ones, and select the one that most matches to the experimental one. Another one is running an exhaustive algorithm to explore the most possible solution space [8]. However, both approaches take a lot of time and computational resources. In Hung’s study [3], a new approach is introduced. He applied LP to vibrational spectra to extract the molecular structure at interfaces. This LP approach helped to return the target orientation distribution information when the mock experimental spectrum consisted of different amino acids. However, when candidates are coming from the same amino acid, LP approach failed to return the target orientation distribution information. The reason why LP failed to return the target composition has not been thoroughly studied in Hung’s study. Whether and how LP approach can be generally applied to different experiment situations have not been explored. My study is to figure out the underlying properties of our LP model. Furthermore, explore the applicability of the LP model to different experimental setting.

## 2.2 Structure of molecules adsorbed to interfaces

(TODO: check with Dennis, how to expand this part.) A picture to display molecules adsorbed to interfaces

## 2.3 Generating model spectra

As mentioned in Chapter 1, before analyzing the vibrational spectra of amino acids, there are a few factors to address. First of all, creating candidate spectra is an essential step. This part of research has been done thoroughly by Hung [3].

To generate these amino acids' vibrational spectra, a molecule's vibration modes need to be modelled in the molecular frame, then transferred to the laboratory frame to work with the systems where interfaces exist. Chapter 2 in Hung's thesis [3] describes how to perform electronic structure calculations using GAMESS [5] to obtain the derivatives of every dipole moment and polarizability. Then he introduced how to use Direction Cosine Matrix (DCM) to transfer these two derivatives from the molecular coordinate system to the laboratory one. After that, Euler angles could be extracted from DCM. Euler angles are used to describe a molecule's coordination at interfaces. They are labelled by  $\theta$ ,  $\phi$  and  $\psi$  as shown in Figure 2.1. They are referred as *tilt*, *azimuthal* and *twist* angles, respectively. Let  $x$ ,  $y$  and  $z$  be lab frame Cartesian coordinates, and  $a$ ,  $b$  and  $c$  be the molecular frame coordinates. *Tilt* angle  $\theta$  is the angle between  $z$  and  $c$ . *Azimuthal* angle  $\phi$  is the rotation about  $z$ . *Twist* angle  $\psi$  is a twist about  $c$  [8]. After three steps of successive rotations of Euler angles, molecule properties can be transferred from the molecular frame to the lab frame.

In order to achieve the above steps, Hung first did a Hessian calculation using GAMESS. Secondly, 7 snapshots of a molecule vibrating in different modes were taken. Thirdly, he did a force field calculation to obtain the derivatives of dipole moment and polarizability for each 7 snapshot moment. Then the derivatives of dipole moment and polarizability are obtained by the interpolation of these 7 snapshot moment. Because the two obtained derivatives are in the molecular frame, Hung used DCM to convert these two derivatives into the lab frame. Then abstracted Euler angles from DCM. After this transformation, he restored the derivatives information into some molecular property files for any further usage. (TODO: double check the

accuracy with Dennis)

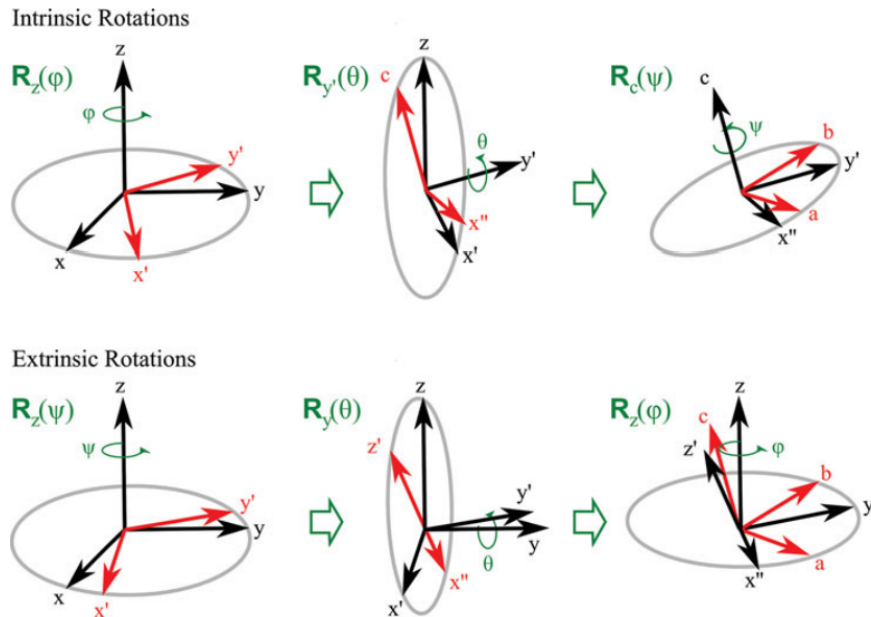


Figure 2.1: The Euler angles represented as the spherical polar angles  $\theta$ ,  $\phi$  and  $\psi$ , and the illustration of the three successive rotations that transform the lab  $x$ ,  $y$ ,  $z$  coordinate system into the molecular  $a$ ,  $b$ ,  $c$  frame intrinsically and extrinsically [8].

In my study, those molecular property files are used to generate the amino acids' spectroscopy information directly. Each molecular property file contains the derivatives of dipole moment and polarizabilities of each vibrational mode. Depends on the number  $N$  of atoms in a molecule, there are  $3N - 6$  vibrational modes. Furthermore, Equation 2.2 to 2.4 are used to generate the amino acids' IR, Raman and SFG spectra.

All the experiments in my study are limited to only consider the *tilt* angle distribution of Euler angles, and assume isotropy on *twist* and *azimuthal* angular distributions. Therefore, *twist* and *azimuthal* angles are integrated to create a uniform distribution. For angle  $\psi$ , it requires the surfaces to be not striped. There can be no anisotropy in the plane of the surface. Because of this, we can limit the candidate number by integrating angle  $\psi$ . On the other hand, for angle  $\phi$ , a uniform distribution implies that the molecule has cylindrical symmetry in its preference of surface. This means that the molecule can be tilted, but has no 'twist' preference. With the integration of these two Euler angles, the number of candidates for one molecule will

be greatly reduced. However, the number of the amino acid candidates is still large when only considering  $\theta$  angle. The possible combinations of all these amino acid candidates are still considered to be excessive (TODO: put these into an approximated number???).

Furthermore, when molecules lay on an interface, the orientation of each molecule varies. To simulate the vibrational spectra, a reasonable orientation distribution for the molecules needed to be studied. The orientation distribution requires either do a molecular dynamic simulation to study the distribution of molecule orientations at the interface, or come up with a analytic orientation distribution function. In my study, the second method is preferred. Moreover, Delta distribution function shown in Equation 2.1 is used to represent the molecule orientation distribution that models the spectrum signals. This means that all the molecules are tilted at one same angle at the interface. This assumption is applied across the whole study.

$$f_{(\theta)} = \delta(\theta - \theta_o) \quad (2.1)$$

Infrared (IR) absorption spectroscopy is a harmonic approximation, its intensity is proportional to the square of the lab-frame dipole moment derivative. For example, the  $x$ -polarized absorption spectrum is given by Equation 2.2.

$$I_x(\omega_{\text{IR}}) = \sum_q \frac{1}{2m_q\omega_q} \left\langle \left[ \frac{\partial u_x}{\partial Q} \right]_q^2 \right\rangle \frac{\Gamma_q^2}{(\omega_{\text{IR}} - \omega_q)^2 + \Gamma_q^2} \quad (2.2)$$

where  $I_x$  represents  $x$ -polarized intensity. The same equation applies to  $I_y$  and  $I_z$ .  $\omega_{\text{IR}}$  is the frequency of the probe radiation.  $\mu$  is the dipole moment.  $m_q$  is the reduced mass.  $\omega_q$  is resonance frequency.  $\Gamma_q$  is the homogeneous line width, is set to 6 in all the experiments.  $Q_q$  is the normal mode coordinate of the  $q$ th vibrational mode. All values of  $\omega_{\text{IR}}$ ,  $\mu$ ,  $m_q$ ,  $Q$  are obtained from the molecular property files. Furthermore, because  $\phi$  and  $\psi$  angles are integrated, the  $y$ -polarized spectrum is identical with the  $z$ -polarized one. Therefore, there are only two unique polarized IR spectra. For simplicity, IR spectra are referred as  $y$  and  $z$  in future experiments.(TODO: need to double check the accuracy with Dennis)

The intensity of Raman scattering is proportional to the square of laboratory-

frame transition polarizability. For example, Raman spectroscopy with an  $x$ -polarized excitation source collects the  $x$ -polarized component of the scattered radiation, which can be approximated from Equation 2.3.

$$I_{xx}(\Delta\omega) = \sum_q \frac{1}{2m_q\omega_q} \left\langle \left[ \frac{\partial\alpha_{xx}^{(1)}}{\partial Q} \right]_q^2 \right\rangle \frac{\Gamma_q^2}{(\Delta\omega - \omega_q)^2 + \Gamma_q^2} \quad (2.3)$$

where  $\Delta\omega$  is the Stokes Raman shift.  $\alpha_{xx}^{(1)}$  is one component of the 9-element polarizability tensor.  $m_q$ ,  $\omega_q$ ,  $\Gamma_q$ , and  $Q_q$  are the same as defined above for IR spectra. All the values of  $\omega_{\text{IR}}$ ,  $\mu$ ,  $m_q$ ,  $Q$  are obtained from the molecular property files. Similar to IR spectroscopy, because of the integration of  $\phi$  and  $\psi$  angles, only 4 unique spectra are obtained from the following polarization:  $xx$ ,  $xy$ ,  $xz$  and  $zz$ . For simplicity, Raman spectra are referred as  $xx$ ,  $xy$ ,  $xz$  and  $zz$  in future experiments (TODO: double check the accuracy of the content with Dennis).

The intensity of SFG spectroscopy is proportional to the squared magnitude of the second-order susceptibility,  $|\chi^{(2)}|^2$ .  $\chi^{(2)}$  is derived from the second-order polarizability,  $\alpha^2$ . Equation 2.4 shows the response intensity of  $I_{xxx}$ .

$$I_{xxx}(\omega_{\text{IR}}) = \sum_q \frac{1}{2m_q\omega_q} \left\langle \left[ \frac{\partial\alpha_{xx}^{(1)}}{\partial Q} \right]_q \left[ \frac{\partial u_x}{\partial Q} \right]_q \right\rangle \frac{1}{\omega_q - \omega_{\text{IR}} - i\Gamma_q} \quad (2.4)$$

where  $I_{xxx}$  is the second-order susceptibility tensor. It is probed by an  $x$ -polarized visible incoming beam at frequency  $\omega_{\text{vis}}$  and a  $x$ -polarized infrared beam incoming with frequency  $\omega_{\text{IR}}$ . Both incoming beams are incident to the sample. Then the  $x$ -component of SFG at frequency  $\omega_{\text{SFG}} = \omega_{\text{vis}} + \omega_{\text{IR}}$  is selected for detection. As  $i = \sqrt{-1}$  is in the denominator,  $\chi^{(2)}$  is a complex value [3]. The SFG response is the imaginary component of the second-order susceptibility. Same as IR and Raman spectroscopy, all the values of  $\omega_{\text{IR}}$ ,  $\mu$ ,  $m_q$ ,  $Q$  are obtained from the molecular property files. Because of the integration of  $\phi$  and  $\psi$  angles, only 3 unique non-zero spectra are obtained from the following polarizations:  $yyz$ ,  $zyz$  and  $zzz$ . For simplicity, SFG spectra are referred as  $yyz$ ,  $zyz$  and  $zzz$  in future experiments. (TODO: double check the accuracy of the content with Dennis).

With these equations and the molecular property files, IR, Raman and SFG spec-

tra can be generated for a candidate. A candidate in my study is a specific amino acid with specific  $\theta$  value. Taking Methionine as an example, Figure 2.2 displays  $x$ -polarized IR spectra of the following candidates: Methionine with  $\theta$  of  $0^\circ$ ,  $20^\circ$ ,  $40^\circ$  and  $60^\circ$ . Their spectra are prefixed with *candidate\_* in the labels. *ir\_x\_* indicates the spectroscopy technique, “number” indicates the  $\theta$  angle’s value. The spectra labelled as *target\_ir\_x*, is generated by combining 10% of *candidate\_ir\_x\_0*, 50% *candidate\_ir\_x\_20* and 40% *candidate\_ir\_x\_40*.

Similarly, Figure 2.3, 2.4 and 2.5 depict the spectra of the same candidates and targets for  $z$ -polarized IR,  $xx$ -polarized Raman and  $yyz$ -polarized SFG spectrum respectively. In Figure 2.2, the biggest differences among the candidates exist at each vibrational mode. The valid range for wavenumber is from 1000 to 2000. Each polarization of IR, Raman or SFG, there are 200 data points can be extracted in the interval of 5 wavenumber. With these data points, the corresponding LP model is constructed as described in Chapter 3.

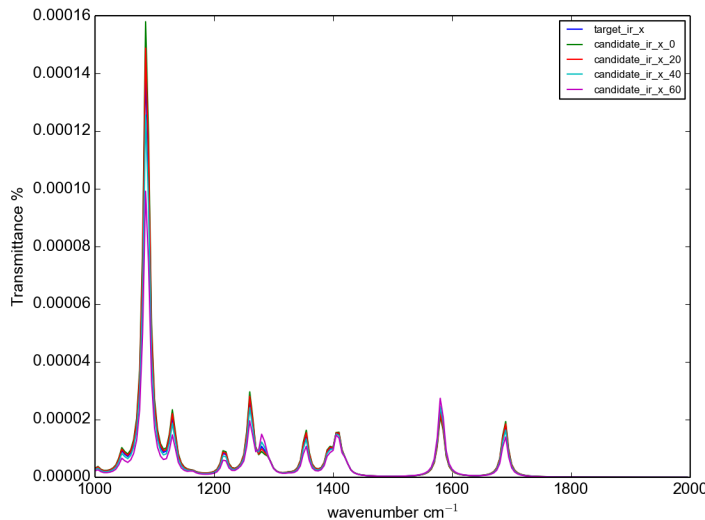


Figure 2.2: IR  $x$ -polarized spectra of methionine’s four candidates and target. The candidates are with  $\theta$  of  $0^\circ$ ,  $20^\circ$ ,  $40^\circ$  and  $60^\circ$ . The target is produced by combining 10% of *candidate\_ir\_x\_0*, 50% *candidate\_ir\_x\_20* and 40% *candidate\_ir\_x\_40*.

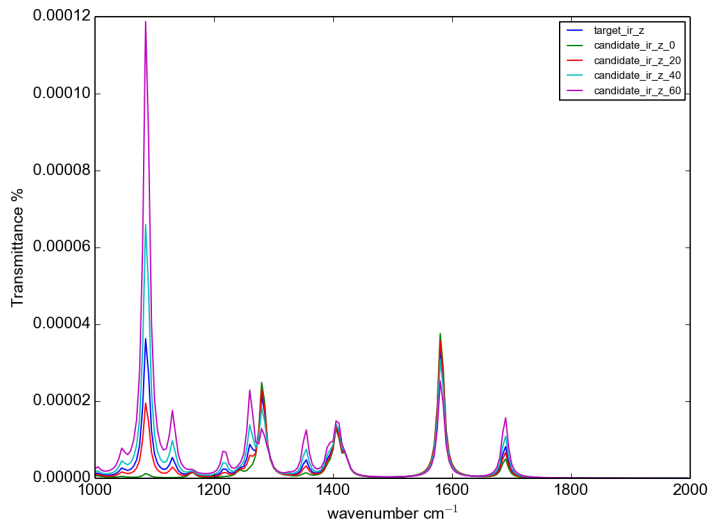


Figure 2.3: IR  $z$ -polarized spectra of methionine's four candidates and target. The candidates are with  $\theta$  of  $0^\circ$ ,  $20^\circ$ ,  $40^\circ$  and  $60^\circ$ . The target is produced by combining 10% of *candidate\_ir\_x\_0*, 50% *candidate\_ir\_x\_20* and 40% *candidate\_ir\_x\_40*.

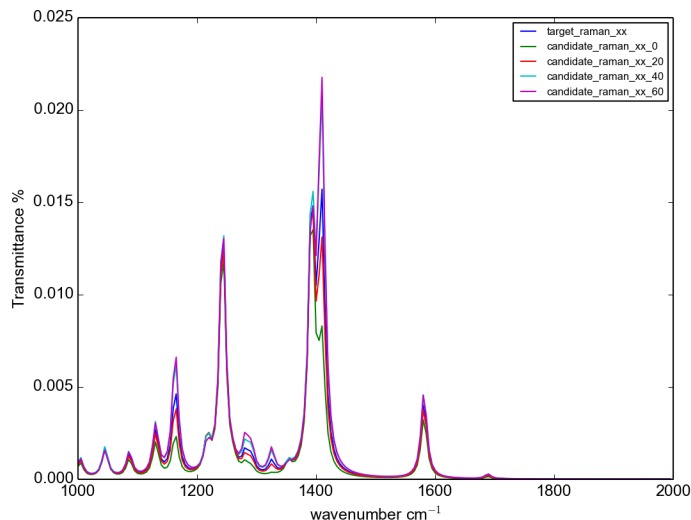


Figure 2.4: Raman  $xx$ -polarized spectra of methionine's four candidates and target. The candidates are with  $\theta$  of  $0^\circ$ ,  $20^\circ$ ,  $40^\circ$  and  $60^\circ$ . The target is produced by combining 10% of *candidate\_ir\_x\_0*, 50% *candidate\_ir\_x\_20* and 40% *candidate\_ir\_x\_40*.

## 2.4 The Properties of the LP Models

Chapter 2 explains what are the current approaches to extract molecular structure at interfaces, and how to produce IR, Raman and SFG spectra theoretically. In Chapter

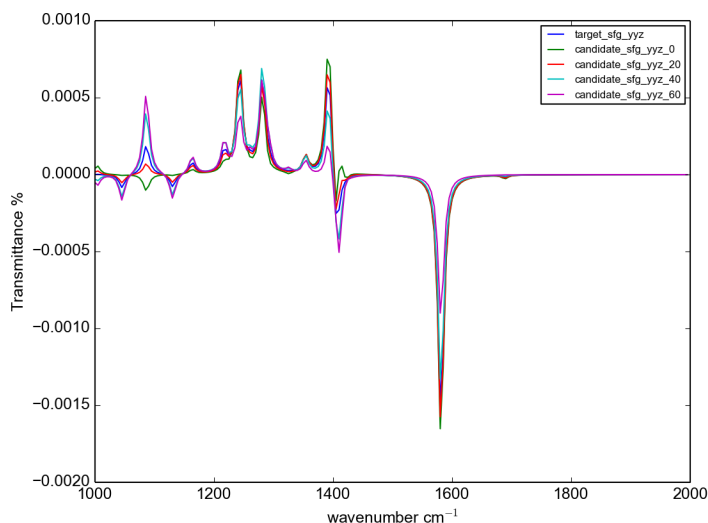


Figure 2.5: SFG  $yyz$ -polarized spectra of methionine's four candidates and target. The candidates are with  $\theta$  of  $0^\circ$ ,  $20^\circ$ ,  $40^\circ$  and  $60^\circ$ . The target is produced by combining 10% of *candidate\_ir\_x\_0*, 50% *candidate\_ir\_x\_20* and 40% *candidate\_ir\_x\_40*.

3, the properties of the LP model are studied. It is conducted by using a toy model to gain an insight of the behaviours our LP approach. The motivation of creating a toy model is to create a molecule as simple as possible, so that only the properties of the LP model is focused. With the further information gained in Chapter 3, further experiences will be conducted to real molecules in Chapter 4, 5 and 6.



## Chapter 3

# Simplified Molecular Model

### 3.1 Description

The goal of Chapter 3 is to introduce the formulas used to describe our LP model. As well as exploring the properties of our LP model by using a toy molecule. This toy molecule contains limited vibration modes. By doing so, the nature of the LP model we use to study the spectral information can be carefully analyzed. Our goal is to figure out with the the spectral information available, could LP model we use output any valuable information.

The toy molecule contains 4 vibration modes. Theses vibrational peaks are at frequencies of 2850, 2960, 3050 and 3200. The widths of the peaks are 5, 10, 5 and 15  $cm^{-1}$ , respectively. The amplitudes of the peak are 1, 0.7,  $-0.2$  and 0.5  $cm^{-1}$ , respectively. The comparing angles of the peaks are 15, 90, 0 and 60. (TODO: check with Dennis, how to further explain those comparing angles?)

For the toy model, only IR spectroscopy is considered. Because we want to limit the complication that comes from the parameters needed to describe the real molecules. Equation 3.1 is used to generate the *cosine*-polarized IR spectrum. Both  $\phi$  and  $\psi$  Euler angles are integrated, only the difference on angle  $\theta$  is considered.

$$f_{\theta}(x) = \sum_{q=1}^4 A_q^2 * \cos^2(\theta - \theta_q) \frac{\gamma^2}{(x - \omega_q)^2 + \gamma^2} \quad (3.1)$$

where  $A$  is the amplitude,  $\theta_q$  is the comparing angle,  $\gamma$  is the width, and  $\omega_q$  is the frequency. (TODO: Double check the correct meaning of each symbol) Ten candidates are produced with 10 different  $\theta$  values as follows:  $0^\circ$ ,  $10^\circ$ ,  $20^\circ$ ,  $30^\circ$ ,  $40^\circ$ ,  $50^\circ$ ,  $60^\circ$ ,  $70^\circ$ ,  $80^\circ$ ,  $90^\circ$ . Their spectra are shown in Figure 3.1. The 10 candidates have peaks at the same frequencies. The spectral signal for candidates is comparatively strong at each peak.

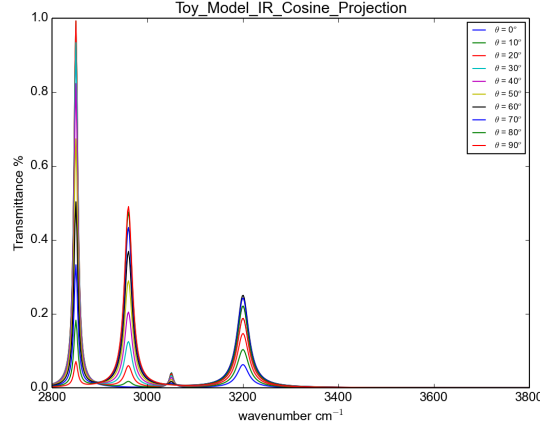


Figure 3.1: *cosine*- polarize IR spectra of toy model candidates

## 3.2 Linear Programming Model for Spectral Study

Equation 3.2 is used to construct our LP model. The optimal solution returned by the LP solver is then compared with the target composition to see if they matches each other. This equation has also been used to study the composition of Ribonucleic acid (RNA) with ultraviolet (UV) spectra [6] and other UV spectroscopy studies [7] back in the 60s.

$$\underset{p_c}{\text{minimize}} \quad \sum_{n=1}^{\# \text{ points}} \left| \text{Target} - \sum_{c=1}^{\# \text{ candidates}} p_c f_{\theta}(x) \right| \quad (3.2)$$

where  $p_c$  are the unknown percentages for each candidate, which are the decision variables.  $n$  is the number of points selected along the wavenumber, both for candidates

and target spectra. *Target* refers to the corresponding data points selected in target spectra. For each data point, the absolute residual between the target spectrum and the one composed by the decision variables is calculated. The objective function minimizes the sum of the absolute residuals over all the data points.

Because Equation 3.2 subjects to no restrictions, and the objection function is not in standard form. Getting rid of the absolute signs in the objective function is needed in order to use an LP approach. To eliminate the absolute sign is achieved by introducing one more variable  $X$  and two more constraints for each data point as shown in Equation 3.3. Then the previous model in Equation 3.2 is converted into the one in Equation 3.4 that can be solved by an LP solver. At last, one more constraint is introduced to restrict the sum of the percentages to be 1, as shown in Equation 3.4.

For each point in the range of valid wavenumbers:

$$\begin{aligned}
 X &= \left| Target - \sum_{c=1}^{\# \text{ candidates}} p_c f_{\theta}(x) \right| \\
 X &\geq Target - \sum_{c=1}^{\# \text{ candidates}} p_c f_{\theta}(x) \\
 X &\geq -Target + \sum_{c=1}^{\# \text{ candidates}} p_c f_{\theta}(x)
 \end{aligned} \tag{3.3}$$

$$\begin{aligned}
& \text{minimize} \quad \sum_{n=1}^{\# \text{ points}} X_p \\
& X_1 - \text{Target}_1 + \sum_{c=1}^{\# \text{ candidates}} p_c f_{\theta}(x_1) \geq 0 \\
& X_1 + \text{Target}_1 - \sum_{c=1}^{\# \text{ candidates}} p_c f_{\theta}(x_1) \geq 0 \\
& \dots \\
& X_n - \text{Target}_n + \sum_{c=1}^{\# \text{ candidates}} p_c f_{\theta}(x_n) \geq 0 \\
& X_n + \text{Target}_n - \sum_{c=1}^{\# \text{ candidates}} p_c f_{\theta}(x_n) \geq 0 \\
& \sum_{c=1}^{\# \text{ candidates}} p_c = 1
\end{aligned} \tag{3.4}$$

Note that our LP model exactly describes our problem to be solved. Assuming that we can obtain sufficiently precise data, solving the LP will yield the target composition. Recall if the solution space is feasible and bounded, then there is a unique optimum solution.

### 3.3 Linear Programming Model Implementation

Next, we describe how to solve Equation 3.4 by implementing our LP model. Code is written to generate a file that contains all the candidates' spectral information needed for the experiments. For this step, the molecular properties files are used. For a specific candidate, given a molecular properties file and a  $\theta$  value, the candidate's spectral information is obtained. For toy model, only the value of  $\theta$  is needed, then Equation 3.1 is used to synthesize the spectral information. To further illustrate, a candidate class is written. This class defines candidate's  $x$ - and  $z$ -polarized IR spectra;  $xx$ -,  $xy$ -,  $xz$ -, and  $zz$ -polarized Raman spectra;  $yyz$ -,  $zyz$ -,  $zzz$ -polarized SFG spectra. Given a candidate's molecular properties and a  $\theta$  value, an instance of this specific candidate is created. For the toy model, it only contains IR spectral information. Therefore, one candidate only contains *cosine*- and *sine*-polarized IR spectra.

In the second step, more code is written to generate a target composition of a list of needed candidates. Then the target composition is used to generate the target spectra. The probe range, which is the range of the wavenumber, is from 2800 to 3300 for toy model. The probe arrange is from 2000 wavenumber to 3000 wavenumber for real molecules. The target spectral information is generated in the same text file as candidate’s spectral information. Depends on the experiment setting, code can be used to generate text files that contain different spectral information.

In the third step, the LP model is constructed by using the spectral information text file generated in the second step. This part of the code was written by Hung [3]. It reads all the candidates and target spectral information, and builds the LP model as shown in Equation 3.4, then creates CPLEX LP input file.

In the fourth step, we use LP solver “GNU linear programming tool kit” (GLPK) to read the CPLEX LP input file, then obtain the result.

### 3.4 Experiments

In Experiment 1 and 2, 4 candidates are selected, the detailed setting is shown in Table 3.1. In Experiment 1, there are 4 candidates with  $\theta$  of  $0^\circ$ ,  $10^\circ$ ,  $20^\circ$ , and  $30^\circ$ . In Experiment 2, the four candidates are with  $\theta$  values of  $0^\circ$ ,  $5^\circ$ ,  $10^\circ$ , and  $15^\circ$ . Instead of having a 10 degree variance in  $\theta$ , 5 degree difference is applied on  $\theta$  in Experiment 2. This means that when the candidates become more similar to each other than the ones in Experiment 1 as their spectra are more similar. In both experiments, 100 data points are selected evenly along the wavenumber from the spectra of *cosine*-polarized IR. The target composition of the candidates are the same for both experiments. In Experiment 1, the return composition is the same as the target one, however, the return composition for Experiment 2 does not match the target one.

In order to figure out why the return composition in Experiment 2 is different from the target one, the spectra generated by the return composition is plotted together with the target spectra as shown in Figure 3.2. Note that the result spectra is almost identical to the target one. The residual between them is almost 0. In order

Experiment index	1	2
Number of Candidates	4	4
Candidates	[0, 10, 20, 30]	[0, 5, 10, 15]
Target Composition	[0.1, 0.5, 0.4, 0]	[0.1, 0.5, 0.4, 0]
Number of Data Points	100 from <i>cosine</i> -polarized IR	100 from <i>cosine</i> -polarized IR
Return Composition	[0.1, 0.5, 0.4, 0]	[0, 0.796962, 0.103038, 0.1]

Table 3.1: Experiment 1 and 2 setting using toy model

to see whether this observation is a general case, Experiment 3 is set up in Table 3.2. Experiment 3 contains more candidates than Experiments 1 and 2. 10 candidates are included with  $\theta$  values ranging from  $0^\circ$  to  $90^\circ$ .

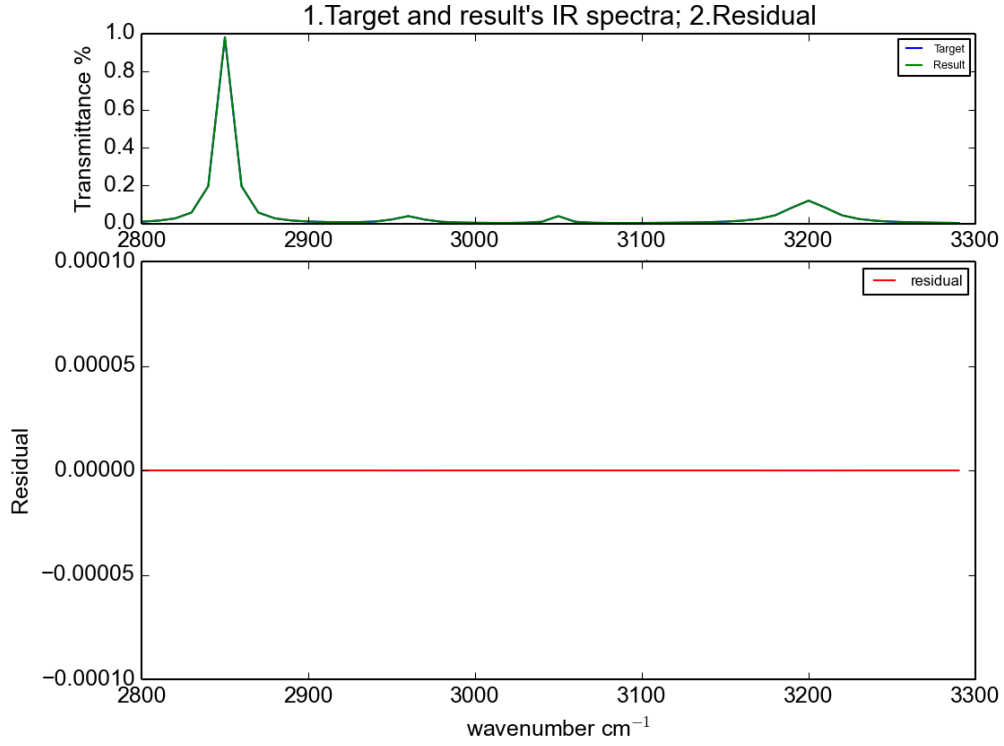


Figure 3.2: Toy model Experiment 2 resulting *cosine*-polarized IR spectrum plotted with the target spectrum; and the residual plot between the spectra.

Table 3.2 indicates the return composition of Experiment 3 is different from the target one. Figure 3.3 shows that the spectrum produced by the return composition is almost identical to the one generated by the target composition in Experiment 3. The residual is negligible as well. This observation is the same as Experiment 2.

Among Experiment 1, 2 and 3, only the return composition of Experiment 1 matches its target one. However, in Experiment 2, the difference in  $\theta$  value among the candidates is smaller than Experiment 1. In Experiment 3, the number of the candidates is larger than Experiment 1. Both effects increase the complexity of the experiments. In both Experiment 2 and 3, the spectrum constructed by the return composition matches to the one built by the target composition.

Experiment index	3
Number of Candidates	10
Candidates	[0, 10, 20, 30, 40, 50, 60, 70, 80, 90]
Target Composition	[0.1, 0, 0.5, 0, 0.4, 0, 0, 0, 0, 0]
Number of Data Points	100 from <i>cosine</i> -polarized IR
Return Composition	[0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0]

Table 3.2: Experiment 3 setting of toy model



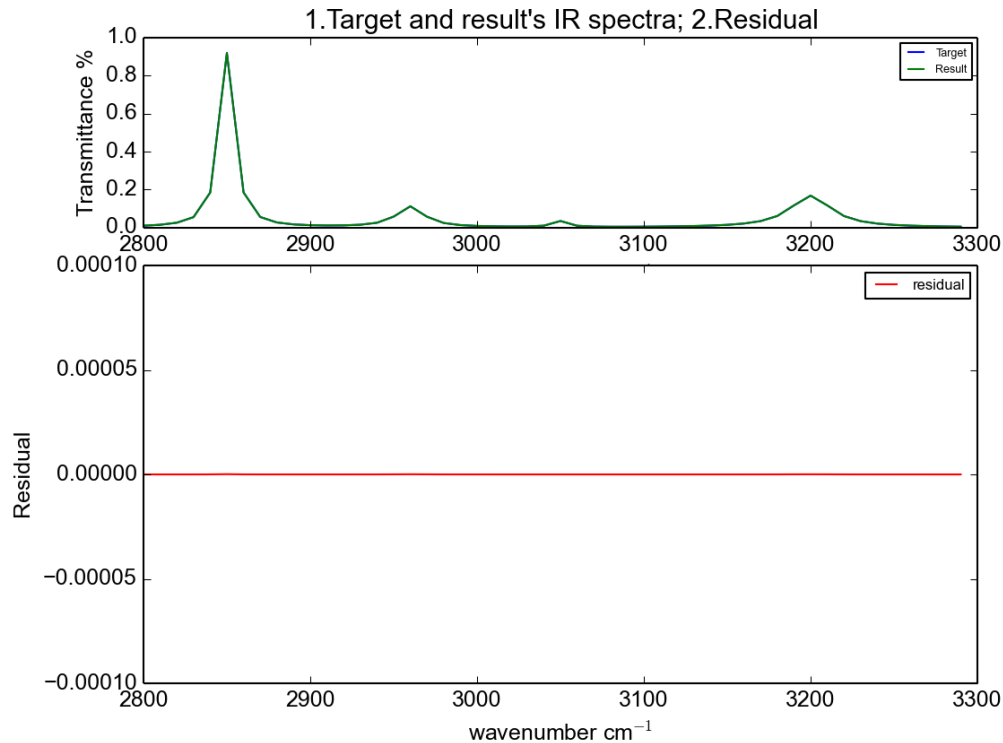


Figure 3.3: Toy model Experiment 3 resulting *cosine*-polarized IR spectrum plotted with target spectrum; and the residual plot between the two spectra

The above observation demonstrates that there are multiple compositions can achieve in constructing the spectrum that are close to the target one. The numerical limitation helps the LP solver to converge to a unique optimum solution. The reason for Experiment 1 to return a composition that matches to the target one, is that the spectral information used to construct the LP model is competent. The constraints constructed in the LP model of Experiment 1 eventually converge to the target composition.

In order to add necessary information to construct the constraints in our LP model, IR's second polarization is introduced to the toy model: the *sine* polarization. Figure 3.4 describes how the *sine*-polarized spectra presented for 10 candidates. Experiment 4 and 5 include both polarizations' spectral information in the LP model. In Table 3.3, Experiment 4's setting is based on Experiment 2, with *sine*-polarized IR spectral information added. 100 data points are selected from this additional spectrum, then

converted to additional decision variables and constraints in the LP model. Experiment 5, it is based on Experiment 3, with *sine*-polarized IR spectral information added. In both Experiment 4 and 5, the return composition matches to the target one. This further proves that as long as we have sufficing information for our LP model, the LP solver returns a composition matches to the target one.

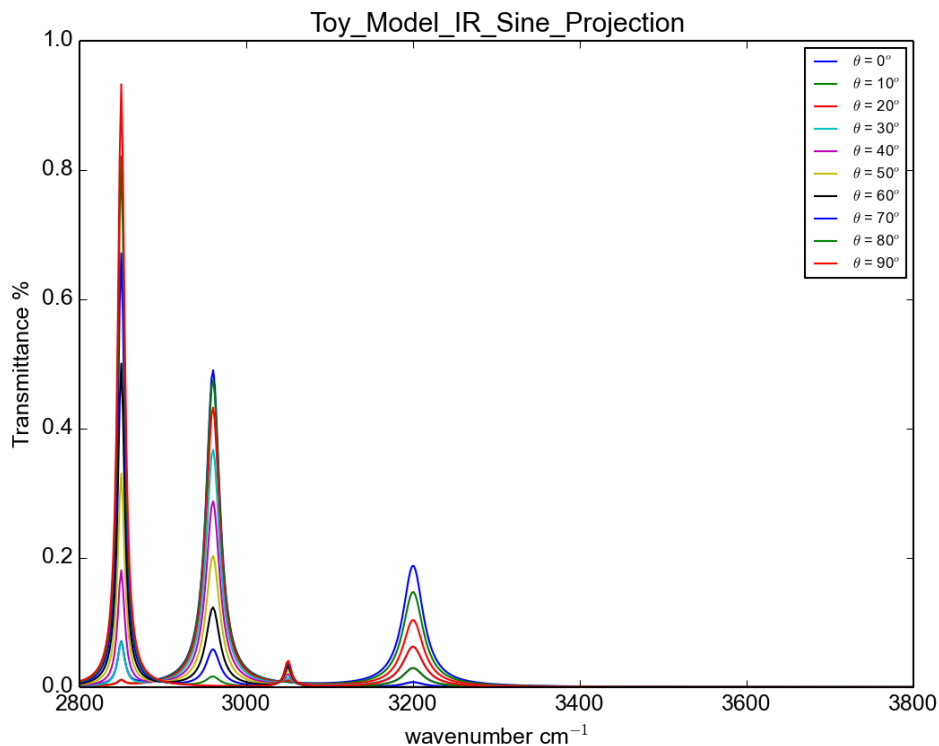


Figure 3.4: *sine*-polarized IR spectra of toy model candidates with  $\theta$  value expanded from  $0^\circ$  to  $90^\circ$

### 3.5 Constraint Study Based on Experiment 4

From Experiment 1 to 5 of toy model, we know having sufficient information in our LP model is the key to obtain the target composition. Having sufficient information means having enough constraints to help LP model converge to a desired result. The information is coming from the valuable data points selected along the spectra. This leads us to do a more detailed study on the constraints in order to see how many data

Experiment index	4	5
Number of Candidates	4	10
Candidates	[0, 5, 10, 15]	[0, 10, 20, 30, 40, 50, 60, 70, 80, 90]
Target Composition	[0.1, 0.5, 0.4, 0]	[0.1, 0, 0.5, 0, 0.4, 0, 0, 0, 0, 0]
Number of Data Points	100 from <i>cosine</i> -polarized IR + 100 from <i>sine</i> -polarized IR	100 from <i>cosine</i> -polarized IR + 100 from <i>sine</i> -polarized IR
Return Composition	[0.1, 0.5, 0.4, 0]	[0.1, 0, 0.5, 0, 0.4, 0, 0, 0, 0, 0]

Table 3.3: Experiment 4 and 5 setting of toy model

Experiment #	# Data Points	Points Selection	Return Composition
6	10	[2800, 3300, 50], <i>cosine</i>	[0, 0.796962, 0.103038, 0.1]
7	20	[2800, 3300, 25], <i>cosine</i>	[0, 0.796962, 0.103038, 0.1]
8	25	[2800, 3300, 20], <i>cosine</i>	[0, 0.796962, 0.103038, 0.1]
9	32	[2800, 3300, 15], <i>cosine</i>	[0, 0.796962, 0.103038, 0.1]
10	50	[2800, 3300, 10], <i>cosine</i>	[0, 0.796962, 0.103038, 0.1]
11	100	[2800, 3300, 5], <i>cosine</i>	[0, 0.796962, 0.103038, 0.1]
12	100 + 1	[2800, 3300, 5], <i>cosine</i> [2800, 3300, 500], <i>sine</i>	[0, 0.796962, 0.103038, 0.1]
13	100 + 5	[2800, 3300, 20], <i>cosine</i> [2800, 3300, 100], <i>sine</i>	[0, 0.796962, 0.103038, 0.1]
14	100 + 10	[2800, 3300, 20], <i>cosine</i> [2800, 3300, 50], <i>sine</i>	[0, 0.796962, 0.103038, 0.1]
15	100 + 50	[2800, 3300, 20], <i>cosine</i> [2800, 3300, 10], <i>sine</i>	[0.1, 0.5, 0.4, 0]
16	100 + 100	[2800, 3300, 20], <i>cosine</i> [2800, 3300, 5], <i>sine</i>	[0.1, 0.5, 0.4, 0]

Table 3.4: Constraint study based on Experiment 4

points are enough to get the desired composition.

Based on Experiment 4, experiments about formulating the LP model with different data information are conducted in Table 3.4. The number of data points indicates how many data points are selected. Points Selection shows how data points are selected. For example, [2800, 3300, 50] means along wavenumber from 2500 to 3300, every 50 wavenumber a data point is selected along a spectrum. *cosine* and *sine* indicate the corresponding polarization of IR spectrum.

As Table 3.4 indicates, the return compositions of Experiment 6 to 14 are the same. To the contrary, from Experiment 15, the return composition matches the target one. In Figure 3.5 displays the spectra conducted by [0, 0.796962, 0.103038, 0.1] and [0.1, 0.5, 0.4, 0], both *sine*- and *cosine*-polarized IR spectra generated by these two compositions are identical.

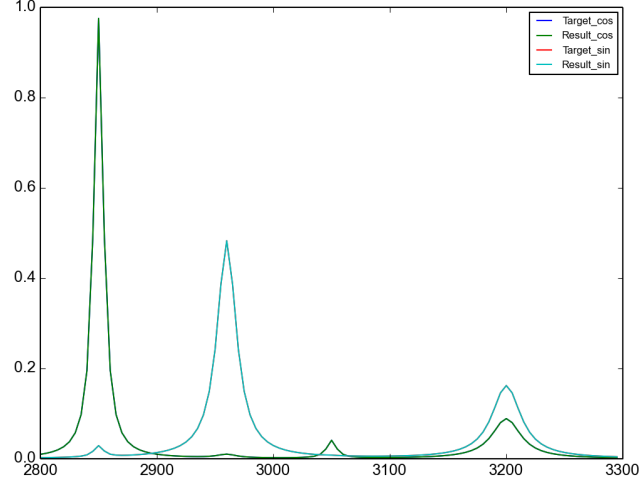


Figure 3.5: IR spectra plotted by the return compositions from the constraint study based on Experiment 4 of toy model

Experiment #	# of Data Points	Point Selection	Return Composition
17	10	[2800, 3300, 50], <i>cosine</i>	[0.156758, 0, 0, 0.825977, 0, 0, 0, 0, 0, 0.017265]
18	25	[2800, 3300, 20], <i>cosine</i>	[0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0, 0]
19	50	[2800, 3300, 10], <i>cosine</i>	[0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0, 0]
20	100	[2800, 3300, 5], <i>cosine</i>	[0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0, 0]
21	500	[2800, 3300, 1], <i>cosine</i>	[0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0, 0]
22	100 + 1	[2800, 3300, 5], <i>cosine</i> [2800, 3300, 500], <i>sine</i>	[0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0, 0]
23	100 + 10	[2800, 3300, 5], <i>cosine</i> [2800, 3300, 50], <i>sine</i>	[0.361587, 0, 0.312061, 0.326352, 0, 0, 0, 0, 0, 0]
24	100 + 20	[2800, 3300, 5], <i>cosine</i> [2800, 3300, 25], <i>sine</i>	[0.174023, 0, 0, 0.791447, 0, 0, 0.0345301, 0, 0, 0]
25	100 + 25	[2800, 3300, 20], <i>cosine</i> [2800, 3300, 20], <i>sine</i>	[0.174023, 0, 0, 0.791447, 0, 0, 0.0345301, 0, 0, 0]
26	100 + 50	[2800, 3300, 5], <i>cosine</i> [2800, 3300, 10], <i>sine</i>	[0, 0, 0.753209, 0, 0.146791, 0, 0.1, 0, 0, 0]
27	100 + 84	[2800, 3300, 5], <i>cosine</i> [2800, 3300, 6], <i>sine</i>	[0.174023, 0, 0, 0.791447, 0, 0, 0.0345301, 0, 0, 0]
28	100 + 100	[2800, 3300, 5], <i>cosine</i> [2800, 3300, 5], <i>sine</i>	[0.1, 0, 0.5, 0, 0.4, 0, 0, 0, 0, 0]

Table 3.5: Constraint study based on Experiment 5 of toy model

### 3.6 Constraint Study Based on Experiment 5

Based on Experiment 5, similar constraint study is conducted as displayed in Table 3.5, and the same observation is obtained as the experiments in Table 3.4. When the result composition  $[0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0]$  and target one  $[0.1, 0, 0.5, 0, 0.4, 0, 0, 0, 0, 0]$  are used to plot the spectra, the produced spectra are almost identical as shown in Figure 3.6.

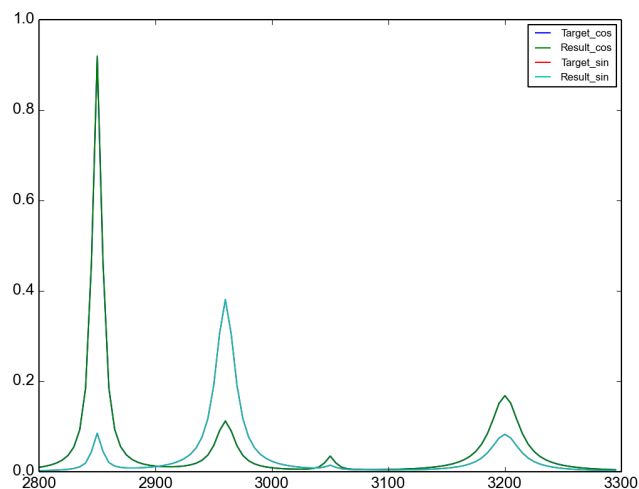


Figure 3.6: IR spectra plotted by the return compositions from the constraint study based on Experiment 5 of toy model

### 3.7 Discussion and Conclusion

Recall that our LP model, for the right data set is expected to return the target composition. We can conclude that, if the target composition is not returned correctly, then the data we collect is not sufficient to describe the experiments to our LP model.

However, when the target composition is not returned correctly, the return composition does build spectra that are almost identical to the target ones. This means that there are more than one composition can build the spectra that are almost identical to the target ones. Because of the numerical limitation, an unique optimum solution is always obtained.

The above conclusion leaves us a new question: how do we know there is sufficient spectral information in order to obtain the target composition of candidates at interfaces? To answer this question, further experiments are conducted by applying the spectral information of real molecules to the LP model. The goal is to investigate with all the spectral information we can obtain for real molecules, can the LP model return the target composition of candidates at interfaces. If this goal can be achieved, can this approach be applied systematically to various circumstances?

## Chapter 4

# Real Molecular Model

### 4.1 Description

After experimenting with the toy problem, lacking sufficient spectral information is the key cause for the failure of obtaining the correct target composition. First of all, in the toy model, there are only four vibrational modes, and the range of the wavenumber is limited. Therefore, the number of data points selected is limited. Secondly, the similarity among the candidates is high, as all the candidates are coming from one same molecule. Third, only IR spectra is considered.

In this chapter, experiments are conducted using real molecules. In addition to IR, both Raman and SFG are introduced to the real molecules, which makes the study one step closer to the overall goal and scope. The real molecule focused on this chapter is Methionine amino acid.

Same as the toy model, in order to limit the possible candidate space of Methionine, *twist* and *azimuthal* angular distributions are assumed to be isotropic, which are integrated. Only  $\theta$  in Euler angles is considered in Methionine's surface orientation distribution function. In Chapter 2 section Generating model spectra, how a molecule's IR, Raman and SFG spectra are generated have been explained. Two unique IR spectra can be obtained from  $x$ , and  $z$  polarizations. Four unique Raman spectra can be obtained from  $xx$ ,  $xy$ ,  $xz$  and  $zz$  polarizations. Three unique SFG spectra can be obtained from  $yyz$ ,  $zyz$  and  $zzz$  polarizations.

Experiment #	1	2	3	4
# Candidates	4	4	4	4
Candidates	[0, 20, 40, 60]	[0, 20, 40, 60]	[0, 20, 40, 60]	[0, 20, 40, 60]
Target Composition	[0.1, 0.5, 0.4, 0]	[0.1, 0.5, 0.4, 0]	[0.1, 0.5, 0.4, 0]	[0.1, 0.5, 0.4, 0]
# Data Points	200(irx)	200(irz)	200(irx) + 200(irz)	200(irx) + 200(ramanxx)
Return Composition	[0.701654, 0, 0, 0.298346]	[0.701654, 0, 0, 0.298346]	[0.701654, 0, 0, 0.298346]	[0.1, 0.5, 0.4, 0]

Table 4.1: Experiment 1 to Experiment 4 setting for methionine candidates

The goal is to see if those spectral information is sufficient for the LP model to return the correct target composition of the candidates of one type molecule at interfaces. If yes, we need to figure out which spectral information is needed for the LP model. If no, we need to check if the cause of the failure is the same as the toy model.

## 4.2 Experiments

Table 4.1, four experiments are set up with four candidates and one same target composition. These four candidates each has  $\theta$  of the following degree:  $0^\circ$ ,  $20^\circ$ ,  $40^\circ$  and  $60^\circ$ . The only difference among these four experiments is the spectroscopy information we select to construct the LP model, and it is indicated by the Number of Data Points. In Experiment 1, only IR  $x$ -polarization spectral information is used. This means that only data points from IR  $x$ -polarization are selected to build the LP model. Same for Experiment 2, data points are obtained from spectra of IR’s  $z$ -polarization. In Experiment 3, the spectral information of IR’s  $x$  and  $z$ -polarizations are combined. At last, in Experiment 4, spectral information of IR  $x$ -polarization and Raman  $xx$ -polarization are combined. The LP model we build for each experiment is different as the data points are selected differently. As the return composition indicates, Experiment 4 contains the most abundant information, as its return composition matches to the target one.

When merely using IR information, the return composition is the same for Experiment 1, 2 and 3. Figure 4.1 displays the resulting spectra generated by using the return composition obtained from the first three experiments. The resulting spectra is almost identical to the target ones. It indicates that with only IR spectral information is not sufficient to get the target composition. However, the return composition could perfectly re-product the target spectra. This means that further information is needed to build the constraints of the LP model. The more constraints are intro-



duced, the more accurate the return composition will be.

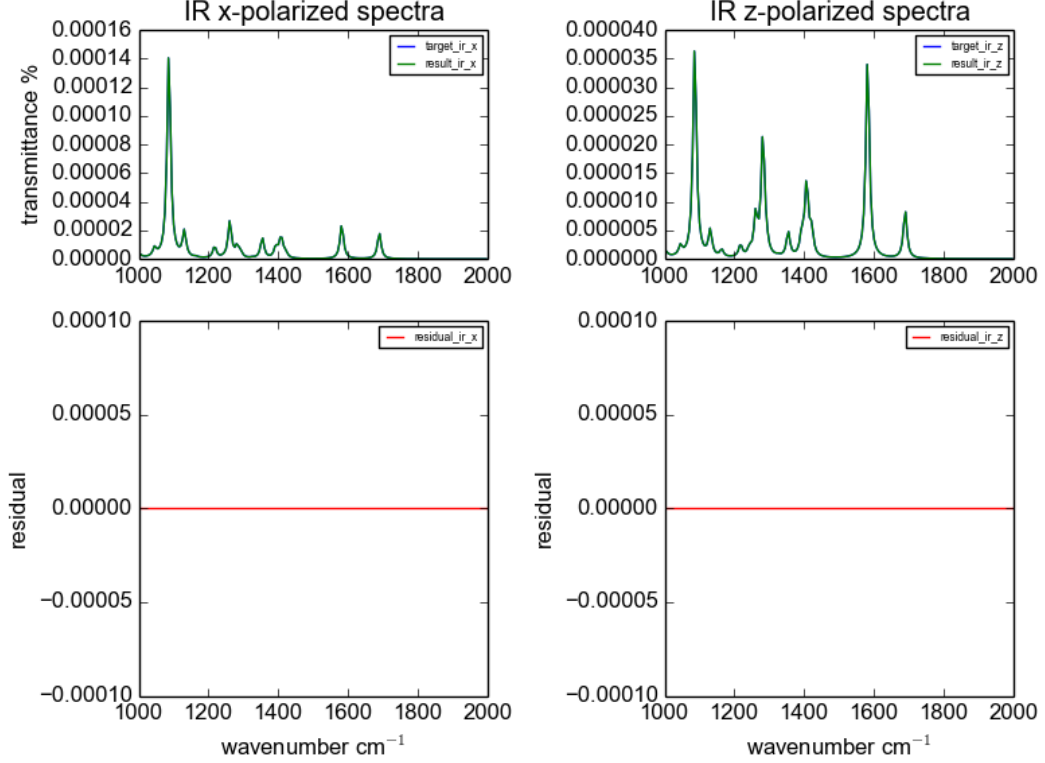


Figure 4.1: Compare target IR spectra with the ones generated by the return composition of Experiment 1, 2 and 3

In Experiment 4, the LP model that constructed from combining IR and Raman spectral information is sufficient to obtain the target composition. When the difference in  $\theta$  degree for candidates decreases from  $20^\circ$  to  $10^\circ$ . Checking if Raman and IR together is still sufficient enough to derive the target composition is desired. Therefore, the following experiments are conducted as shown in Table 4.2.

Experiment 5 shows that the LP model constructed by merely using IR spectral information is not sufficient to derive the target composition. Experiment 6 indicates that combining IR and Raman spectral information helps to derive the target composition. What's more, Experiment 7, 8 and 9, illustrate that Raman spectral information itself is sufficient to obtain the target composition as well.

For experiment setting in Table 4.1 and Table 4.2, combining IR and Raman

# Candidates	4	
Candidates	[0, 10, 20, 30]	
Target Composition	[0.1, 0.5, 0.4, 0]	
Experiment index	# Data Points	Result Composition
5	200(irx) 200(irz)	[0.752528, 0, 0, 0.247472]
6	200(irx) 200(irz) 200(ramanxx)	[0.1, 0.5, 0.4, 0]
7	200(ramanxx) 200(ramanxy) 200(ramanxz)	[0.1, 0.5, 0.4, 0]
8	200(ramanxx) 200(ramanxy) 200(ramanzz)	[0.1, 0.5, 0.4, 0]
9	200(ramanxx) 200(ramanxy) 200(ramanxz) 200(ramanzz)	[0.1, 0.5, 0.4, 0]

Table 4.2: Experiment 5 to Experiment 9 setting for methionine candidates

spectral information to construct a LP model is sufficient enough to obtain the target composition. In order to study the limitation of the LP model, the complexity of the experiment setting needed to be increased. Therefore, another group of experiments are designed as shown in Table 4.3. There are 5 candidates included in the experiments. Each candidate has  $\theta$  with the following degree:  $0^\circ$ ,  $10^\circ$ ,  $20^\circ$ ,  $30^\circ$  and  $40^\circ$ . The target composition is more complex than previous experiments, each candidate takes 20% in the mixture.

Experiment 10 uses only IR spectral information to construct the LP model, and the return composition does not match the target one. Experiment 11 uses only Raman spectral information, and the return composition does not match to the target neither. Same for Experiment 12 that uses only SFG spectral information. From Experiment 13, different kinds of spectral information are combined. In Experiment 13, IR and Raman spectral information is used to produce the LP model, still the return composition is different from the target one. Experiment 14 combines Raman and SFG, Experiment 15 uses IR and SFG, Experiment 16 cooperates all the three spectral information, however, none of them returns a composition that matches the target one.

The results of Experiment 10 to 16 indicate that even combining all the spectral information of IR, Raman and SFG, it is still not sufficient to attain the target composition for the experiments set up in Table 4.3. Our LP model is showing its limitation in these experiments. In order to confirm the reason causing the LP model to returning the target composition is because of insufficient information, further experiments are conducted in Table 4.4.

### 4.3 Experiments to Explain the Limitation of LP Model for Methionine Molecule

In order to further explore the reason that LP model reaches its limitation for the real molecule, Experiment 17 and 18 are conducted. Methionine candidates are still used. To make the study case more general than Experiment 1 to 16, candidates'  $\theta$  values are expanded from  $0^\circ$  to  $80^\circ$ . In total, there are 9 candidates. Because the SFG spec-

Number of Candidates	5	
Candidates	[0, 10, 20, 30, 40]	
Target Composition	[0.2, 0.2, 0.2, 0.2, 0.2]	
Experiment index	Constraints	Result
10	200(irx) 200(irz)	[0.607766, 0, 0, 0, 0.392234]
11	200(ramanxx) 200(ramanxy) 200(ramanxz) 200(ramanzz)	[0.247792, 0, 0.502139, 0, 0.250069]
12	200(sfgyyz) 200(sfgzyz) 200(sfgzzz)	[0.321014, 0, 0.31018, 0.163041, 0.205764]
13	200(irx) 200(irz) 200(ramanxx) 200(ramanxy) 200(ramanxz) 200(ramanzz)	[0.247792, 0, 0.502139, 0, 0.250069]
14	200(ramanxx) 200(ramanxy) 200(ramanxz) 200(ramanzz) 200(sfgyyz) 200(sfgzyz) 200(sfgzzz)	[0.321014, 0, 0.31018, 0.163041, 0.205764]
15	200(irx) 200(irz) 200(sfgyyz) 200(sfgzyz) 200(sfgzzz)	[0.321014, 0, 0.31018, 0.163041, 0.205764]
16	200(irx) 200(irz) 200(ramanxx) 200(ramanxy) 200(ramanxz) 200(ramanzz) 200(sfgyyz) 200(sfgzyz) 200(sfgzzz)	[0.321014, 0, 0.31018, 0.163041, 0.205764]

Table 4.3: Experiment 10 to Experiment 16 setting for methionine candidates

# Candidates	9	
Candidates	[0, 10, 20, 30, 40, 50, 60, 70, 80]	
Target Composition	[0.2201, 0.28905, 0.05201, 0.08251, 0.35633, 0, 0, 0, 0]	
Experiment #	# of Data Points	Result Composition
17	each 5 wavenumber of IR, Raman and SFG spectra	[0.158921, 0.388434, 0.0, 0.0985466, 0.354099, 0.0, 0.0, 0.0, 0.0]
18	each 500 wavenumber of IR, Raman and SFG spectra	[0.397991, 0.0, 0.203394, 0.0357663, 0.362848, 0.0, 0.0, 0.0, 0.0]

Table 4.4: Experiment 17 and 18 to explain the limitation of our LP model for methionine molecule

tra for  $\theta$  of  $90^\circ$  is a straight line, it is excluded from all the experiments related to real molecules. For target composition, five candidates are randomly selected to be presented. The difference between Experiment 17 and 18 is that different amount of data points are selected to build the LP model. From all three spectroscopy techniques' spectral information, every 5 wavenumber a data point is selected for Experiment 17. Every 500 wavenumber a data point is selected for Experiment 18. As a result, Experiment 17 and 18 each returns a different composition. Both compositions do not match to the target one.

However, for both Experiment 17 and 18, when the return composition is used to generate the IR, Raman and SFG spectra, then these spectra are plotted together with the spectra created by the target composition. All the spectra are almost identical for IR, Raman and SFG. Figure 4.2, 4.3 and 4.4 display the spectra plotted by using the return composition and the target one of Experiment 17. Every spectrum is almost identical to each other as shown in the figures. Same for Experiment 18 as shown in Figure 4.5, 4.6 and 4.7. These figures prove again that there are more than one composition that can perfectly construct the target spectra. The data information used to construct the LP model is not sufficient to converge to the return composition exactly matches to the target one. This conclusion exactly fits the result obtained from the experiments we have done with the toy model.

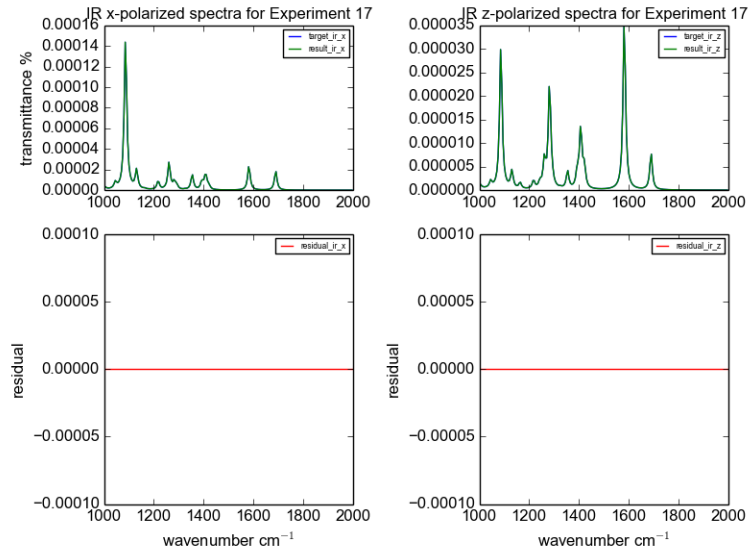


Figure 4.2: IR spectra plotted by using target composition and return composition of Experiment 17

## 4.4 Conclusion

## 4.5 Extra Experiments

TODO: this part of experiments are similar as what are done in Chapter 5 and 6. Think how to involve this part properly.

From Experiment 1 to 18, LP model helps to return the target composition for some cases, and not for others. We want to figure out if there a clean line indicating the information used to generate the LP model is not sufficient to obtain the target composition for one molecule. In order to answer this question, more systematic experiments needed to be organized. Therefore, the following experiments are conducted. The Methionine candidate space is the same as Experiment 17 and 18. spread from  $0^\circ$  to  $80^\circ$  on  $\theta$ .

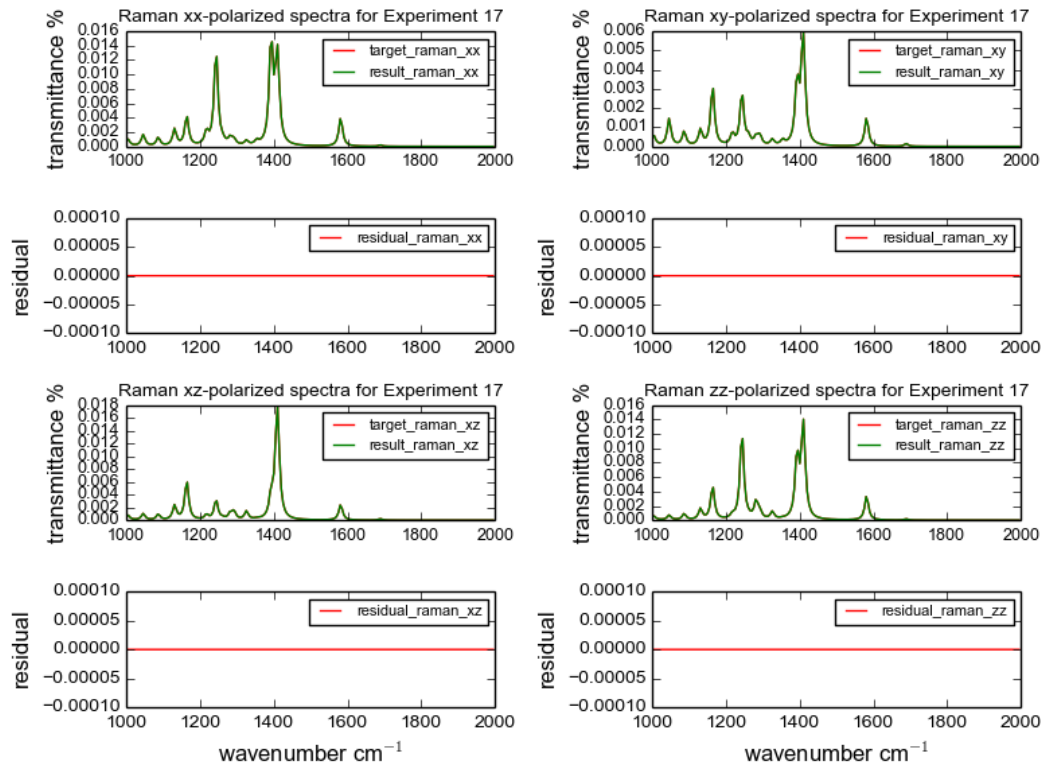


Figure 4.3: Raman spectra plotted by using the target composition and the return composition of Experiment 17

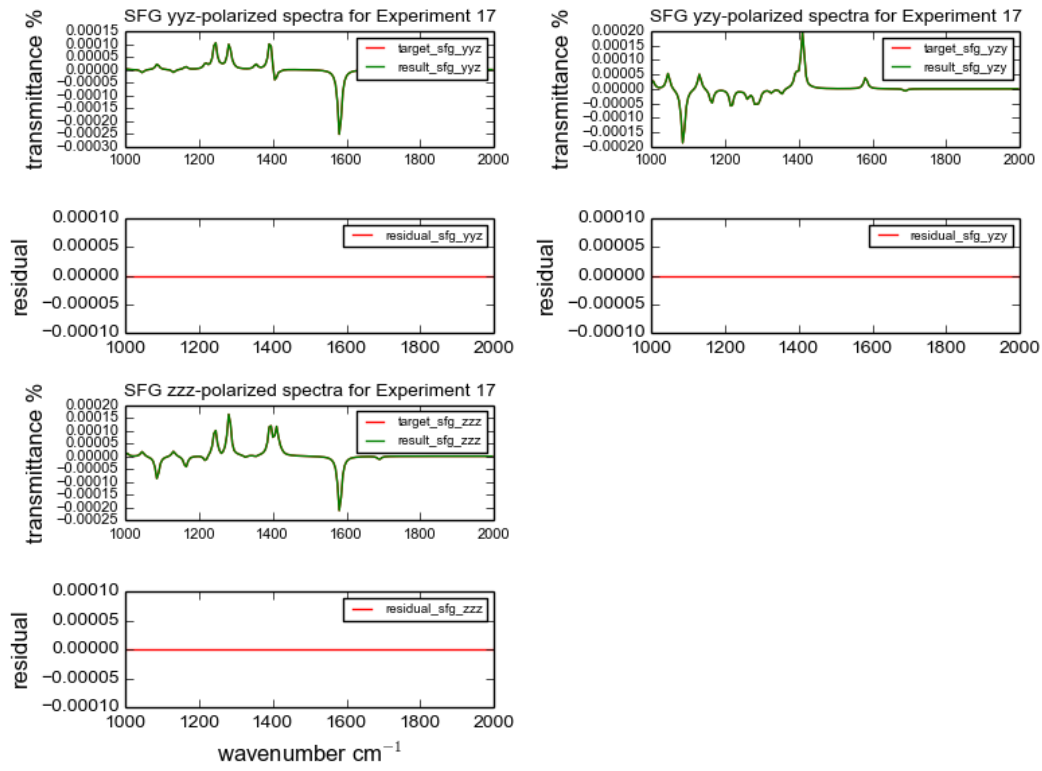


Figure 4.4: SFG spectra plotted by using the target composition and the return composition of Experiment 17



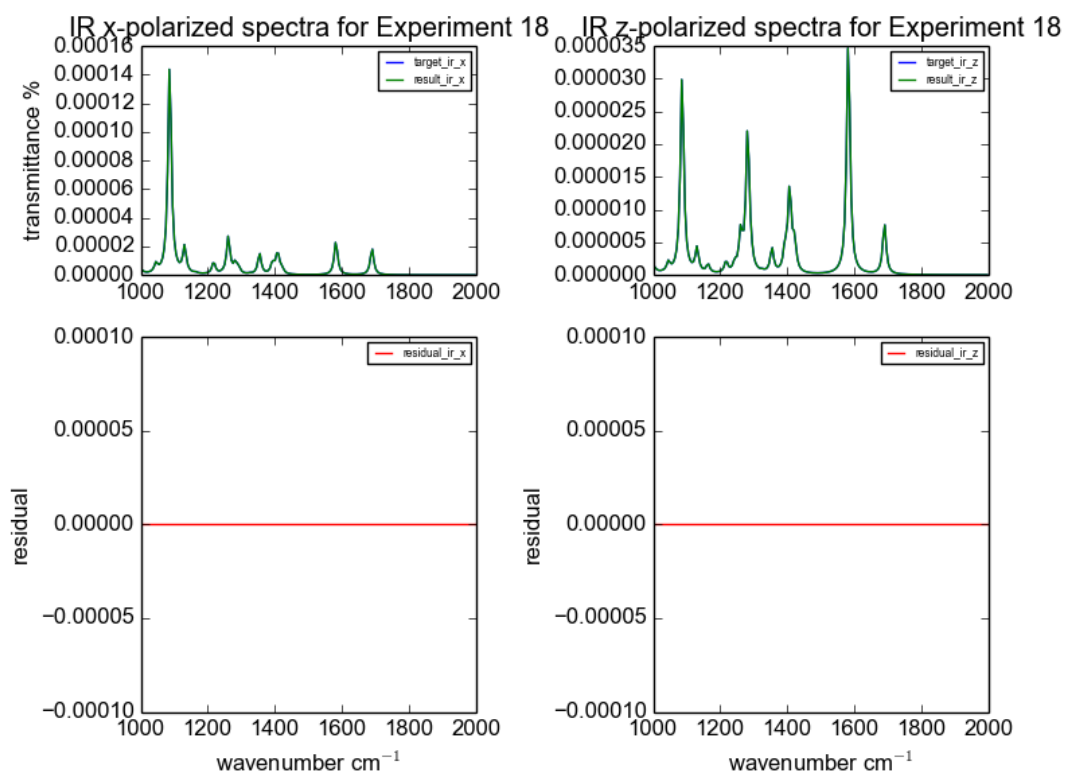


Figure 4.5: IR spectra plotted by using the target composition and the return composition of Experiment 18

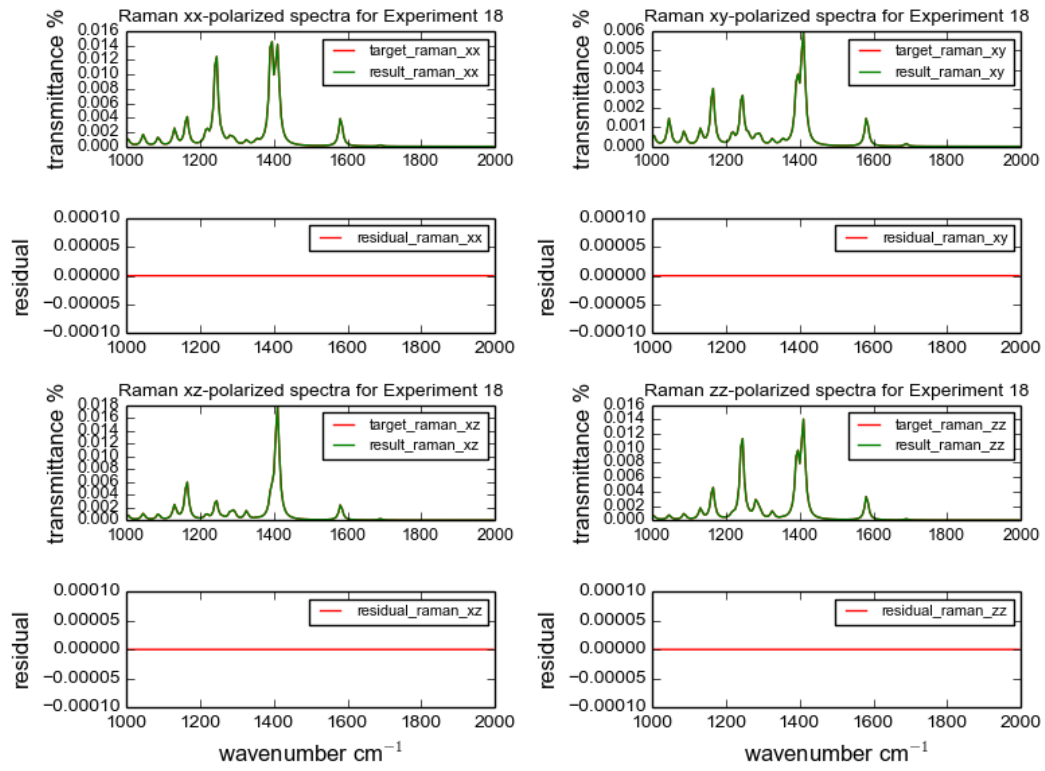


Figure 4.6: Raman spectra plotted by using the target composition and the return composition of Experiment 18

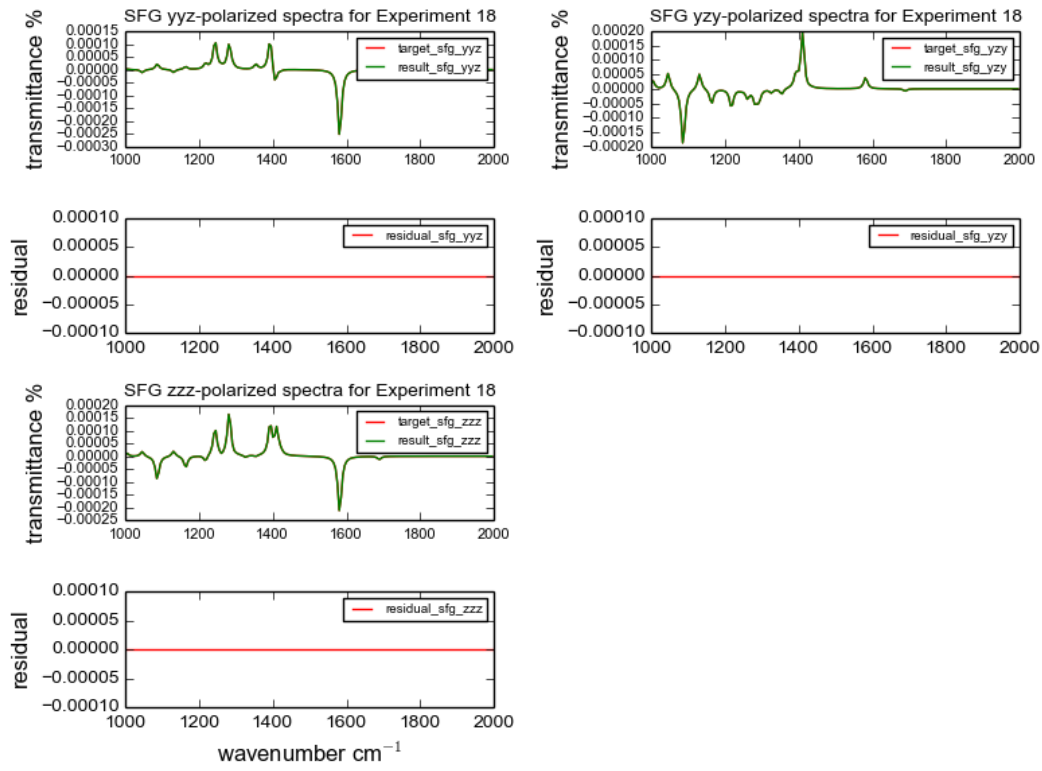


Figure 4.7: SFG spectra plotted by using the target composition and the return composition of Experiment 18

## Chapter 5

# Mixture of Molecules

### 5.1 Description

In Chapter 4, experiments indicate that for one type of molecule at interfaces, even combining all the three spectral information, the constructed LP model cannot return the target composition in most cases. The existing spectral information is not adequate to obtain the target composition of one type of molecule at interfaces. Multiple return compositions can build the spectra that are almost exactly the same as the target ones. These compositions are returned by the LP models that use different amounts of spectral information. This indicates that even extracting data information from three spectroscopy techniques, it is still not sufficient to obtain the target composition for one type of molecule at interfaces. Besides one type of molecule at interfaces, we are also interested in the case where different molecules at interfaces. For a mixture of different molecules at interfaces, we want to figure out whether our LP model can help to obtain the target composition. If the LP model success in obtaining the target composition with certain spectral information, we want to know which spectral information. Moreover, we want to know the efficiency of the spectral information in obtaining the target composition.

## 5.2 Experiments

### 5.2.1 Experiments Considering Each Amino Acid Candidates from $0^\circ$ to $80^\circ$ on $\theta$ in the Mixture

To achieve the study of the coordination distribution of various molecules at interfaces, further experiments are constructed. These experiments have the following common settings.

First, there are six different amino-acids in the mixture: methionine, leucine, isoleucine(ile), alanine, threonine and valine. For each amino acid, only  $\theta$  difference is considered, the other two Euler angles are integrated. Each amino acid molecule has 9 candidates in the mixture, they have  $\theta$  of the following values:  $0^\circ$ ,  $10^\circ$ ,  $20^\circ$ ,  $30^\circ$ ,  $40^\circ$ ,  $50^\circ$ ,  $60^\circ$ ,  $70^\circ$  and  $80^\circ$ . Because when  $\theta$  equals  $90^\circ$ , the SFG spectra is a straight line. The corresponding candidate is excluded from all the experiments. As a result, there are 54 candidates in the mixture.

Second, the target composition need to be generated. The operation includes two steps: randomly pick one candidate from each amino acid's 9 candidates, then randomly generate a percentage for the selected candidate. The target composition is made of six randomly selected candidates with assigned percentage coming from the amino acids. The rest 48 candidates have 0 percentage in the target composition. Namely, six selected candidate makes 100% component of the mixture.

Third, the IR, Raman and SFG spectra need to be generated for all the 54 candidates and the target.

Table 5.1 displays a set of experiments each experiment contains different spectral information. In Experiment 1, candidates  $x$ - and  $z$ -polarized IR spectra are obtained. The target's IR spectra are generated by the dot product of the target composition and all the candidates' spectral data. Then the corresponding LP model is conducted using Equation 3.4. Therefore, we claim that the LP model in Experiment 1 only contains IR information.

Similarly, Experiment 2 contains only Raman spectral information of the follow-

Experiment Index	Spectral Information
Experiment 1	$x$ and $z$ polarized IR spectra
Experiment 2	$xx$ , $xy$ , $xz$ and $zz$ polarized Raman spectra
Experiment 3	$yyz$ , $zyy$ and $zzz$ polarized SFG spectra
Experiment 4	$x$ and $z$ polarized IR spectra $xx$ , $xy$ , $xz$ and $zz$ polarized Raman spectra
Experiment 5	$x$ and $z$ polarized IR spectra $yyz$ , $zyy$ and $zzz$ polarized SFG spectra
Experiment 6	$xx$ , $xy$ , $xz$ and $zz$ polarized Raman spectra $yyz$ , $zyy$ and $zzz$ polarized SFG spectra
Experiment 7	$x$ and $z$ polarized IR spectra $xx$ , $xy$ , $xz$ and $zz$ polarized Raman spectra $yyz$ , $zyy$ and $zzz$ polarized SFG spectra

Table 5.1: Detailed experiment set setting for the mixture of amino acids

ing four polarizations:  $xx$ ,  $xy$ ,  $xz$  and  $zz$ . Experiment 3 contains only SFG spectral information of  $yyz$ ,  $zyy$  and  $zzz$  three polarizations.

Starting from Experiment 4, spectral information of different spectroscopy techniques are combined. In Experiment 4, IR spectral information is combined with Raman. In Experiment 5, IR spectral information is combined with SFG. In Experiment 6, Raman and SFG spectral information are incorporated. At the end, in Experiment 7, all three spectral information are put together: IR, Raman and SFG.

Finally, this experiment set is run 100 times in order to see which experiment in the set returns the target composition with the highest accuracy. This accuracy is measured by the time of each experiment returns the target composition. The scoring mechanism to measure whether a return composition matches to the target one is described in the next section.

### 5.2.2 Scoring method

At the first glance, the sum of residuals between the spectra composed by the return composition and the target one can be used to measure the accuracy of the return composition. However, in most experiments conducted earlier, the spectra generated by the return composition are almost identical to the ones created by the target composition. The sum of residuals between these spectra is negligible, which makes it appropriate to use as a scoring criteria.

Another way to measure the accuracy of the return composition, is to compare it directly with the target one. Calculating the sum of the residuals between a target composition and a return one directly can be a fast approach to evaluate the accuracy of each experiment. The shortage of this approach is that it cannot be used to measure in real experiments where the target composition is unknown. However, in the current experiments, this approach can be a way to evaluate the return composition for all the experiments where the target compositions are known in advance.

The return composition of each experiment in the set is obtained for each run. Each return composition is compared with the target one to calculate the sum of

the residuals. If the sum is smaller than a certain threshold, which is  $10^{-7}$ , then the return composition is considered to be the same as the target one.

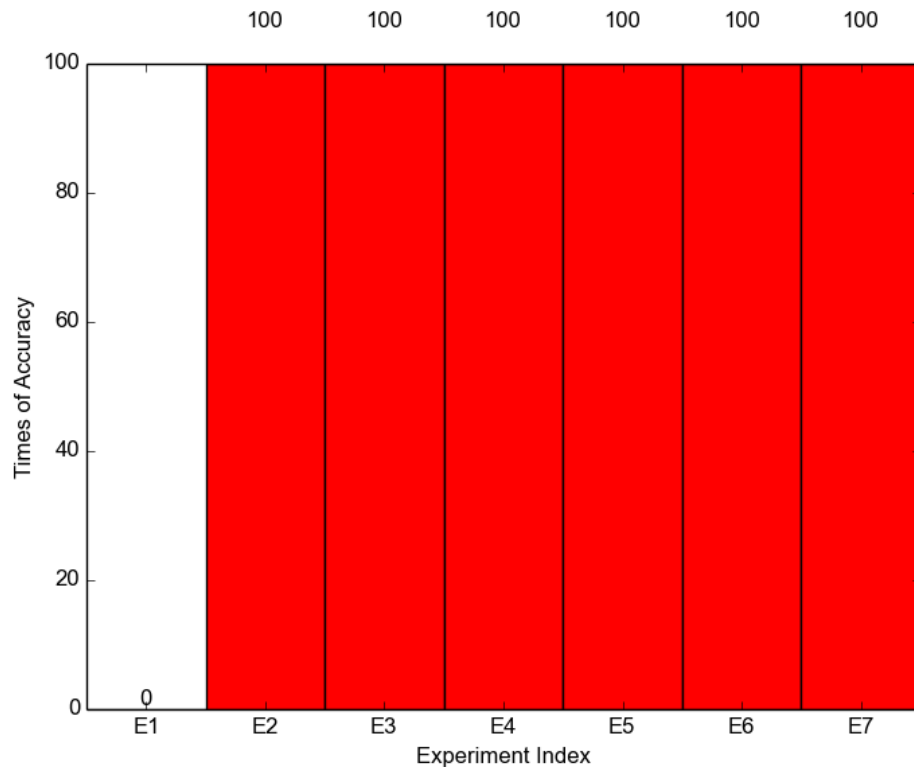


Figure 5.1: Accuracy analysis for experiments considering a mixture of amino acids with candidates from  $0^\circ$  to  $80^\circ$  on  $\theta$  for each amino acid. Accuracy indicates how many times each experiment in the set return a composition matches the target one.

The experiment set is run 100 times based on the scoring method, the result is shown in Figure 5.1. Experiment 2, 3, 4, 5, 6 and 7, all return the target composition in the 100 runs. This result indicates that Raman or SFG alone is sufficient to obtain the target composition. For a mixture of amino acids with candidates from  $0^\circ$  to  $80^\circ$  on  $\theta$  for each amino acid. Any experiments that contain Raman and SFG result in the same accuracy.

The only exception is Experiment 1. The accuracy is fairly low, which indicates that IR spectra do not contain sufficient information in order to obtain the target composition.



To take a further insight into the return composition of Experiment 1, the experiment set is re-run 100 times, only the return composition of Experiment 1 is analyzed and focused. In each run, IR  $x$ - and  $z$ -polarized spectra are plotted both by the returned composition and the target one. The result is that these spectra conducted by the two different compositions are very close to each other in each run. Randomly take one run as an example, Figure 5.2 displays the plotted spectra, and they are almost identical to each other. The residual is very small for the data points where these two spectra are not overlapped. This indicates that the optimum composition returned by the LP model conducted with only IR spectral information has achieved its best in obtaining a composition that best fit the target spectra.

(TODO: rewrite or remove this paragraph) Comparatively, SFG has three unique polarizations, and Raman has four unique polarizations. From each projection’s spectrum, we evenly select 200 data points. This means that one more projection will bring in 200 more constraints or 400 more (when we take the absolute sign off) constraints to the LP model. This would make a huge difference in the LP model, in term of further refining the candidate selection in target composition. However, it is still too early for us to say that Raman has more coordination information because it has four unique polarizations. Because for Raman’s any polarization, the spectrum of candidate with  $\theta$  equals to one degree is identical to the one of candidate with this  $\theta$  degree’s complementary. For example, the Raman spectra for candidate with  $\theta$  of  $10^\circ$ , is the same as candidate with  $\theta$  of  $170^\circ$ . And for IR, it is the same case. Only SFG tells the differences between these two degrees, as the spectra for candidate with  $\theta$  of one degree is symmetric to its complementary along wavenumber as shown in Figure 5.8.

### 5.2.3 Experiments Considering Each Amino Acid Candidates from $0^\circ$ to $180^\circ$ on $\theta$ in the Mixture

To further study the capacity of the LP models built for the mixture of molecules, the candidate pool is expanded from  $0^\circ$  to  $180^\circ$  in terms of the  $\theta$  value. Therefore, each amino acid has 18 candidates. In total, there are 108 candidates in the mixture. The same set of experiments in Table 5.1 is used. The only difference is instead of randomly select one candidate from 9 candidates, it is selected from 9. All 108 candidates’ IR, Raman and SFG spectra need to be generated. Figure 5.3 illustrates

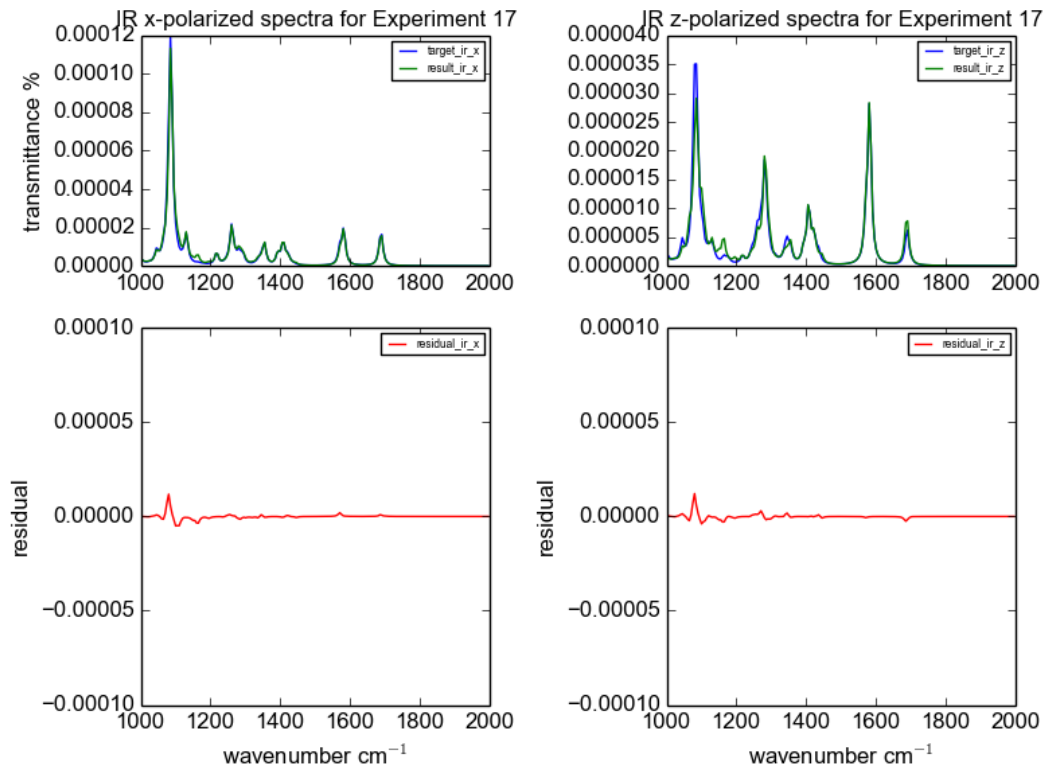


Figure 5.2: IR Spectra Plotted by Result Composition and Target Composition.

the results obtained in 100 runs. The accuracy in Experiment 1 is still low. This is not surprising as the complexity of the candidates has increased. Moreover, IR spectra for candidate with  $\theta$  of one degree is identical to the one with  $\theta$  of this degree's complementary, as shown in Figure 5.7. This also increases the difficulty for the LP model using IR spectral information to return the target composition.

However, it should be noticed that the accuracy for Experiment 2 has dramatically dropped. This can be caused by the Raman spectra for one candidate with a  $\theta$  is identical to the one of this  $\theta$  value's complementary as displayed in Figure 5.8.

In Figure 5.3, the accuracy for Experiment 3 is no longer high neither. After increasing the number of amino acid candidates from 9 to 18, the complexity of the corresponding LP model has increased. Although the added candidates' SFG spectra are symmetric along wavenumber which may greatly increase the uniqueness of the candidates as shown in Figure 5.9. The SFG spectral information is still insufficient

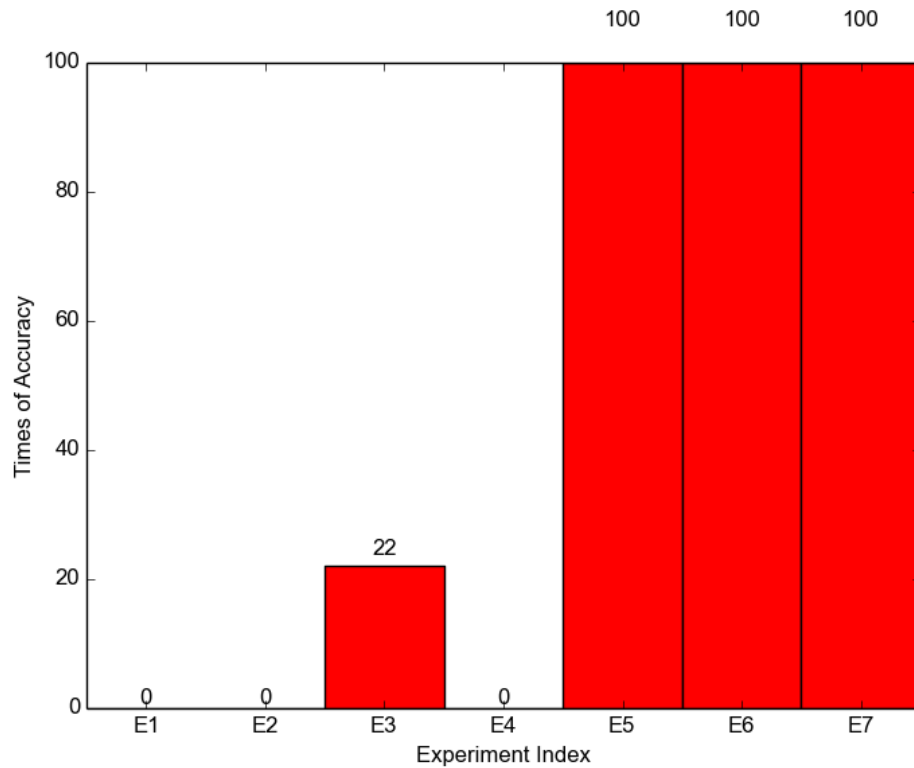


Figure 5.3: Accuracy analysis for experiments considering a mixture of amino acids with candidates from  $0^\circ$  to  $180^\circ$  on  $\theta$  for each amino acid. Accuracy indicates how many times each experiment in the set return a composition matches the target one.

to obtain the target composition.

The good result starts to emerge when using the combinations of IR and SFG or Raman and SFG. Figure 5.3 shows that Experiment 5, 6, and 7 all have 100% accuracies. This phenomenon can be explained as follow: SFG helps to distinguish a candidate from its complementary on  $\theta$  value. The extra spectral information coming from IR or Raman helps to further refine the LP model, which can then converge the return composition to the target one.

Although the accuracy in Experiment 2 is low. There is still some noticable result in the return composition: for each amino acid, the percentage assigned is correct; however, the candidate presented may be always be correct. It is either the correct one, or the correct one's complementary. Randomly select one experiment run as an example, Figure 5.4 displays the target composition. Figure 5.5 displays the return

composition of Experiment 2 in the run. Figure 5.6 is the return composition of Experiment 6. Figure 5.4 and 5.6 are identical, which means the return composition of Experiment 6 is the same as the target one. The values in Figure 5.5 are the same as Figure 5.4. However, the position of each value is not the same in two the figures. For example, the percentage value 0.299586 of methionine is for  $\theta = 120^\circ$  in Figure 5.4, but is for  $\theta = 60^\circ$  in Figure 5.5. These two angles are complementary. This observation is the same for isoleucine, alanine, threonine, and valine in the figure. This observation is a general case across all the runs of the experiment set. The return composition of Experiment 6 matches the target one. However, the return composition of Experiment 2 fails to pick each amino acid's correct candidate from this candidate's complementary. This can be explained as the Raman spectra for one  $\theta$  are the same as its complementary.

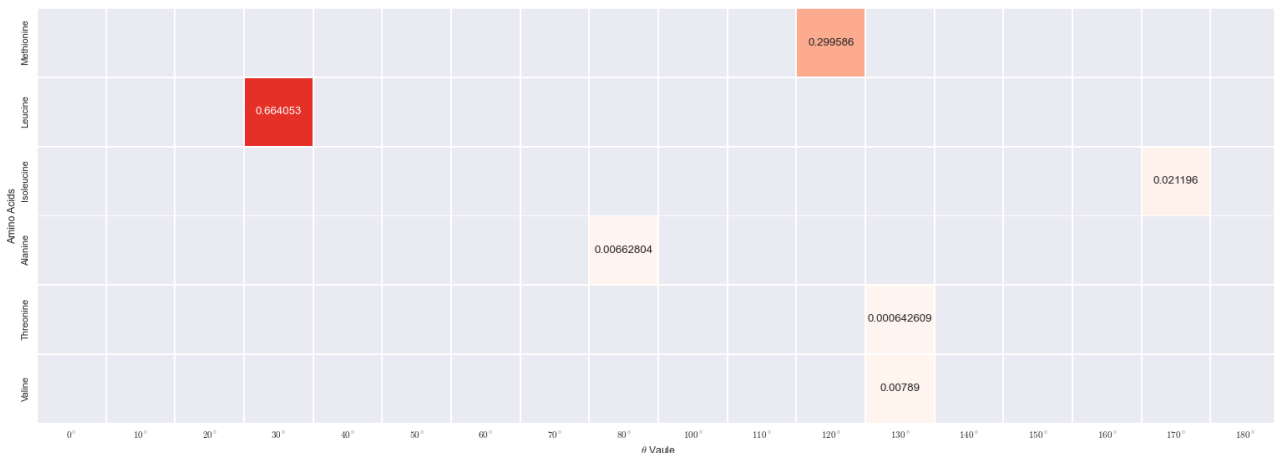


Figure 5.4: Target composition of one random run of six mixed amino acids with candidates expanded from  $0^\circ$  to  $180^\circ$  on  $\theta$  for each amino acid. More detailed data of this target composition can be found in A.1 in the Appendix.

The return composition of Experiment 4 is the same as the one of Experiment 2, which means combining IR spectra information with Raman is not sufficient for this experiments setting. This is because the IR spectra for one  $\theta$  degree are also the same as its complementary. Spectral information from SFG is needed in order to study the cases that having  $\theta$  expanded from  $0^\circ$  to  $180^\circ$ .

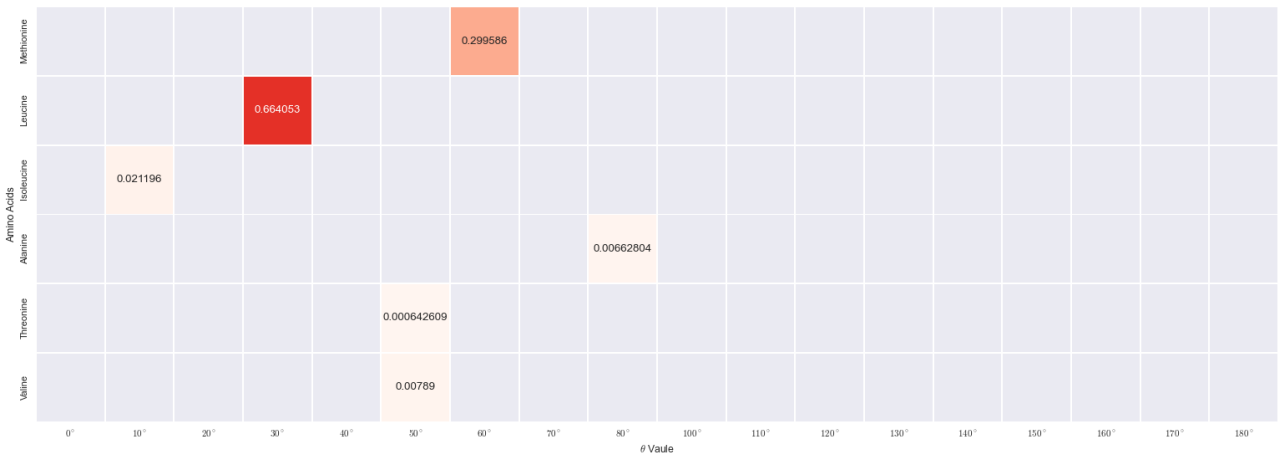


Figure 5.5: Return composition of experiment 2 for one random run of six mixed amino acids with candidates expanded from  $0^\circ$  to  $180^\circ$  on  $\theta$ . More detailed data of this target composition can be found in A.2 in the Appendix.

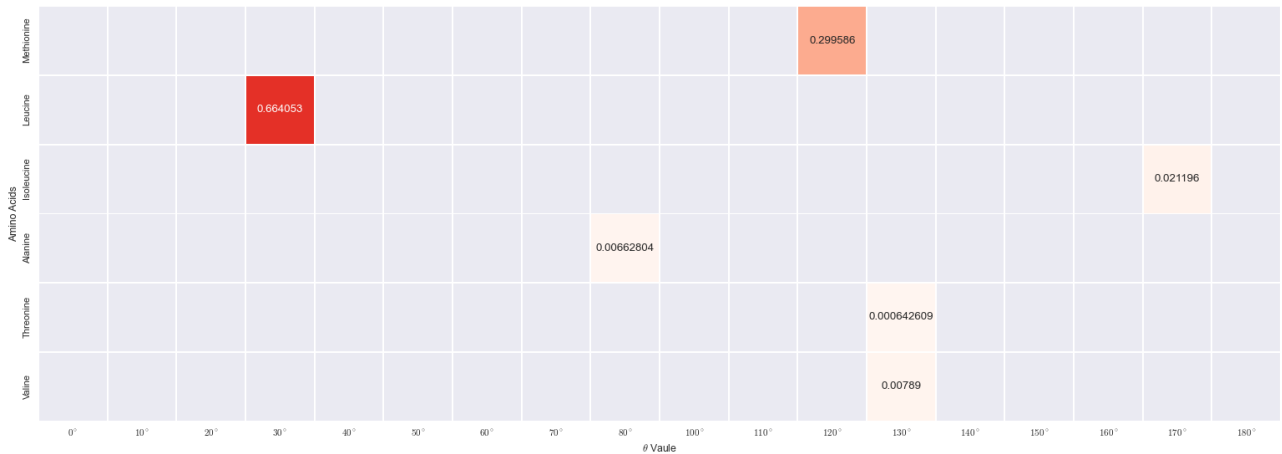


Figure 5.6: Return composition of experiment 6 for one random run of six mixed amino acids with candidates expanded from  $0^\circ$  to  $180^\circ$  on  $\theta$ . More detailed data of this target composition can be found in A.3 in the Appendix.

### 5.3 Conclusion

Raman and SFG spectral information alone is sufficient to obtain the target composition, when considering a mixture of amino acids with candidates expanded from  $0^\circ$  to  $80^\circ$  on  $\theta$  for each amino acid.

When the candidates are expanded from  $0^\circ$  to  $180^\circ$  on  $\theta$ , SFG spectral information

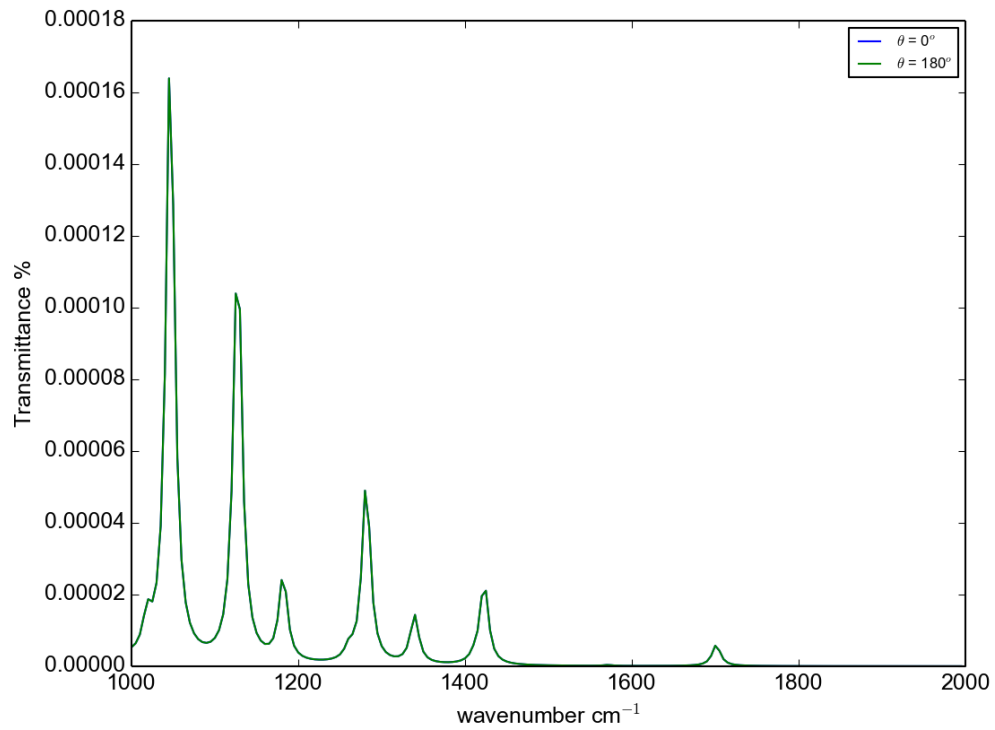


Figure 5.7: IR  $z$  projection spectrum for alanine candidate with  $\theta$  of  $0^\circ$  is identical to alanine candidate with  $\theta$  of  $180^\circ$

needs to combine with IR or Raman in order to obtain the target composition.

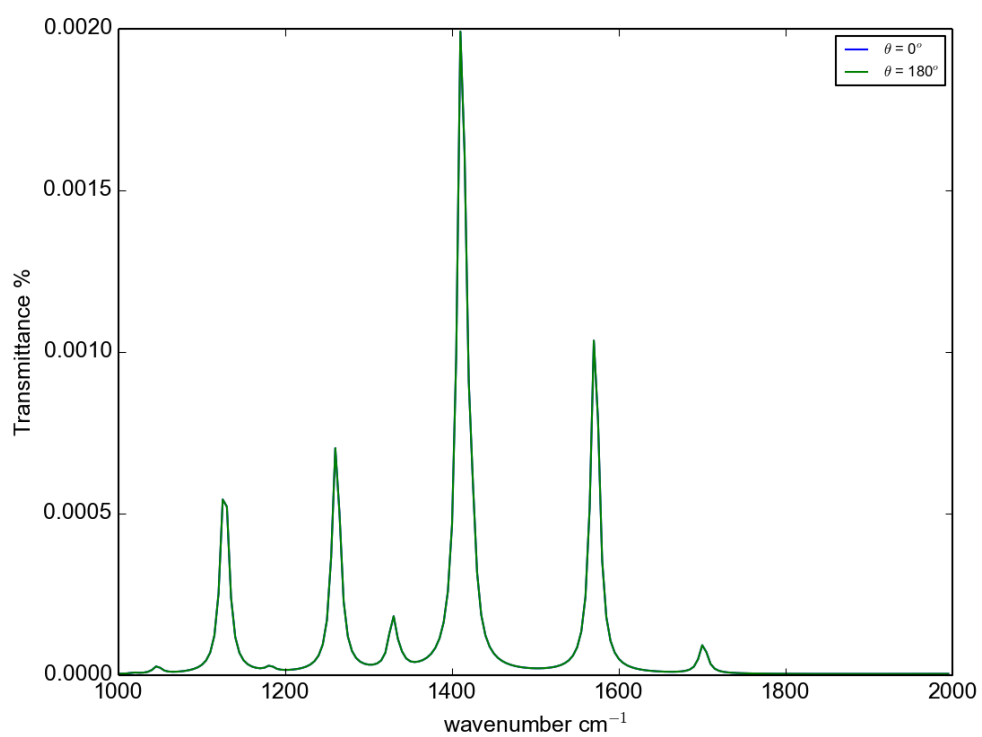


Figure 5.8: Raman  $zz$  projection spectrum for alanine candidate with  $\theta$  of  $0^\circ$  is identical to alanine candidate with  $\theta$  of  $180^\circ$

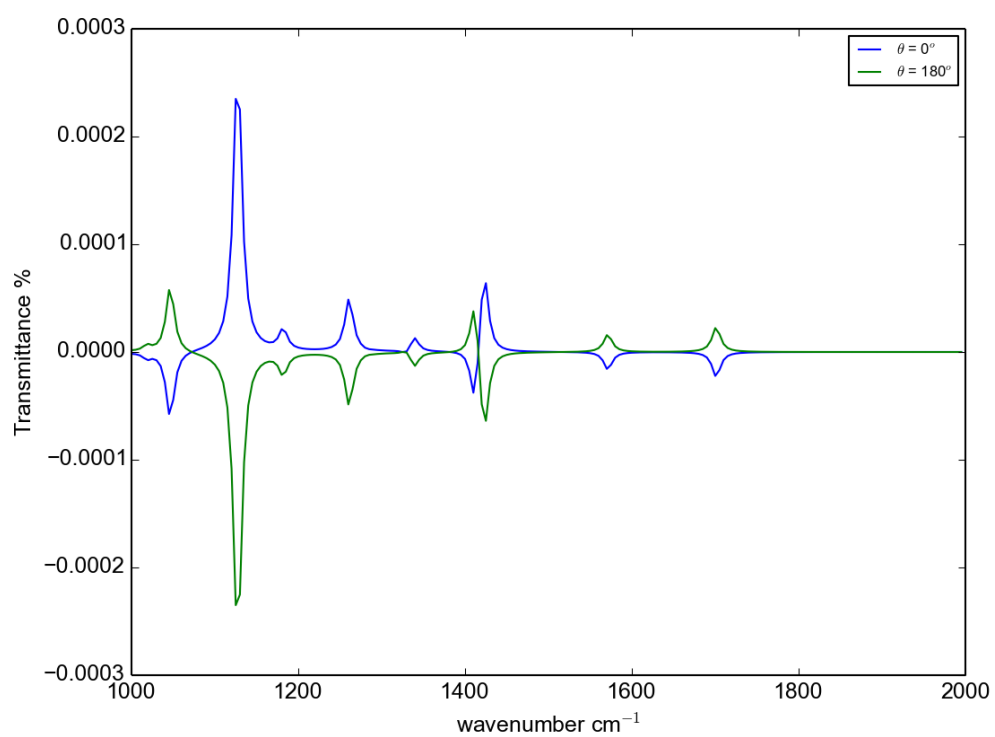


Figure 5.9: SFG  $zzz$  projection spectrum for alanine candidate with  $\theta$  of  $0^\circ$  is not identical to alanine candidate with  $\theta$  of  $180^\circ$ , but symmetric along wavelength



## Chapter 6

# Possibilities for Treating Experimental Data

### 6.1 Description

The experimental spectra obtained from IR, Raman or SFG techniques have an amplitude scaling factor when compared to the candidate spectra generated mathematically. This means that between candidates' theoretical spectra and the experimental one, there is an unknown scaling factor. Within one particular spectroscopy technique, this scaling factor is the same for any polarization. Take IR as an example. The scaling factor for the spectrum of  $x$  polarization is the same as the one for the spectrum of  $z$  polarization. It is necessary to introduce this scaling factor to our LP model.

### 6.2 Experiment

#### 6.2.1 Experiments with Scaling Factor Considering Each Amino Acid Candidates from $0^\circ$ to $80^\circ$ on $\theta$ in the Mixture

In Chapter 5, the LP model constructed by Experiments 2 to 7 in Table 5.1 for  $\theta$  ranged from  $0^\circ$  to  $80^\circ$  do well in retrieving the target composition for the mixed amino acids. Therefore, based on these experiments, we investigate the LP equations can be applied directly to the real experimental data for the same  $\theta$  range.

Therefore, the same experiment settings in Table 5.1 are used for the following experiments. The goal is the same, that is to figure out which spectral information helps to retrieve the target composition for the mixture of six amino acids' candidates. The only difference is that, in each run of the experiment set, an arbitrary scaling factor is generated for IR, Raman and SFG, respectively. Therefore, the target spectra are not only composed by the target composition of all candidates, but also need to multiple by the randomly generated scaling factors of each spectroscopy technique.

To start with, we limit the scaling factors to be smaller than 1.

After a few runs of the experiment set, it is observed that the returned compositions always contains one extra variable in every experiment. For Experiment 2, 4, 6 and 7, the returned composition contains the correct selected candidates. However, the percentage values of the selected candidates are different from the target composition. The ratio between the returned percentage and the target percentage are the same for all the selected candidates. Furthermore, when this ratio adds up the extra variable, it equals 1. Randomly select one experiment run as an example. Figure 6.2 displays the target composition, only the selected candidates are annotated with assigned percentage. Figure 6.2 displays the return composition of Experiment 2. The selected candidates in the return composition are correct. However, each percentage value is different from the one in the target composition. There is one extra value in Figure 6.2 with a value of 0.4.

Moreover, Equation 6.1 shows the ratio between the percentage of the selected candidates in the return composition and the target one is the same for all the amino acids. The value of this ratio is 0.6. When this ratio is added up with the extra variable (referred to as slack variable (SV) in LP) 0.4, the total is 1. As the scaling factors are pre-generated in the experiment set, the value is known, which is 0.6 for Raman spectra. In conclusion, the SV is returned by LP. Then the scaling factor (SF) equals to  $1 - SV$ . From the scaling factor, the ratio between the return composition and the target one is known. At the end, the target composition can be re-built from the ratio and the return composition. The re-constructed target composition matches to the original one.

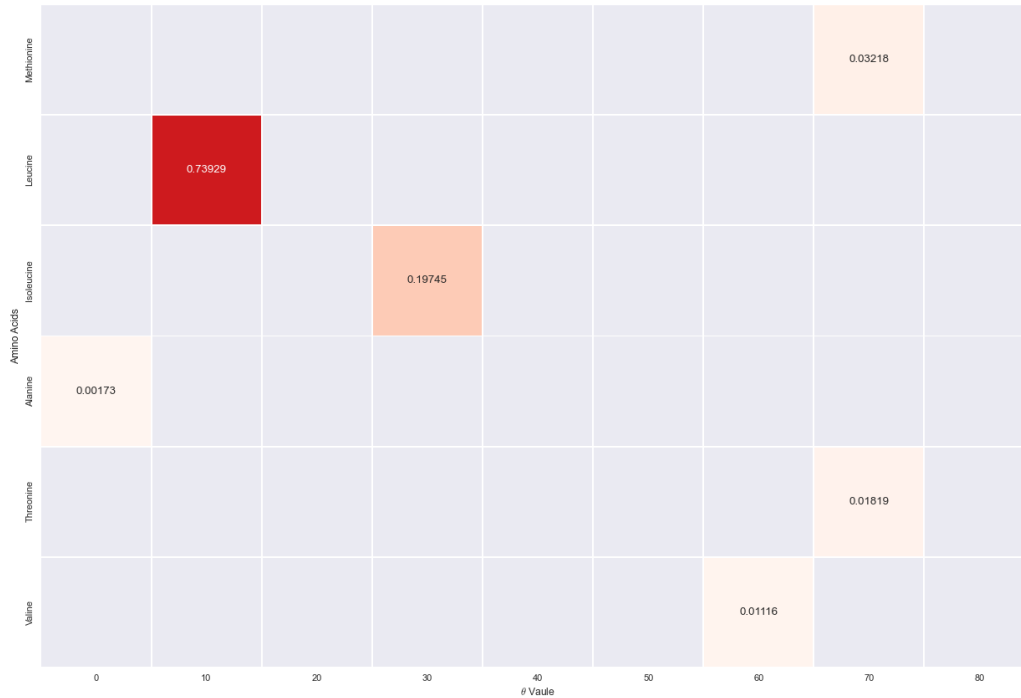


Figure 6.1: Target composition for one random run of the experiment set with scaling factor for mixed amino acids, with  $\theta$  expanded from  $0^\circ$  to  $80^\circ$ .

$$\frac{0.019308}{0.03218} = \frac{0.443574}{0.73929} = \frac{0.11847}{0.19745} = \frac{0.001038}{0.00173} = \frac{0.010914}{0.01819} = \frac{0.006696}{0.01116} = 0.6 \quad (6.1)$$

To verify if the above observation is a general case, the experiment set in Table 5.1 is run 100 times with randomly generated scaling factors in each run. Figure 6.3 indicates the experiment result. Experiment 2, 4, 6 and 7 hit the above observation with almost 100% frequency. This indicates that even with the scaling factor, Raman spectral information alone is sufficient to study the mixed molecules' coordination distribution at interfaces when each amino acid's candidates expanded from  $0^\circ$  to  $80^\circ$  on  $\theta$ . The target composition can be re-constructed correctly from the return slack variable and the return composition. Figure 6.3 also illustrates that Experiment 3 does not hit the above observation with high frequency. With the scaling factor as the addition, SFG spectral information is not sufficient to obtain the target composition. Experiment 5 indicates that even combining IR and SFG spectral information, the constructed LP model cannot help to reconstruct the target composition. This can cause by the different scaling factors of these two spectroscopy techniques.

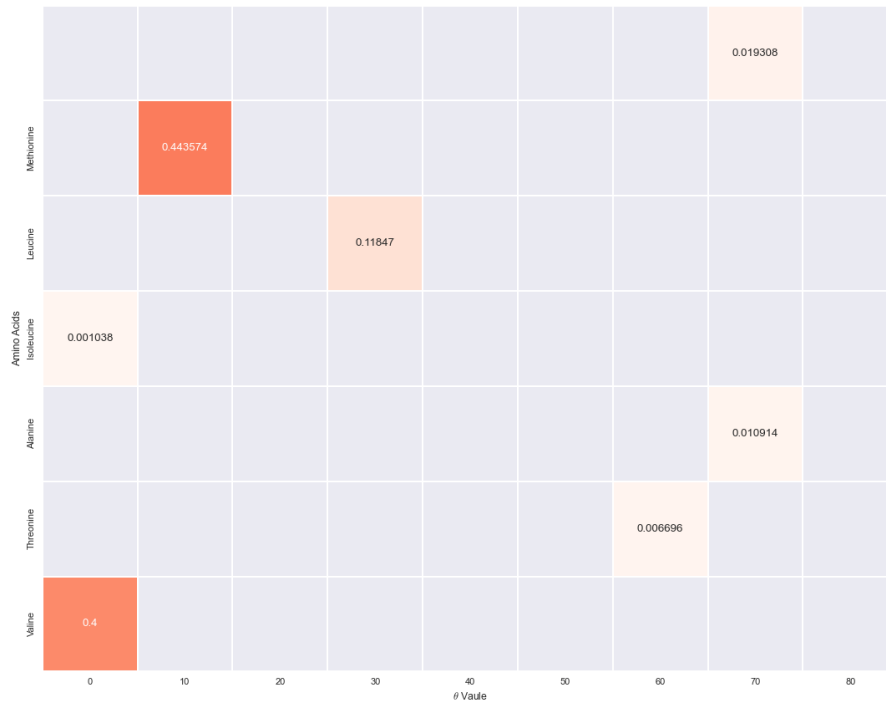


Figure 6.2: Return composition of Experiment 2 for one random run of the experiment set with scaling factor for mixed amino acids, with  $\theta$  expanded from  $0^\circ$  to  $80^\circ$ .

### 6.2.2 Experiments with Scaling Factor Considering Each Amino Acid Candidates from $0^\circ$ to $180^\circ$ on $\theta$

When each amino acid's candidates are expanded from  $0^\circ$  to  $180^\circ$  on  $\theta$ , the same experiment set is applied 100 times with randomly generated scaling factors in each run. The experiment result from the 100 run illustrates that all experiment in the set meets the above observation with zero frequency.

However, when further analyze the return compositions of Experiment 2 and 6, there are few other observations to be noted. To facilitate the explanation, one random run is picked as an explicit example. Figure 6.4 is the target composition. Figure 6.5 and Figure 6.6 are the return compositions of Experiment 2 and Experiment 6. The generated scaling factor for IR, Raman and SFG are 0.863411, 0.770505 and 0.239947.

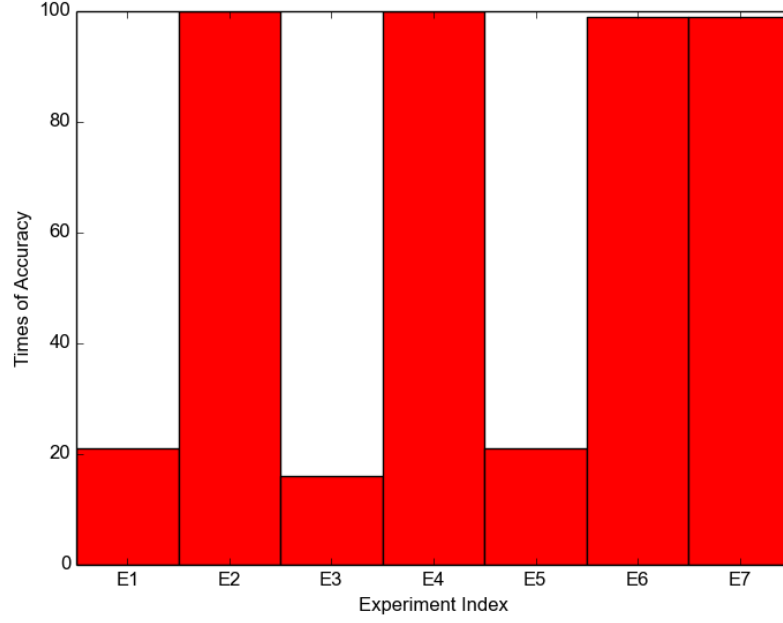


Figure 6.3: Experiment accuracy analysis for experiments using experimental spectra data that contains scaling factor that is smaller than 1 and candidates with  $\theta$  from  $0^\circ$  to  $80^\circ$

In Figure 6.5, in the return composition of Experiment 2, the slack variable equals  $1 - SF = 1 - 0.770505 = 0.229495$ . For each amino acid, the selected candidate in the return composition may not be the exact one as shown in the target composition. However, this selected candidate is always either the correct one, or the correct one's  $\theta$  complimentary. Moreover, the ratios between the percentage of each selected candidate in Figure 6.5 and Figure 6.4 are the same as shown in Equation 6.2. These ratios all equal to the scaling factor of Raman.

In Figure 6.5, for each amino acid, there are two selected candidates in the return composition. These two selected candidates are the correct one and its  $\theta$  complimentary. When the percentages of these two selected candidates are added, it equals to the percentage returned for the amino acid in Figure 6.4.  $0.27162 + 0.142619 = 0.414239$ . Between these two selected candidates, the correct one's percentage is always bigger than its  $\theta$  complement.  $0.27162 > 0.142619$ . In conclusion, Experiment 2 achieves in telling the slack variable, the scaling factor, and the ratio between the returned candidates and the target ones. However, in order to distinguish the exact candidate

of each amino acid, the extra information from Experiment 6 is required. Experiment 6 tells the correct candidate from its complement on  $\theta$ . Together with the return information from Experiment 2 and 6, the target composition can be obtained. These observations can be applied to every run of the experiment set.

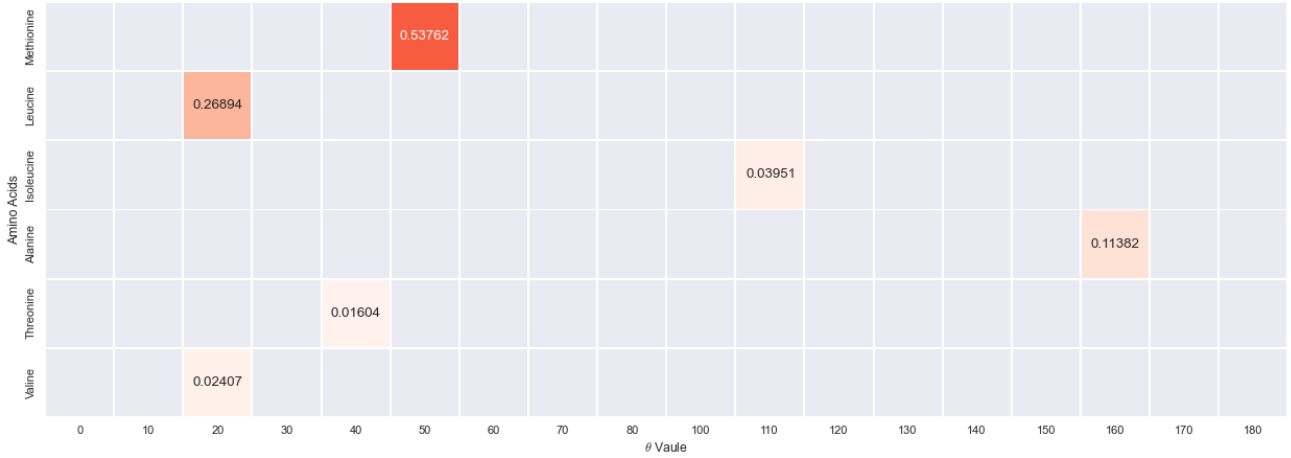


Figure 6.4: Target composition of one random run of experiments containing scaling factor and the mixed amino acids' candidates with  $\theta$  expended from  $0^\circ$  to  $180^\circ$

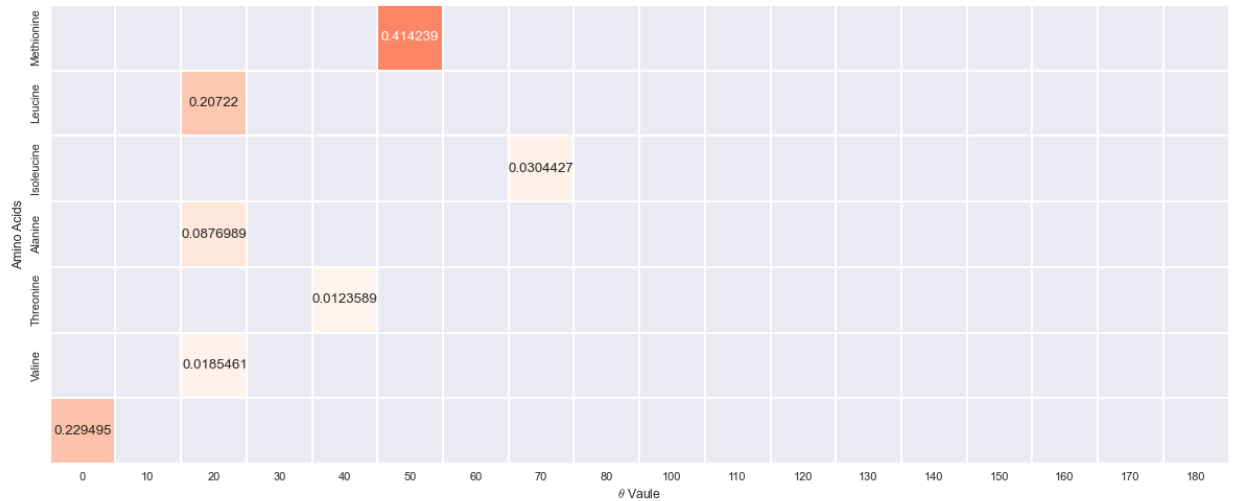


Figure 6.5: Return composition of Experiment 2 for one random run of experiments containing scaling factor and the mixed amino acids' candidates with  $\theta$  expended from  $0^\circ$  to  $180^\circ$

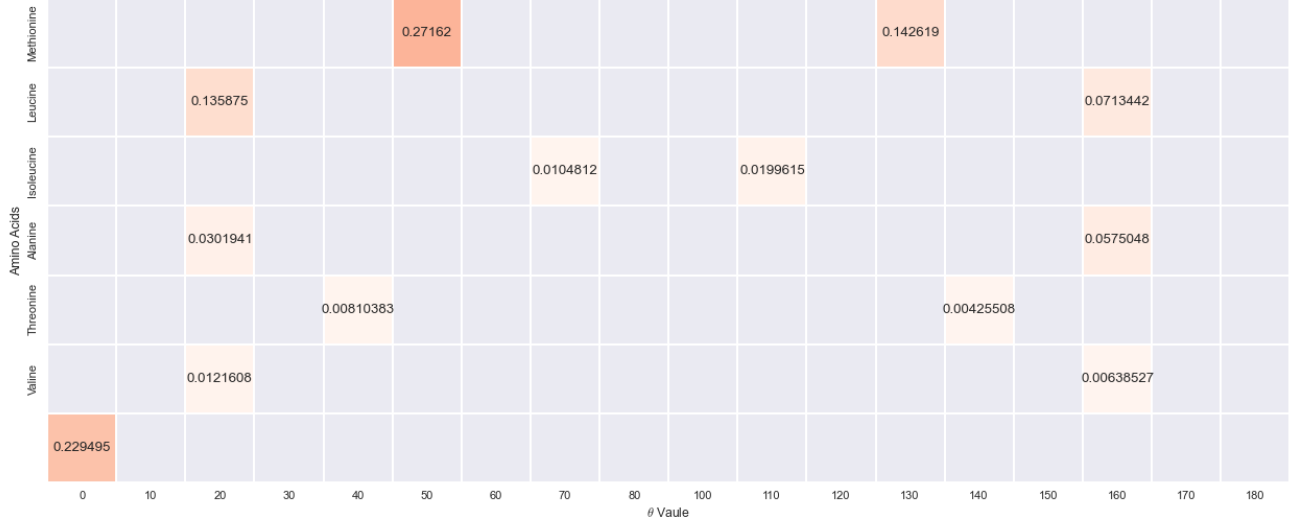


Figure 6.6: Return composition of Experiment 6 for one random run of experiments containing scaling factor and the mixed amino acids' candidates with  $\theta$  expanded from  $0^\circ$  to  $180^\circ$

$$\frac{0.414239}{0.53762} = \frac{0.20722}{0.26894} = \frac{0.0304427}{0.03951} = \frac{0.0876989}{0.11382} = \frac{0.0123589}{0.01604} = \frac{0.0185461}{0.02407} = 0.770505 \quad (6.2)$$

### 6.3 Conclusion

With Scaling factor introduced to different spectroscopy techniques, Raman spectral information alone is sufficient to obtain the target composition, when considering a mixture of amino acids with candidates expanded from  $0^\circ$  to  $80^\circ$  on  $\theta$ . The target composition can be re-constructed from the return SV and composition. The SF equals 1 minus SV.

When each amino acid's candidates are expanded from  $0^\circ$  to  $180^\circ$ , both return compositions from Experiment 2 and 6 are needed to obtain the target composition.

## Chapter 7

# Conclusion and Future Work

### 7.1 Conclusion

#### 7.1.1 Contributions

### 7.2 Future Work



# Appendix A

## Additional Information

This is a good place to put tables, lots of results, perhaps all the data compiled in the experiments. By avoiding putting all the results inside the chapters themselves, the whole thing may become much more readable and the various tables can be linked to appropriately.

The main purpose of an Appendix however should be to take care of the future readers and researchers. This implies listing all the housekeeping facts needed to continue the research. For example: where is the raw data stored? where is the software used? which version of which operating system or library or experimental equipment was used and where can it be accessed again?

Ask yourself: if you were given this thesis to read with the goal that you will be expanding the research presented here, what would you like to have as housekeeping information and what do you need? Be kind to the future graduate students and to your supervisor who will be the one stuck in the middle trying to find where all the stuff was left!

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.299586 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.664053 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.021196 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.00662804 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.000642609 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.00789 & 0 & 0 & 0 & 0 \end{bmatrix} \quad (\text{A.1})$$

$$\begin{bmatrix}
0 & 0 & 0 & 0 & 0 & 0 & 0.299586 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0.664053 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0.021196 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.00662804 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0.000642609 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0.00789 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0
\end{bmatrix} \quad (\text{A.2})$$

$$\begin{bmatrix}
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.299586 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0.664053 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.021196 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.00662804 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.000642609 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.00789 & 0 & 0 & 0 & 0 & 0 & 0
\end{bmatrix} \quad (\text{A.3})$$

$$\begin{bmatrix}
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.03218 & 0 \\
0 & 0.73929 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0.19745 & 0 & 0 & 0 & 0 & 0 \\
0.00173 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.01819 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0.01116 & 0 & 0
\end{bmatrix} \quad (\text{A.4})$$

$$\begin{bmatrix}
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.019308 & 0 \\
0 & 0.443574 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0.11847 & 0 & 0 & 0 & 0 & 0 \\
0.001038 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.010914 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0.006696 & 0 & 0 \\
0.4 & & & & & & & & 
\end{bmatrix} \quad (\text{A.5})$$

Return composition of Result 2 [0.0, 0.0, 0.0, 0.0, 0.0214195, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.00498874, 0

Return composition of Result 6 [0.0, 0.0, 0.0, 0.0, 0.0214195, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.00498874, 0

Target composition [0, 0, 0, 0, 0.08497, 0, 0, 0, 0, 0, 0, 0.01979, 0, 0, 0, 0, 0, 0, 0.41463, 0, 0, 0, 0, 0, 0,

$g_{10.639037067053}g_{20.252083642737}g_{30.168511414065}$



# Bibliography

- [1] Sophie Brasselet. Polarization-resolved nonlinear microscopy: application to structural molecular and biological imaging. *Adv. Opt. Photon.*, 3(3):205, Sep 2011.
- [2] Vasek Chvatal. *Linear Programming*. W. H. Freeman and Company, 1983.
- [3] Kuo Kai Hung. Extracting surface structural information from vibrational spectra with linear programming. Master’s thesis, University of Victoria, 2015.
- [4] Bernd Cartner Jiri Maousek. *Understanding and Using Linear Programming*. Springer, 2007.
- [5] S.T.Elbert. M.S.Gordon. J.H.Jensen. S.Koseki. N.Matsunaga. K.A.Nguyen. S.J.Su. T.L.Windus. M.Dupuis. J.A.Montgomery M.W.Schmidt. K.K.Baldrige. J.A.Boatz. *General Atomic and Molecular Electronic Structure System*. Department of Chemistry Iowa State University, July 2016.
- [6] Arnold W. Pratt, J. Nicolet Toal, and George W. Rushizky. Computer assisted analysis of oligonucleotides. *Annals of the New York Academy of Sciences*, 128(3):900–913, 1966.
- [7] William C. Whiten. Marvin B. Shapiro. Arnold W. Pratt. Linear programming applied to ultraviolet absorption spectroscopy. *Communications of the ACM*, 6:66–67, 1963.
- [8] Sandra. Hung Kuo-Kai. Stege Ulrike. Roy and Hore Dennis K. Rotations, projections, direction cosines, and vibrational spectra. *Applied Spectroscopy Reviews*, 49:233–248, May 1999.

- [9] X. Zhuang, P. B. Miranda, D. Kim, and Y. R. Shen. Mapping molecular orientation and conformation at interfaces by surface nonlinear optics. *Phys. Rev. B*, 59:12632–12640, May 1999.