

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^AT_EX 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 | 23 |
| Table 3.3 | Experiment 4 and 5 Setting | 26 |
| Table 3.4 | Constraint Study Based on Experiment 4 | 27 |
| Table 3.5 | Constraint Study Based on Experiment 5 | 29 |
| Table 4.1 | Experiment 1 to Experiment 4 Setting for Methionine Candidates | 32 |
| Table 4.2 | Experiment 5 to Experiment 9 Setting for Methionine Candidates | 34 |
| Table 4.3 | Experiment 5 to Experiment 9 Setting for Methionine Candidates | 36 |
| Table 4.4 | Experiments to Explain the Limitation of LP Model for Methio- nine Molecule | 37 |
| Table 5.1 | Detailed Experiment Group Setting | 45 |

List of Figures

| | | |
|------------|---|----|
| Figure 2.1 | The Euler angles represented as the spherical polar angles θ , ϕ and ψ , 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 [?]. | 10 |
| Figure 2.2 | IR x -polarized spectra of methionine's four candidates and target. The candidates are with θ of 0° , 20° , 40° and 60° . 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 θ of 0° , 20° , 40° and 60° . 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 θ of 0° , 20° , 40° and 60° . 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 θ of 0° , 20° , 40° and 60° . 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 Result Plotting for 4 Candidates on IR Cosine Projection | 22 |
| Figure 3.3 | Toy Model Result Plotting for 10 Candidates on IR Cosine Projection | 24 |
| Figure 3.4 | Toy Model Candidates IR Sine Projection | 25 |
| Figure 3.5 | Toy Model Constraint Study 1 | 28 |

| | |
|---|----|
| Figure 3.6 Toy Model Constraint Study 2 | 30 |
| Figure 4.1 Compare target spectra with spectra generated by composition returned by LP model with only IR spectra of x and z projection | 33 |
| Figure 4.2 IR spectra plotted by using target composition and return com- position of Experiment 17 | 38 |
| Figure 4.3 Raman spectra plotted by using target composition and return composition of Experiment 17 | 39 |
| Figure 4.4 SFG spectra plotted by using target composition and return com- position of Experiment 17 | 40 |
| Figure 4.5 IR spectra plotted by using target composition and return com- position of Experiment 18 | 41 |
| Figure 4.6 Raman spectra plotted by using target composition and return composition of Experiment 18 | 42 |
| Figure 4.7 SFG spectra plotted by using target composition and return com- position of Experiment 18 | 43 |
| Figure 5.1 Accuracy analysis for experiments considering a mixture of amino acids with candidates from 0° to 80° on θ for each amino acid . | 47 |
| Figure 5.2 IR Spectra Plotted by Result Composition and Target Compo- sition. | 49 |
| Figure 5.3 Accuracy analysis for experiments considering a mixture of amino acids with candidates from 0° to 180° on θ for each amino acid | 50 |
| Figure 5.4 Target composition for one random run of six mixed amino acids with θ expanded from 0° to 180° on θ | 52 |
| Figure 5.5 return composition of experiment 2 for one random run of six mixed amino acids with θ expanded from 0° to 180° on θ | 52 |
| Figure 5.6 return composition of experiment 6 for one random run of six mixed amino acids with θ expanded from 0° to 180° on θ | 53 |
| Figure 5.7 IR z projection spectrum for alanine candidate with θ of 0° is identical to alanine candidate with θ of 180° | 54 |
| Figure 5.8 Raman zz projection spectrum for alanine candidate with θ of 0° is identical to alanine candidate with θ of 180° | 55 |
| Figure 5.9 SFG zzz projection spectrum for alanine candidate with θ of 0° is not identical to alanine candidate with θ of 180° , but symmetric along wavelength | 56 |

| | | |
|------------|--|----|
| Figure 6.1 | Target composition for one random run of experiment set with scaling factor for mixed amino acids with θ expended from 0° to 80° | 58 |
| Figure 6.2 | Return composition of Experiment 2 for one random run of experiment set with scaling factor for mixed amino acids with θ expended from 0° to 80° | 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 θ from 0° to 80° | 60 |

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. [?]

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 [?], [?]. 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 μ (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 μ 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 α (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 [?]:

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

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 [?].

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 [?]. However, both approaches take a lot of time and computational resources. In Hung’s study [?], 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 [?].

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 [?] describes how to perform electronic structure calculations using GAMESS [?] 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 θ , ϕ and ψ 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 θ is the angle between z and c . *Azimuthal* angle ϕ is the rotation about z . *Twist* angle ψ is a twist about c [?]. 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)

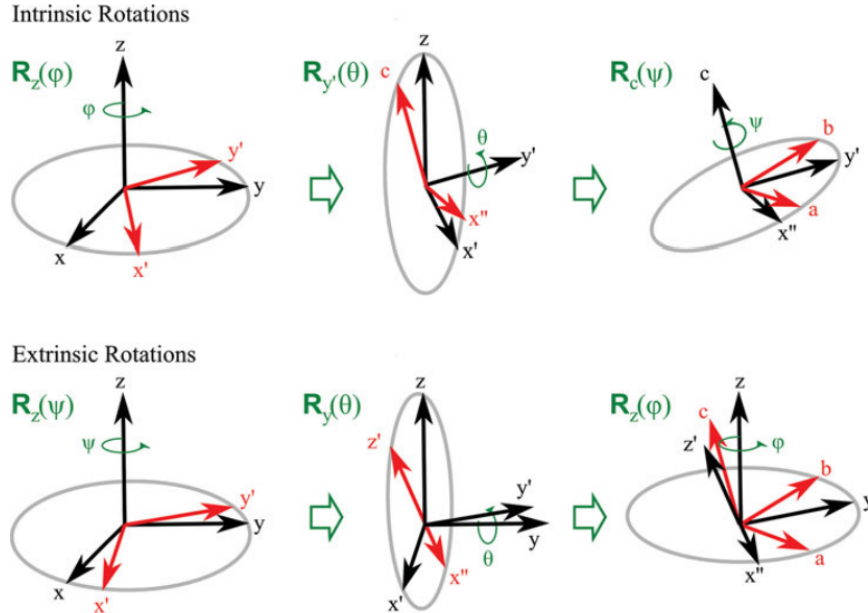


Figure 2.1: The Euler angles represented as the spherical polar angles θ , ϕ and ψ , and the illustration of the three successive rotations that transform the lab x , y , z coordinate system into the molecular a , b , c frame intrinsically and extrinsically [?].

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 ψ , 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 ψ . On the other hand, for angle ϕ , 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 θ 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 . ω_{IR} is the frequency of the probe radiation. μ is the dipole moment. m_q is the reduced mass. ω_q is resonance frequency. Γ_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 ω_{IR} , μ , m_q , Q are obtained from the molecular property files. Furthermore, because ϕ and ψ 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 , ω_q , Γ_q , and Q_q are the same as defined above for IR spectra. All the values of ω_{IR} , μ , m_q , Q are obtained from the molecular property files. Similar to IR spectroscopy, because of the integration of ϕ and ψ 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, α^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 ω_{vis} and a x -polarized infrared beam incoming with frequency ω_{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 [?]. The SFG response is the imaginary component of the second-order susceptibility. Same as IR and Raman spectroscopy, all the values of ω_{IR} , μ , m_q , Q are obtained from the molecular property files. Because of the integration of ϕ and ψ 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 θ value. Taking Methionine as an example, Figure 2.2 displays x -polarized IR spectra of the following candidates: Methionine with θ of 0° , 20° , 40° and 60° . Their spectra are prefixed with *candidate_* in the labels. *ir_x_* indicates the spectroscopy technique, “number” indicates the θ 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.

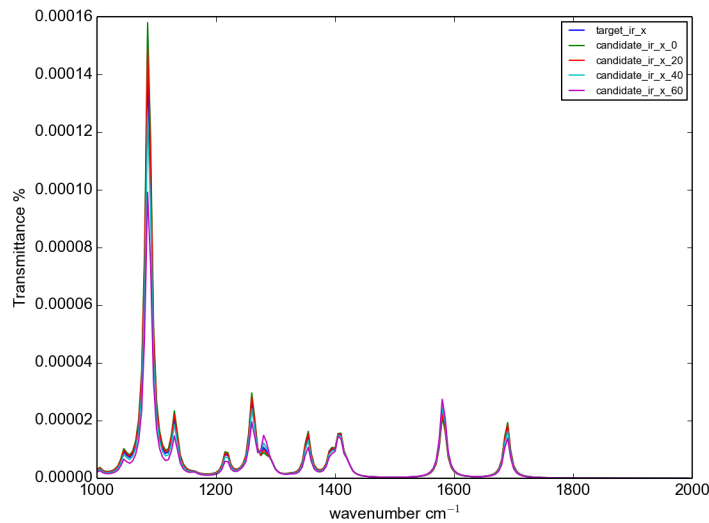


Figure 2.2: IR x -polarized spectra of methionine’s four candidates and target. The candidates are with θ of 0° , 20° , 40° and 60° . The target is produced by combining 10% of *candidate_ir_x_0*, 50% *candidate_ir_x_20* and 40% *candidate_ir_x_40*.

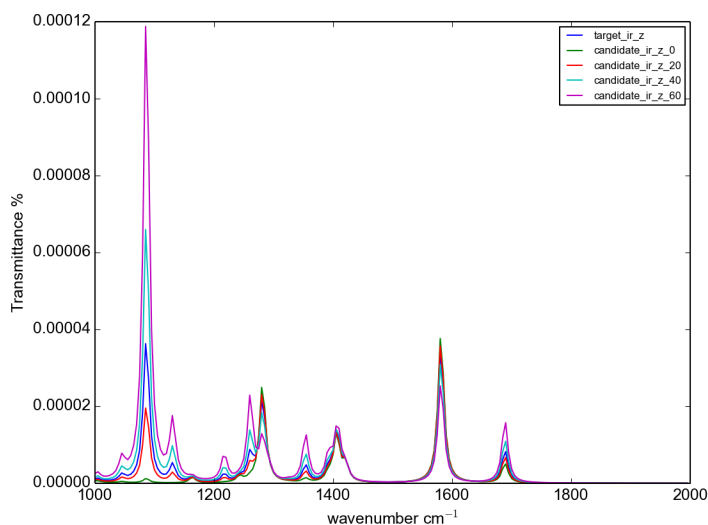


Figure 2.3: IR z -polarized spectra of methionine's four candidates and target. The candidates are with θ of 0° , 20° , 40° and 60° . The target is produced by combining 10% of *candidate_ir_x_0*, 50% *candidate_ir_x_20* and 40% *candidate_ir_x_40*.

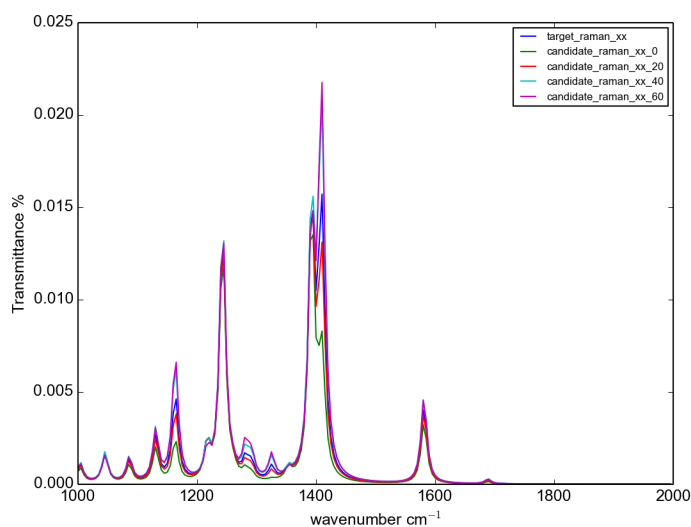


Figure 2.4: Raman xx -polarized spectra of methionine's four candidates and target. The candidates are with θ of 0° , 20° , 40° and 60° . The target is produced by combining 10% of *candidate_ir_x_0*, 50% *candidate_ir_x_20* and 40% *candidate_ir_x_40*.

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

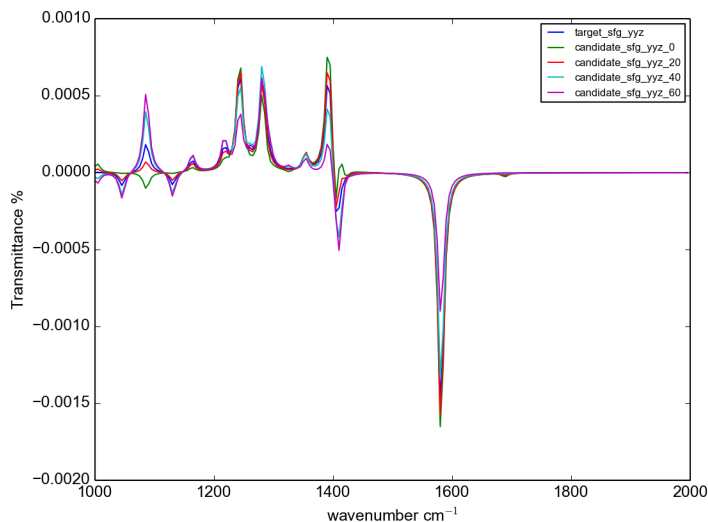


Figure 2.5: SFG *yyz*-polarized spectra of methionine's four candidates and target. The candidates are with θ of 0° , 20° , 40° and 60° . The target is produced by combining 10% of *candidate_ir_x_0*, 50% *candidate_ir_x_20* and 40% *candidate_ir_x_40*.

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 ϕ and ψ Euler angles are integrated, only the difference on angle θ 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, θ_q is the comparing angle, γ is the width, and ω_q is the frequency. (TODO: Double check the correct meaning of each symbol) Ten candidates are produced with 10 different θ values as follows: 0° , 10° , 20° , 30° , 40° , 50° , 60° , 70° , 80° , 90° . 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.

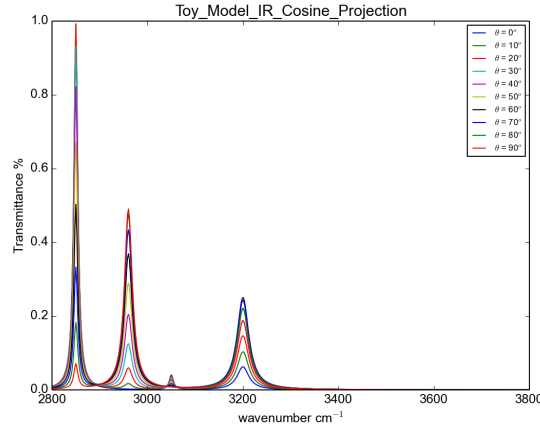


Figure 3.1: *cosine-* polarize IR spectra of toy model candidates

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 [?] and other UV spectroscopy studies [?] 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:

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

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

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 θ value, the candidate's spectral information is obtained. For toy model, only the value of θ 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 θ 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 [?]. 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 θ of 0° , 10° , 20° , and 30° . In Experiment 2, the four candidates are with θ values of 0° , 5° , 10° , and 15° . Instead of having a 10 degree variance in θ , 5 degree difference is applied on θ 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 θ values ranging from 0° to 90° .

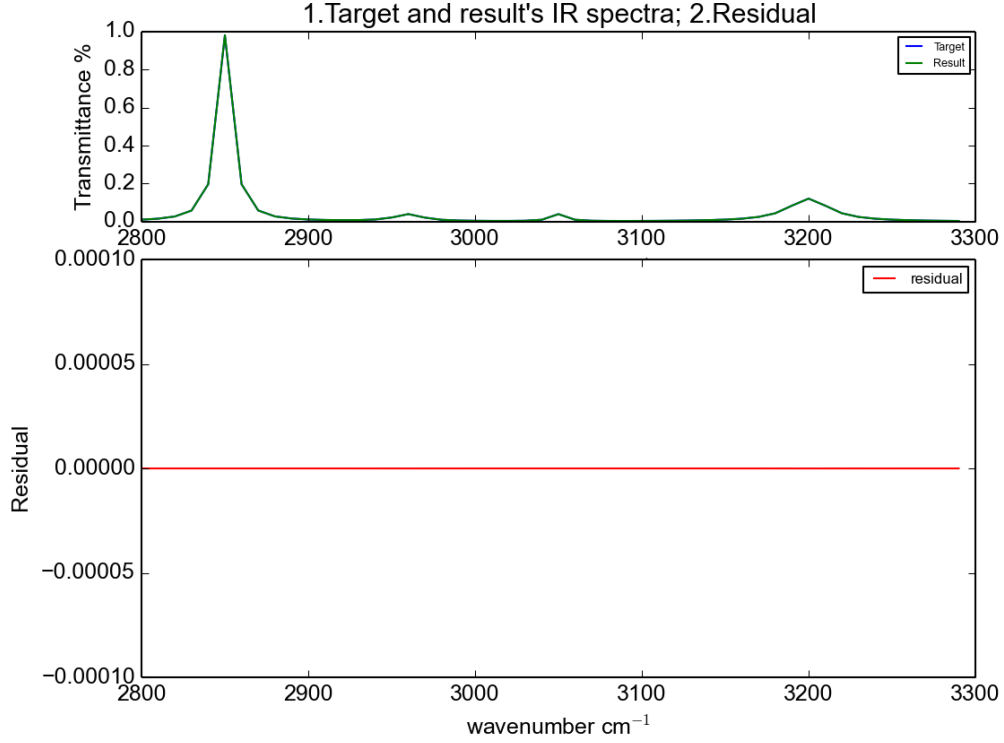


Figure 3.2: Toy model Experiment 2 resulting *cosine*-polarized IR spectrum plotted with the target spectrum; and the residual plot between the spectra.

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 θ 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

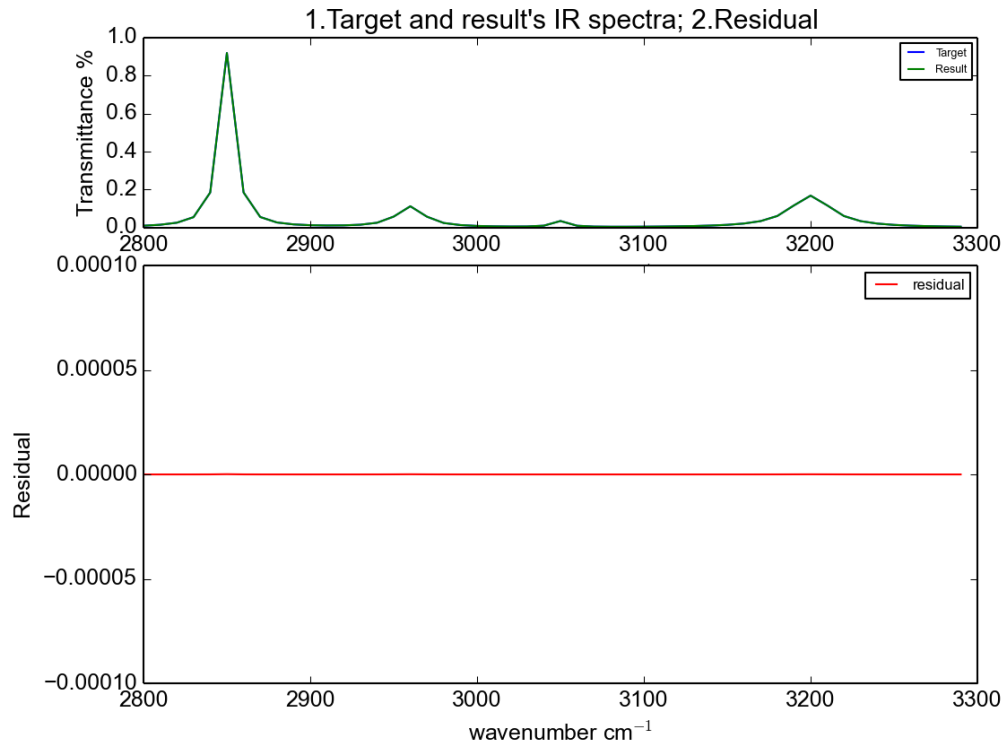


Figure 3.3: Toy model Experiment 3 resulting *cosine*-polarized IR spectrum plotted with target spectrum; and the residual plot between the two spectra

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. Same with Experiment 5, it is based on Experiment 3, with sine-polarization IR spectral information added. For both Experiment 4 and 5, the return composition now matches to the target one. This further proves that as long as we have sufficing information to build the constraints, the LP solver will return a composition matches to the target one.

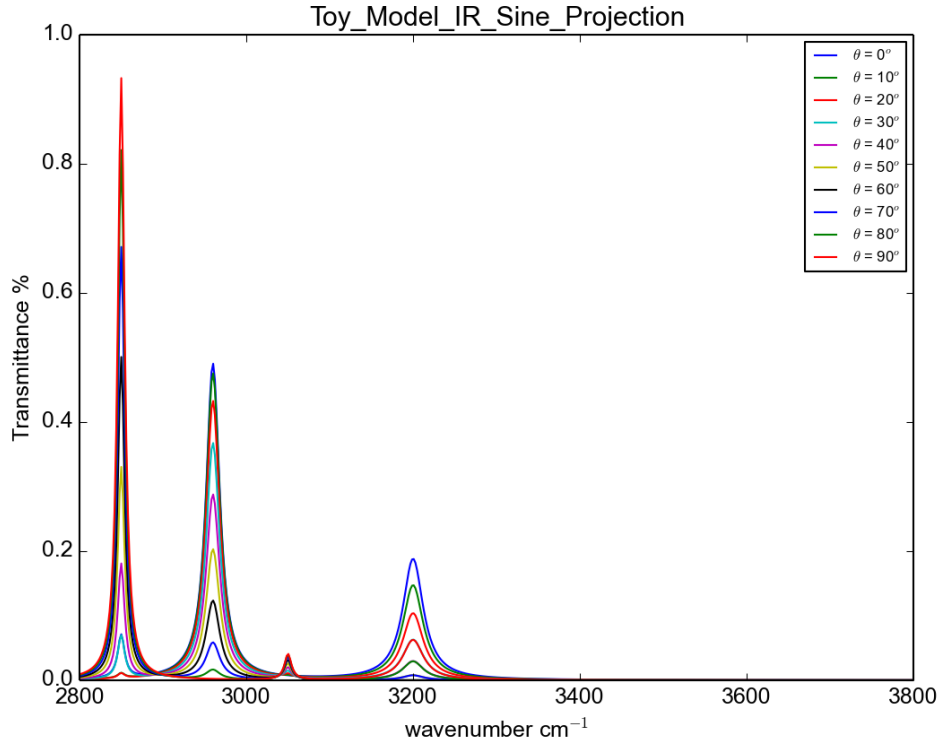


Figure 3.4: *sine*-polaried IR spectra of toy model candidates with θ value expanded from 0° to 90°

3.5 Constraint Study Based on Experiment 4

From Experiment 1 to 5, 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. Moreover, the information is coming from the valuable data points selected along the spectra. This

| | | |
|-----------------------|--|---|
| 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 Index | Number of Data Points | Points Selection | Result |
|------------------|-----------------------|--|------------------------------|
| 6 | 10 | [2800, 3300, 50] | [0, 0.796962, 0.103038, 0.1] |
| 7 | 20 | [2800, 3300, 25] | [0, 0.796962, 0.103038, 0.1] |
| 8 | 25 | [2800, 3300, 20] | [0, 0.796962, 0.103038, 0.1] |
| 9 | 32 | [2800, 3300, 15] | [0, 0.796962, 0.103038, 0.1] |
| 10 | 50 | [2800, 3300, 10] | [0, 0.796962, 0.103038, 0.1] |
| 11 | 100 | [2800, 3300, 5] | [0, 0.796962, 0.103038, 0.1] |
| 12 | 100 + 1 | [2800, 3300, 5], [2800, 3300, 500] | [0, 0.796962, 0.103038, 0.1] |
| 13 | 100 + 5 | [2800, 3300, 20], [2800, 3300, 100] | [0, 0.796962, 0.103038, 0.1] |
| 14 | 100 + 10 | [2800, 3300, 20], [2800, 3300, 50] | [0, 0.796962, 0.103038, 0.1] |
| 15 | 100 + 50 | [2800, 3300, 20], [2800, 3300, 10] | [0.1, 0.5, 0.4, 0] |
| 16 | 100 + 100 | [2800, 3300, 20], [2800, 3300, 5] | [0.1, 0.5, 0.4, 0] |

Table 3.4: Constraint Study Based on Experiment 4

leads us to do a further study on the constraints in order to see how many data 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. [2800, 3300, 50] means along wavenumber from 2500 to 3300, every 25 wavenumber. For example, Experiment 6 contains 10 data points from *cosine*-polarized IR spectrum. Every 50 wavenumber, one data point is selected. Similarly, for Experiment 7, 8, 9, 10, 11, every 25, 20, 15, 10 and 5 wavenumber, one data point is select. From Experiment 12 to 16, data points are selected from both cosine-polarized and sine-polarized IR spectrum.

(TODO: rethink: What can we exactly get from the following two tables? Should we include this study?)

One interesting result from Table 3.4 is that: from Experiment 1 to 9, the result composition is the same. To the contrary, from Experiment 10, the return composition gets changed to the target one. Furthermore, if we plot the return composition of [0,

$[0.796962, 0.103038, 0.1]$ and the target one $[0.1, 0.5, 0.4, 0]$ in Picture 3.5. In this picture, we can see that the spectra generated by these two composition are identical.

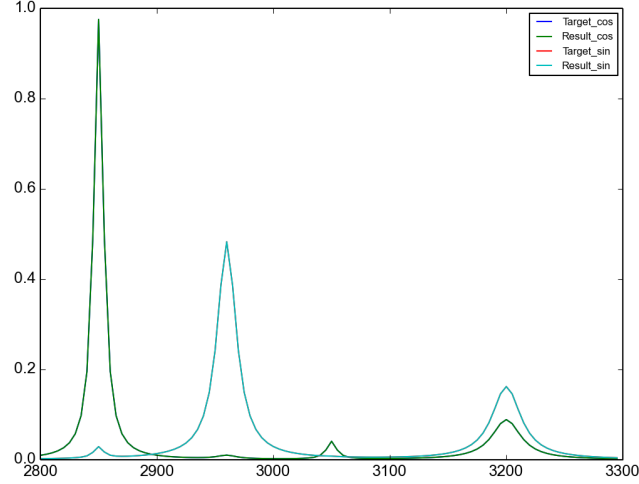


Figure 3.5: Toy Model Constraint Study 1

3.6 Constraint Study Based on Experiment 5

When the same constraint study is applied to the data based on Experiment 5 in Table 3.5, the observation is the same as the experiments in Table 3.4. This further proves that: We can obtain different solutions by have different constraints. 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 plotted together, they are almost identical as well, as shown in Figure 3.6.

3.7 Discussion and Conclusion

With all the experiments conducted with the toy model, we have learnt that the reason, that our LP model does not return a composition that matches the target one, is that the model does not have sufficient information to build the constraints. However, with the limited information, the optimal solution returned by our LP model does build the perfect target spectrum. This means that the solution for the composition that achieves minimum residual of the objective function is not unique. However, in real experiment, because numerical restriction, an unique optimal solution is obtained

| Experiment Index | Points | Point Selection | Result |
|------------------|-----------|------------------------------------|---|
| 17 | 10 | [2800, 3300, 50] | [0.156758, 0, 0, 0.825977, 0, 0, 0, 0, 0, 0.017265] |
| 18 | 25 | [2800, 3300, 20] | [0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0, 0] |
| 19 | 50 | [2800, 3300, 10] | [0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0, 0] |
| 20 | 100 | [2800, 3300, 5] | [0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0, 0] |
| 21 | 500 | [2800, 3300, 5] | [0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0, 0] |
| 22 | 100 + 1 | [2800, 3300, 5], [2800, 3300, 500] | [0, 0, 0.730541, 0, 0.212061, 0, 0, 0.0573978, 0, 0, 0] |
| 23 | 100 + 10 | [2800, 3300, 5], [2800, 3300, 50] | [0.361587, 0, 0.312061, 0.326352, 0, 0, 0, 0, 0, 0] |
| 24 | 100 + 20 | [2800, 3300, 5], [2800, 3300, 25] | [0.174023, 0, 0, 0.791447, 0, 0, 0.0345301, 0, 0, 0] |
| 25 | 100 + 25 | [2800, 3300, 20], [2800, 3300, 20] | [0.174023, 0, 0, 0.791447, 0, 0, 0.0345301, 0, 0, 0] |
| 26 | 100 + 50 | [2800, 3300, 5], [2800, 3300, 10] | [0, 0, 0.753209, 0, 0.146791, 0, 0.1, 0, 0, 0] |
| 27 | 100 + 84 | [2800, 3300, 5], [2800, 3300, 6] | [0.174023, 0, 0, 0.791447, 0, 0, 0.0345301, 0, 0, 0] |
| 28 | 100 + 100 | [2800, 3300, 5], [2800, 3300, 5] | [0.1, 0, 0.5, 0, 0.4, 0, 0, 0, 0, 0] |

Table 3.5: Constraint Study Based on Experiment 5

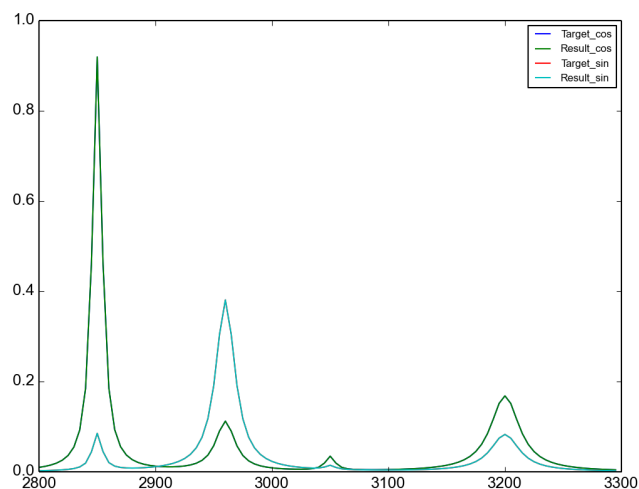


Figure 3.6: Toy Model Constraint Study 2

from the LP model.

Above analysis simulates the following question: how can we know there is enough information to achieve the target composition? In the next step, we will experiment with real molecules. The goal is to investigate with all the spectral information that we can obtain for real molecules, can our LP model return the target composition for the target spectrum? If yes, can we apply the LP model systematically? Furthermore, to maximally explore the capacity of our LP model, and study its limitation. Finally, come up with some general instructions for applying our LP model. These are the main focus for the following chapters.

Chapter 4

Realistic Molecular Model

4.1 Description

After experimenting with our toy problem, lacking sufficient information for the LP model is the key cause for the failure of obtaining the correct composition of the target spectra. First of all, in the toy model, there are only four vibrational modes, and the wavenumber range is limited. The number of data points selected is limited. Therefore, the constraints for the LP model is not sufficient. Secondly, the similarity among the candidates is high, as all the candidates are coming from one same molecule with the only difference in value θ . Third, the data points are only extracted from IR spectra.

In this chapter, real molecules are introduced. They contain more abundant spectroscopy information. In addition to IR, both Raman and SFG are introduced for real molecules, which makes the study one step closer to the overall goal and scope. Similar experiments applied to the toy model are now applied to a real molecule. The real molecule focused on this chapter is an amino acid - Methionine.

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 θ 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

| | | | | |
|-----------------------|----------------------------|----------------------------|----------------------------|-------------------------|
| Experiment index | 1 | 2 | 3 | 4 |
| Number of 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] |
| Number of 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

spectra can be obtained from xx , xy , xz and zz polarizations. Three unique SFG spectra can be obtained from yyz , zyz and zzz polarizations.

All the three spectroscopy techniques are applied in the experiments regarding real molecule. The goal is to see if those spectral information is enough to construct a LP model that returns the correct coordination distribution information about an amino acid at interfaces. If yes, we need to figure out what spectral information is needed to construct 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 θ of the following degree: 0° , 20° , 40° and 60° . 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, for 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 result spectra generated by using the return composition obtained from the first three experiments. The resulting spectra

is almost identical to the target ones. It again proves that with the information only coming from IR spectra is not sufficient to get the target composition. However, the return composition could perfectly re-product the target spectra. This indicates that we would need further information for the constraint of LP model, in order to further refine the result. The more constraints are introduced, the more accurate the return composition will be.

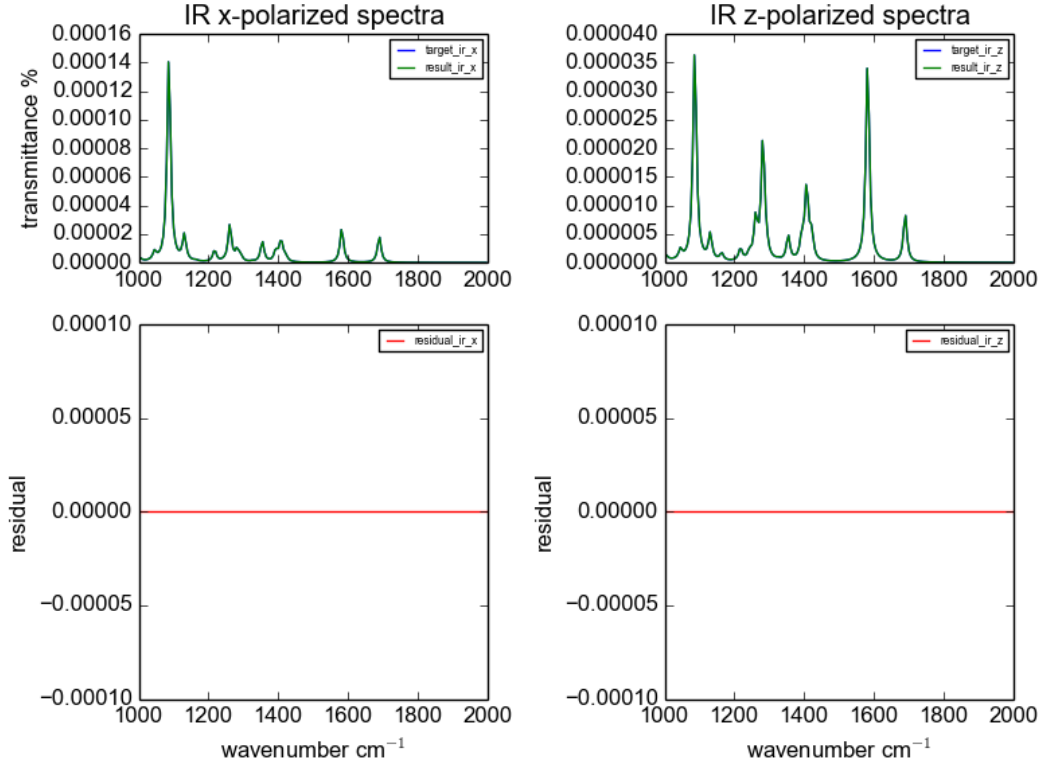


Figure 4.1: Compare target spectra with spectra generated by composition returned by LP model with only IR spectra of x and z projection

For the setting in Experiment 1 to 4, the LP model that constructed from combining IR and Raman spectral information is sufficient to obtain the target composition. When the difference in θ degree for candidates is smaller than 20° , which is 10° , we need to check if Raman and IR together still sufficient enough to derive the target composition. 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

| | | |
|----------------------|--|----------------------------|
| Number of Candidates | 4 | |
| Candidates | [0, 10, 20, 30] | |
| Target Composition | [0.1, 0.5, 0.4, 0] | |
| Experiment index | Number of 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(ramanzx) | [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(ramanzx)+ 200(ramanzz) | [0.1, 0.5, 0.4, 0] |

Table 4.2: Experiment 5 to Experiment 9 Setting for Methionine Candidates

information is not sufficient enough to derive the target composition for the current candidate setting. Experiment 6 indicates that combining IR and Raman spectral information helps to derive the target composition. What’s more, Experiment 7 to 9, illustrates 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 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 have been designed as shown in Table 4.3. There are 5 candidates included in the experiments. Each candidate has θ with the following degree: 0° , 10° , 20° , 30° and 40° . The target composition is more complex than previous experiments, each candidate takes 20% in the mix.

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.

As the result of Experiment 10 to 16 indicates 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. The LP model is showing its limitation in these experiments. In order to confirm if the reason causing the LP model to return a different composition is because of insufficient information. Further experiments are conducted as shown in Table 4.4.

| | | |
|----------------------|---|--|
| 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(ramanzx) + 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(ramanzx) + 200(ramanzz) | [0.247792, 0, 0.502139, 0, 0.250069] |
| 14 | 200(ramanxx) + 200(ramanxy) + 200(ramanzx) + 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(ramanzx) + 200(ramanzz) + 200(sfgyyz) + 200(sfgzyz) + 200(sfgzzz) | [0.321014, 0, 0.31018, 0.163041, 0.205764] |

Table 4.3: Experiment 5 to Experiment 9 Setting for Methionine Candidates

| | | |
|----------------------|--|--|
| Number of 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 index | Number 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: Experiments to Explain the Limitation of LP Model for Methionine Molecule

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' θ values are expanded from 0° to 80° . In total, there are 9 candidates. Because the SFG spectra for θ of 90° is a straight line, it is excluded from all the experiments. 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, and plotted together with the spectra created by 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 return composition and target one for Experiment 17. Every spectrum is almost identical to each other 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 information to construct the LP model is not sufficient to converge to the one exactly matches to the target composition. This conclusion exactly fits the result obtained from the experiments have done with the toy molecule.

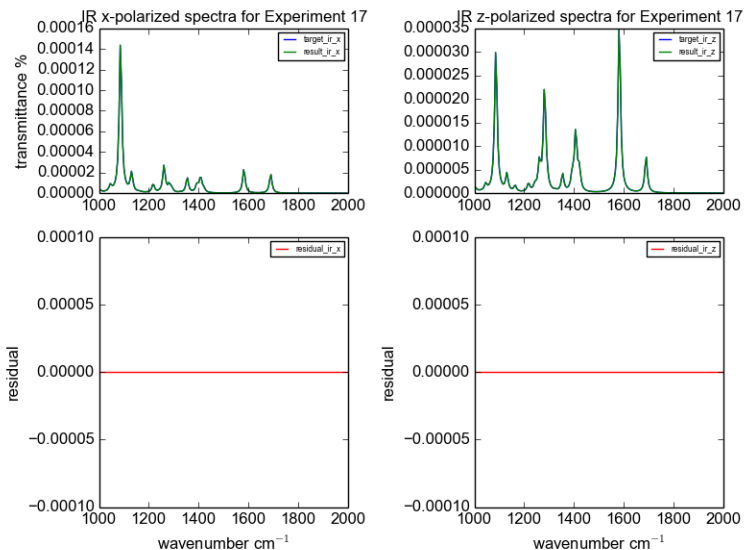


Figure 4.2: IR spectra plotted by using target composition and return composition of Experiment 17

4.4 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° to 80° on θ .

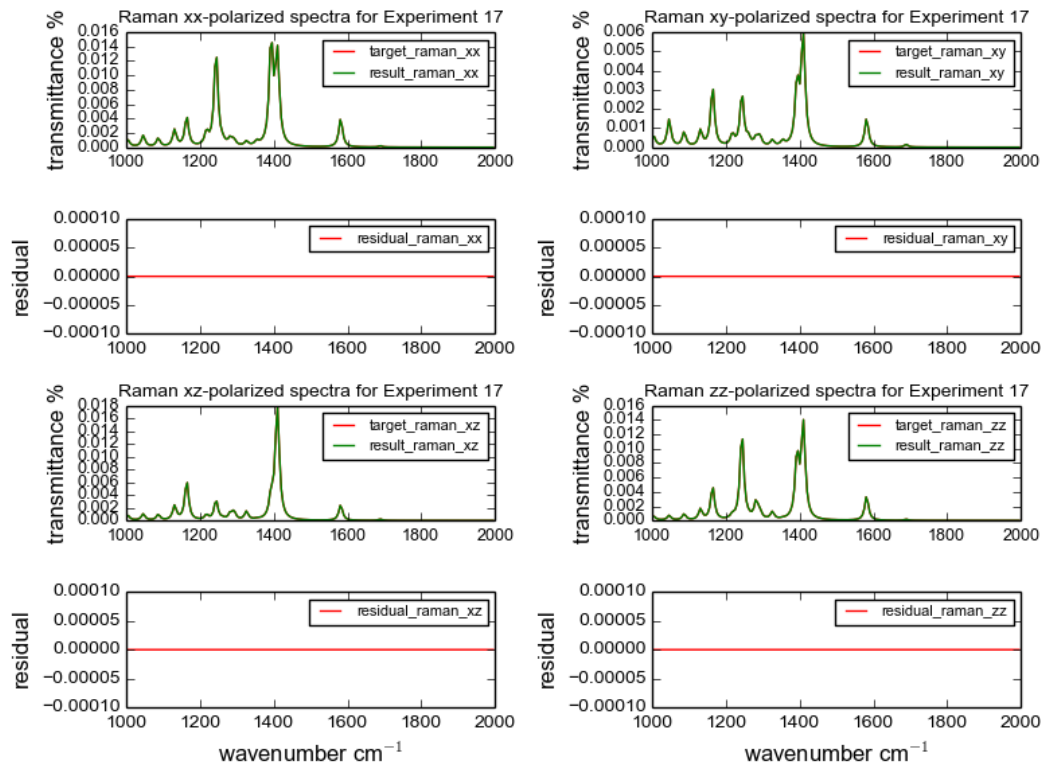


Figure 4.3: Raman spectra plotted by using target composition and return composition of Experiment 17

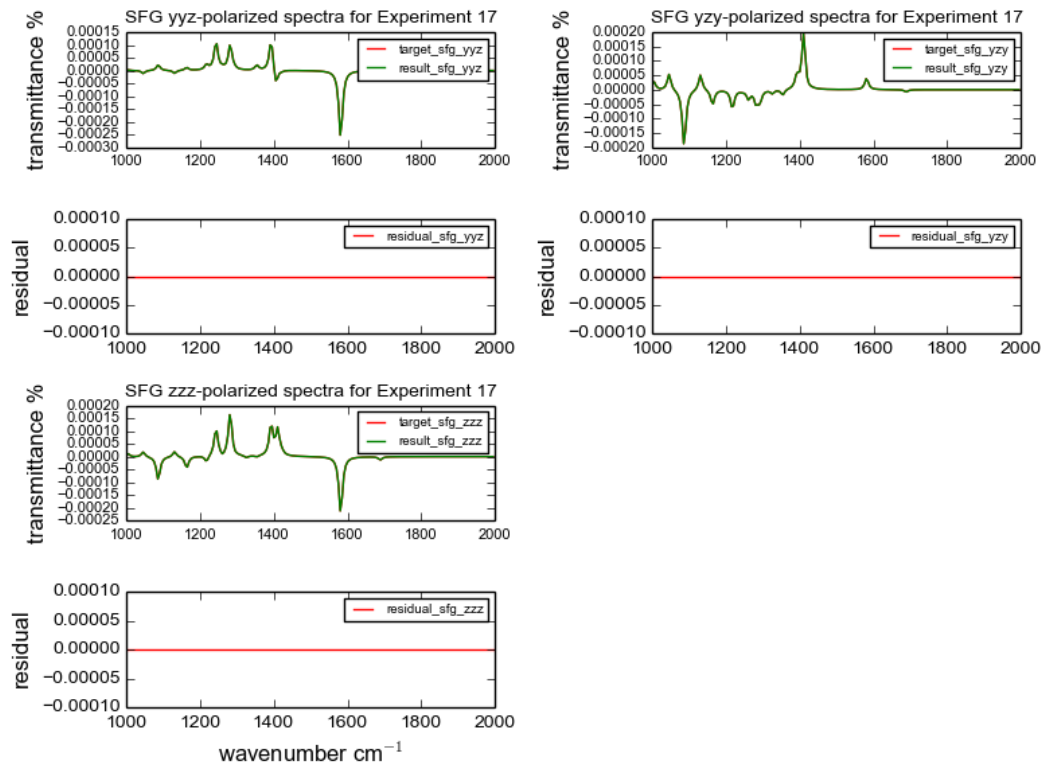


Figure 4.4: SFG spectra plotted by using target composition and return composition of Experiment 17

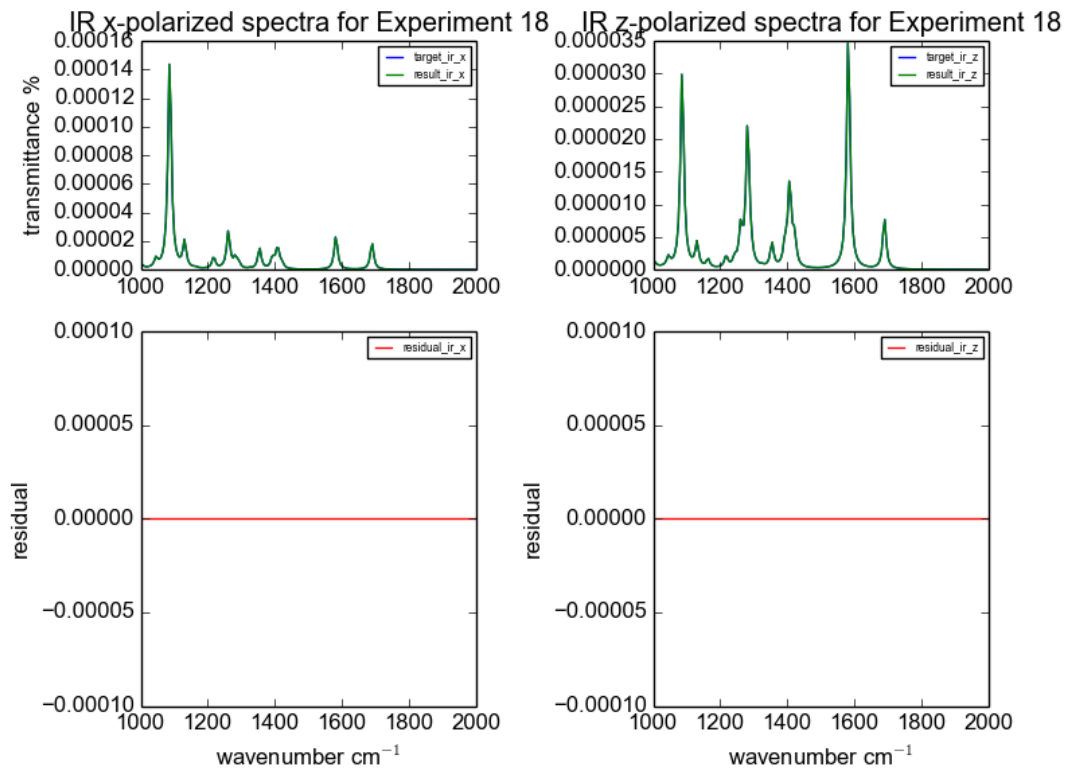


Figure 4.5: IR spectra plotted by using target composition and return composition of Experiment 18

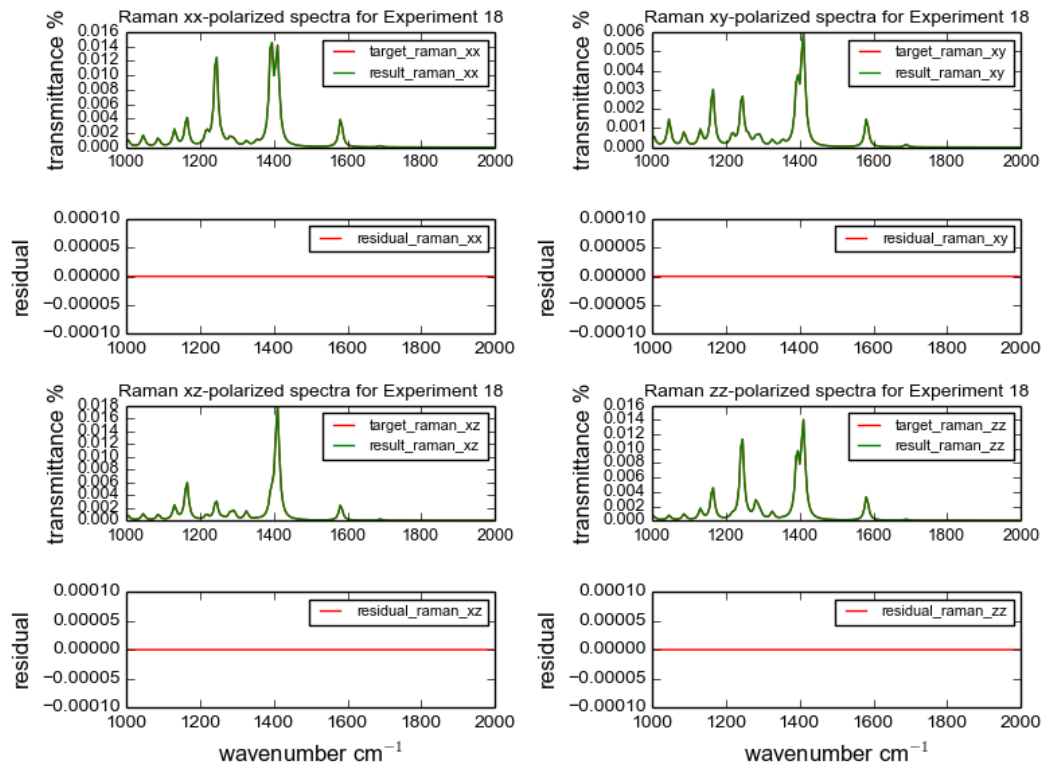


Figure 4.6: Raman spectra plotted by using target composition and return composition of Experiment 18

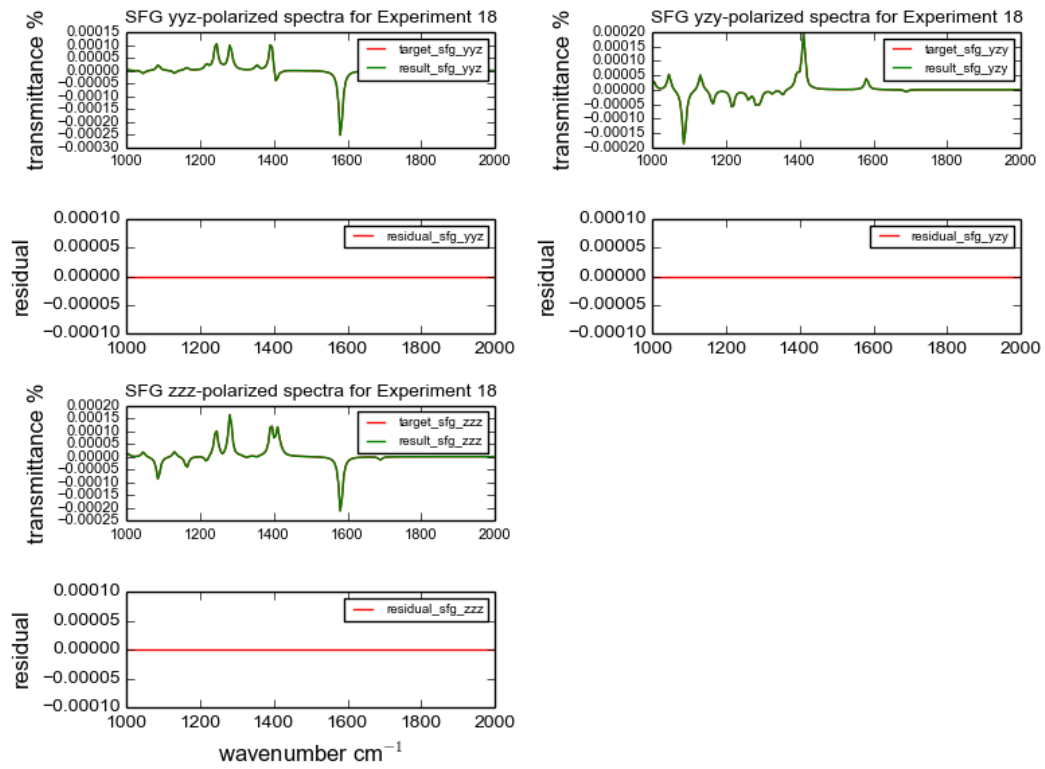


Figure 4.7: SFG spectra plotted by using target composition and return composition of Experiment 18

Chapter 5

Mixture

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's coordination distributions 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. Because of the numerical limitation, each Lp model returns an optimal composition solution. It seems that our LP models have hit their limitation for the case of one molecule. However, there is another factor that is valuable to explore: the coordination of different molecules at interfaces. For a mixture of different molecules at interface, we want to figure out whether our LP models can help to obtain the target composition. If the LP models success in obtaining the target composition, then the rate of accuracy is the key factor of the study as well.

5.2 Experiments

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

| Experiment Index | Spectrum Information |
|------------------|---|
| Experiment 1 | x and z polarized IR spectra |
| Experiment 2 | xx, xy, xz and zz polarized Raman spectra |
| Experiment 3 | yyz, yzy 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, yzy and zzz polarized SFG spectra |
| Experiment 6 | xx, xy, xz and zz polarized Raman spectra; yyz, yzy and zzz polarized SFG spectra |
| Experiment 7 | x and z polarized IR spectra; xx, xy, xz and zz polarized Raman spectra; yyz, yzy and zzz polarized SFG spectra |

Table 5.1: Detailed Experiment Group Setting

First, there are six different amino-acids in the mixture: methionine, leucine, isoleucine(ile), alanine, threonine and valine. For each amino acid, only θ difference is considered, the other two Euler angles are integrated. Each amino acid molecule has 9 candidates in the mixture, they have θ of the following values: 0° , 10° , 20° , 30° , 40° , 50° , 60° , 70° and 80° . Because when θ equals 90° , 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, and randomly generate a percentage for the selected candidate. The target composition is made of six randomly selected candidates coming from six different 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.

Each experiment in the experiment set contains different spectral information as shown in Table 5.1. In Experiment 1, candidates' IR x and z polarization spectra are obtained. The target's IR x and z polarization 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 following 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.

Each LP model of the experiments is using the same formula as shown in Equation 3.4. Because every experiment is built with different spectral information, each LP model is conducted with different spectral information.

Finally, this experiment set is run 100 times in order to see which experiment in Table 5.1 return 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.

5.3 Scoring methods

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 one. Therefore, this sum of residuals is negligible. It is not appropriate to use it as a scoring criteria.

Another way to measure the accuracy of the return composition, is to compare it directly with the target one. Therefore, calculating the sum of the residuals between a target composition and a return one directly is a faster 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 mock experiments where the target compositions are known in advance.

The return composition of each experiment in the experiment 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 $1e-7$. Then the return composition is considered to be the same as the target one.

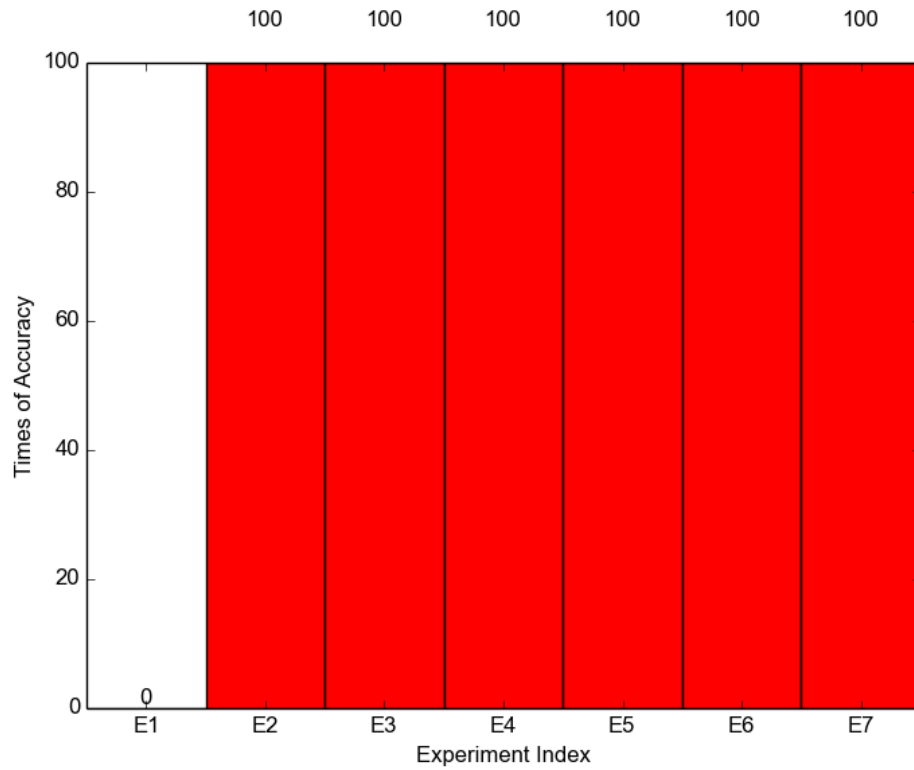


Figure 5.1: Accuracy analysis for experiments considering a mixture of amino acids with candidates from 0° to 80° on θ for each amino acid

The experiment set is ran for 100 times, the result is shown in Figure 5.1. In Experiment 2, the return composition of the LP model constructed by using Raman spectral information meets the target composition 100 times. This means that within this set of experiments, Raman is sufficient to obtain the correct composition of the target spectra. Moreover, the accuracy is 100%.

Experiment 3 is the LP model constructed by using SFG spectral information. Its accuracy is 100% as well. This indicates that SFG spectral information is as abundant

as Raman for this set of experiments.

Experiment 4 to 7, each experiment contains either spectral information of Raman, or SFG, or both. Therefore, the corresponding LP model can help to get the target composition with the same accuracy as Experiment 2 and Experiment 3.

The only exception is Experiment 1. The accuracy is not as high as the other experiments. Its accuracy of the LP model that constructed using IR spectral information is 0. The low performance can be caused by the insufficient spectral information of IR.

When this experiment set is re-run 100 times, only Experiment 1's returned composition 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 two polarizations' spectra conducted by the two different compositions are very close to each other in every run. For example, a random run is picked, then the two polarizations' spectra are plotted in Figure 5.2. The spectra plotted by the return composition are almost the same as the ones plotted by the target composition. 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 θ equals to one degree is identical to the one of candidate with this θ degree's complementary. For example, the Raman spectra for candidate with θ of 10° , is the same as candidate with θ of 170° . And for IR, it is the same case. Only SFG tells the differences between these two degrees, as the spectra for candidate with θ of one degree is symmetric to its complementary along wavenumber as shown in

Figure 5.8.

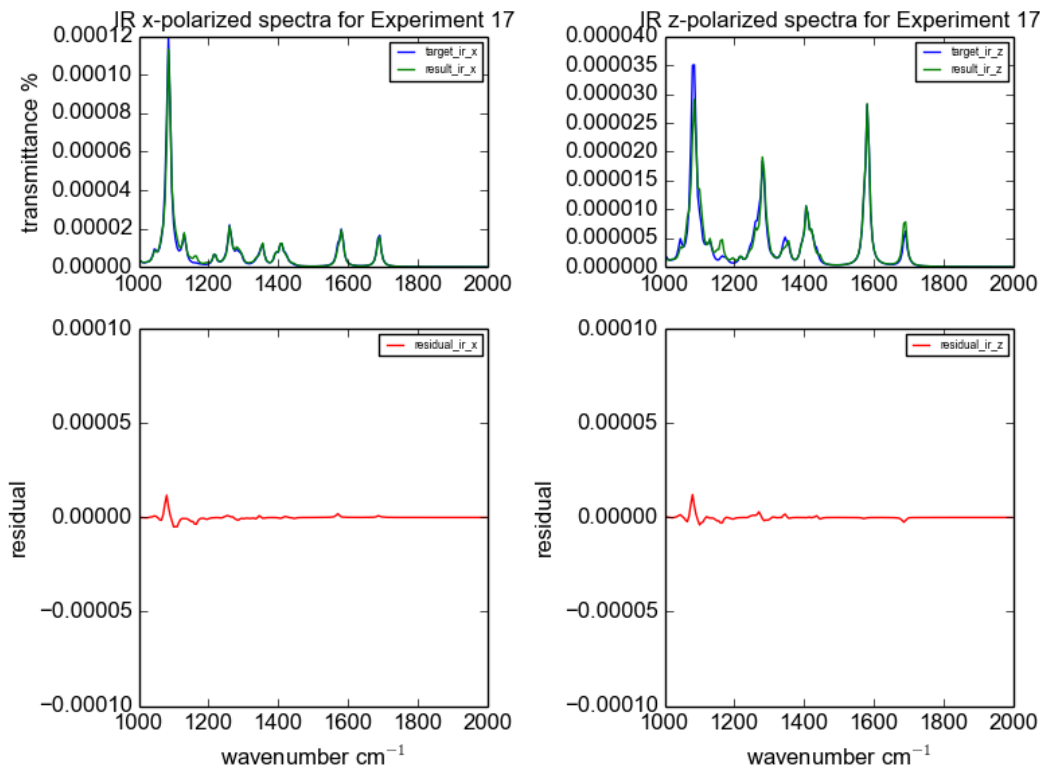


Figure 5.2: IR Spectra Plotted by Result Composition and Target Composition.

To further study the capacity of the LP models built for the mixture of molecules, the candidate pool is expanded from 0° to 180° in terms of the θ 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 randomly select one candidate from 18 candidates, instead of 9. All 108 candidates' IR, Raman and SFG spectra need to be generated. Figure 5.3 illustrates 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 θ of one degree is identical to the one with θ of this degree's complementary, as shown in Figure 5.7. This also increases the difficulty for the LP model constructed by using IR spectral information to return the target composition.

However, it should be noticed that the accuracy for Experiment 2 has dramatically

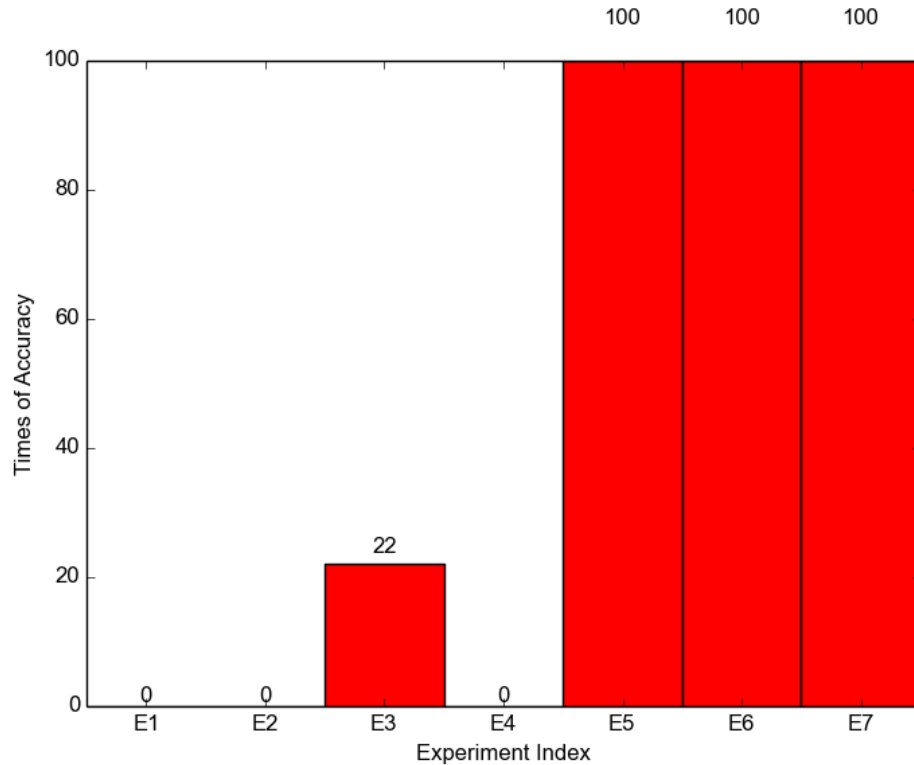


Figure 5.3: Accuracy analysis for experiments considering a mixture of amino acids with candidates from 0° to 180° on θ for each amino acid

dropped. This is because the Raman spectra for one candidate with a θ is identical to the one of this θ 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. From each projection, 200 data points are selected from the target and the candidates' spectra. Therefore, while the number of constraints is the same as before, the number of candidates are twice bigger than before. 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 spectral information is still insufficient to converge the composition to the target one.

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, Experiment 6, Experiment 7

all have 100% accuracies. This phenomenon can be explained as follow: SFG helps to distinguish a candidate from its complementary on θ 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 when each amino acid's candidates spread from 0° to 180° on θ . There are still some noticable result in the return composition: for each amino acid, the percentage assigned is correct; however, the candidate presented may be the one with the correct degree, or the one with the correct degree's complementary. For example, a random run is selected. Figure 5.4 displays the target composition and Figure 5.5 displays the return composition of Experiment 2. Figure 5.6 is the return composition of Experiment 6. From the three figures, when extracting the non-zero values to generate a list, the three lists are the same. However, when overlapping Figure 5.4 with Figure 5.5, the position of each non-zero value is not identical. For example, value of 0.299586 appears at θ of 150° for Valine in Figure 5.4. In Figure 5.5, it appears at θ of 30° for Valine. Same observation for values of 0.021196, 0.00662804, 0.000642609, and 0.00789 in Figure 5.4 and 5.5. The LP model of Experiment 2 fails to tell which candidate is the exact one between the correct one and its complementary on θ . This observation is a general case across all the experiment groups in the case of Experiment 2. From Experiment 6, as long as the spectra data from SFG is plugged into the LP model, the return composition is the same as the target one.

From the above analysis, Experiment 2 appears the ability of limiting the number of candidates to 2 for each amino acid. These two candidates are complementary on θ degree, with one of them to be the correct one for the target composition. 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. Spectral information from SFG is needed in order to study the cases that having θ expanded from 0° to 180° .

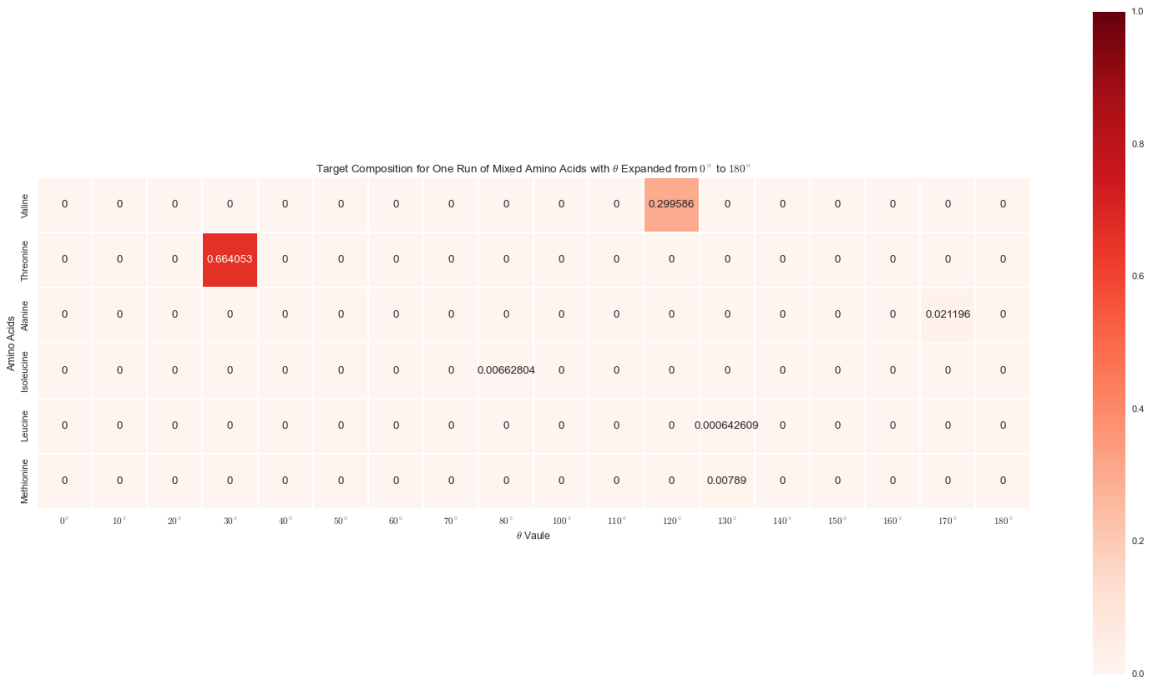


Figure 5.4: Target composition for one random run of six mixed amino acids with θ expanded from 0° to 180° on θ

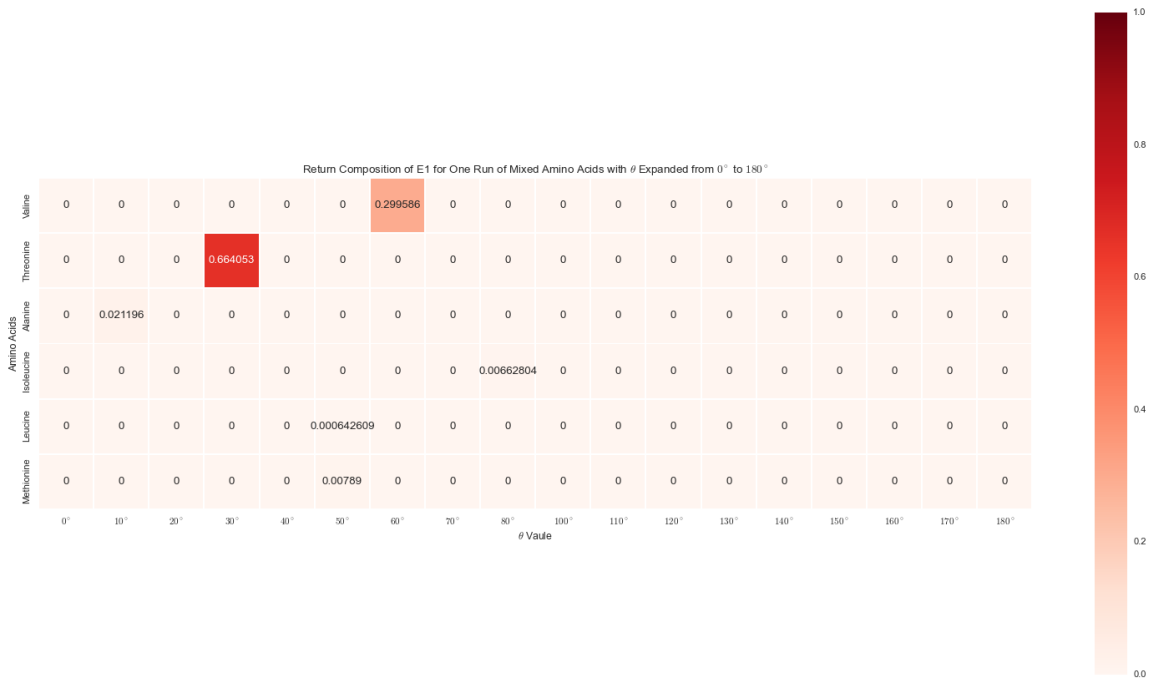


Figure 5.5: return composition of experiment 2 for one random run of six mixed amino acids with θ expanded from 0° to 180° on θ

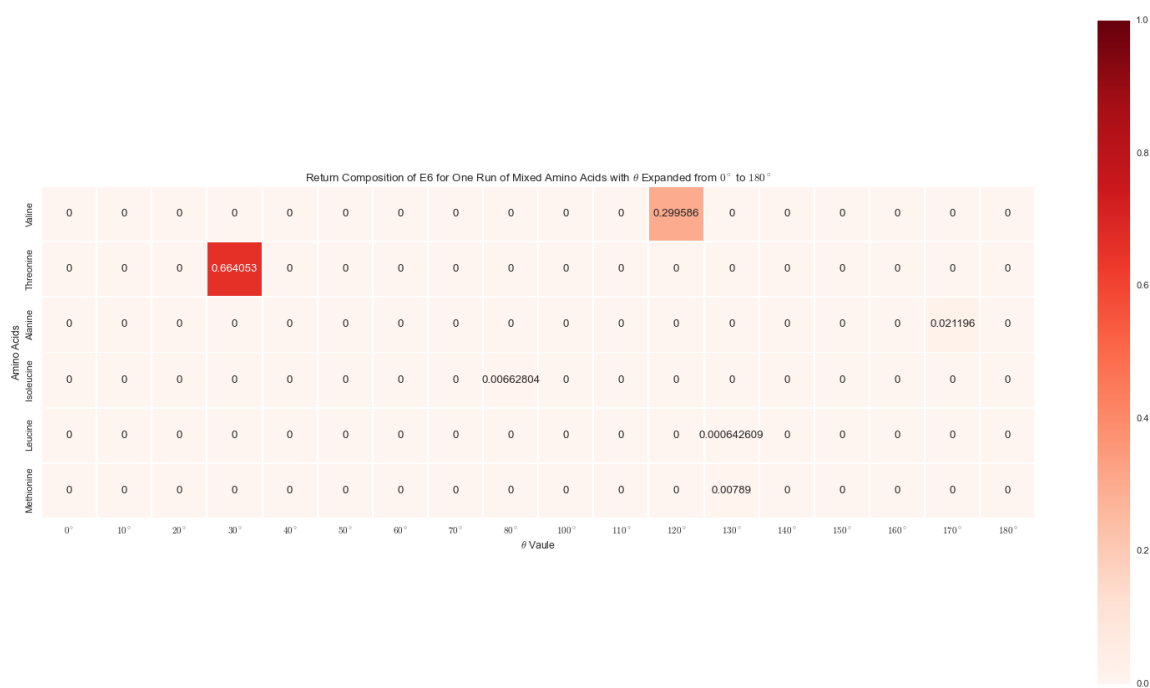


Figure 5.6: return composition of experiment 6 for one random run of six mixed amino acids with θ expanded from 0° to 180° on θ

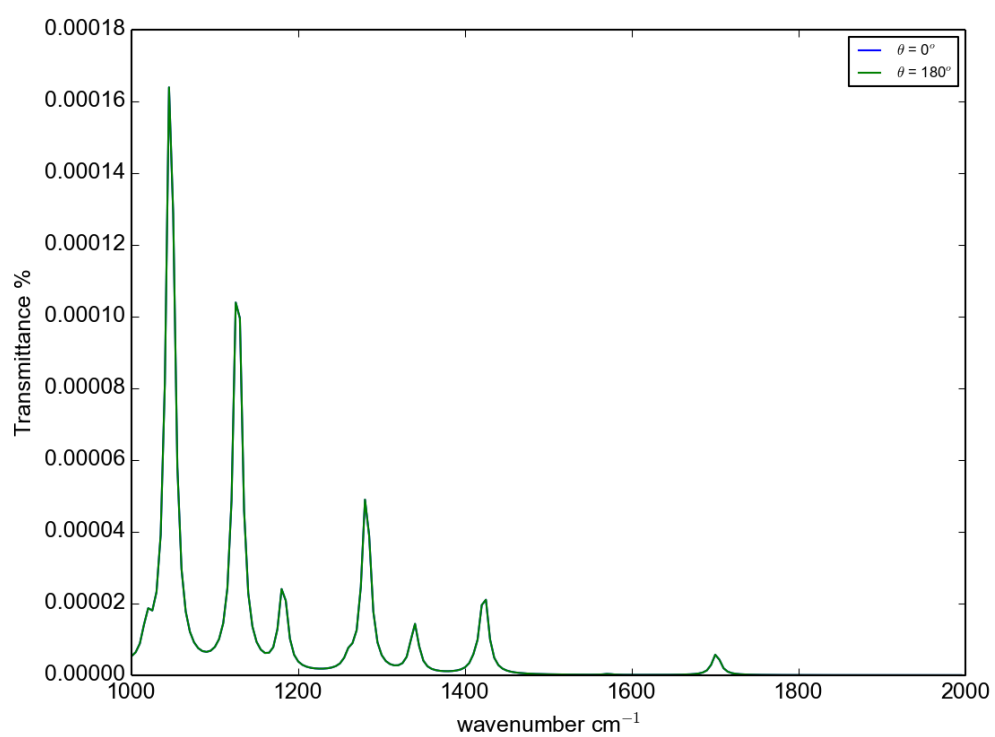


Figure 5.7: IR z projection spectrum for alanine candidate with θ of 0° is identical to alanine candidate with θ of 180°

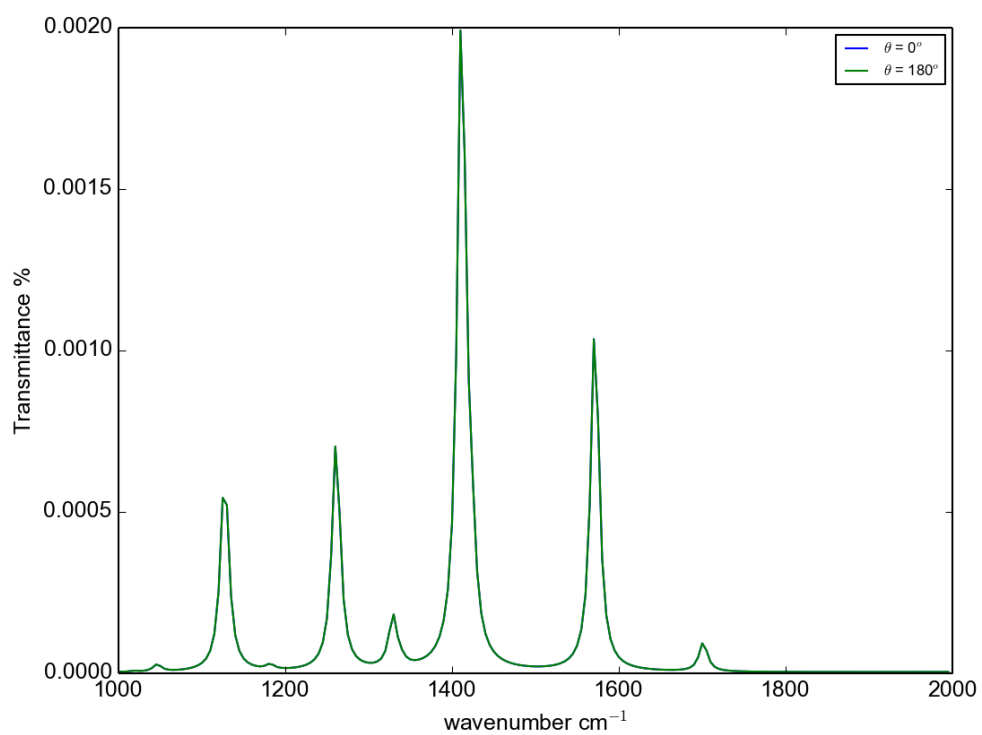


Figure 5.8: Raman zz projection spectrum for alanine candidate with θ of 0° is identical to alanine candidate with θ of 180°

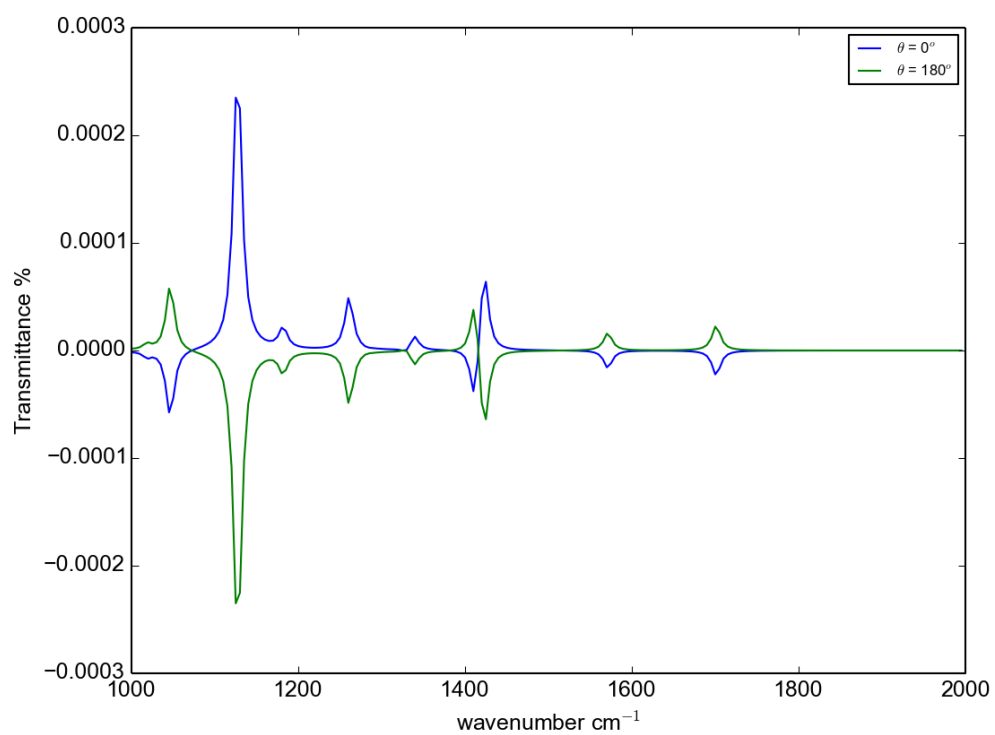


Figure 5.9: SFG *zzz* projection spectrum for alanine candidate with θ of 0° is not identical to alanine candidate with θ of 180° , but symmetric along wavelength

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 comparing 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 the LP models. The LP models constructed by Experiment 2 to 7 in Table 5.1 for θ ranged from 0° to 80°) are doing well in retrieving the target composition for the mixed amino acids. Therefore, based on these experiments, we would like to know if the same LP models can be applied directly to the real experimental data for the same θ range.

Therefore, the same experiment setting in Table 5.1 are used for the following experiments. The goal is the same, 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 is 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 right selected candidates. However, the percentage values of the candidates are different from the target composition. The ratio between the returned percentage and the target percentage are the same for the selected candidates. Furthermore, when this ratio adds up the extra variable, it equals 1. Randomly select one experiment run, then take Experiment 2 as an example. Figure 6.2 displays the target composition generated, 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 in the right position, however, each percentage value is different from the one in the target composition. There are one extra value in Figure 6.2 with value of 0.4.

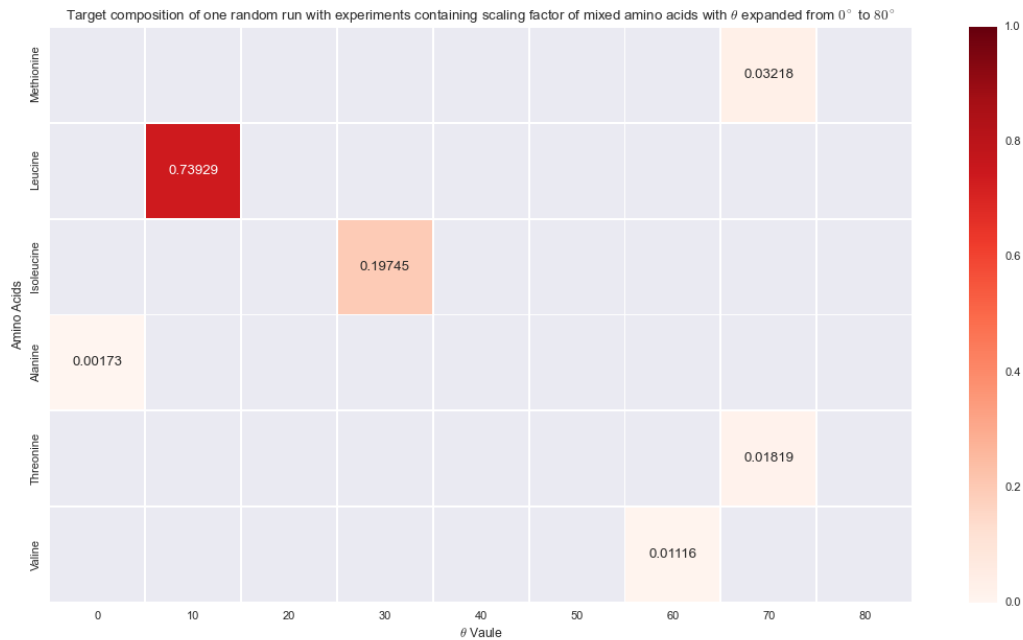


Figure 6.1: Target composition for one random run of experiment set with scaling factor for mixed amino acids with θ expended from 0° to 80°

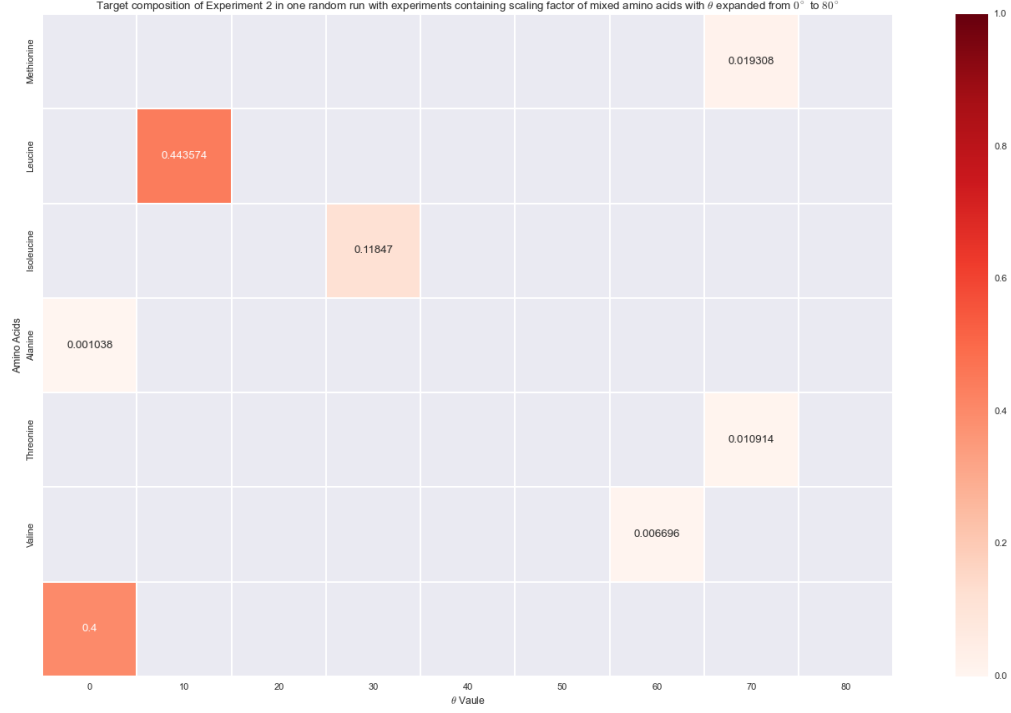


Figure 6.2: Return composition of Experiment 2 for one random run of experiment set with scaling factor for mixed amino acids with θ expended from 0° to 80°

Moreover, as 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 as slack variable 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. In conclusion, the slack variable (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.

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

To check whether the above observation is a general case, the experiment set in

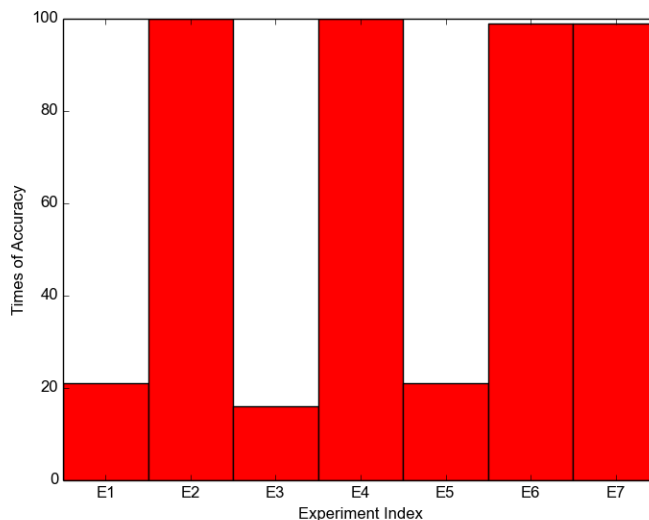


Figure 6.3: Experiment Accuracy Analysis for Experiments using experimental spectra data that contains scaling factor that is smaller than 1 and candidates with θ from 0° to 80°

Table 5.1 is run 100 times with randomly generated scaling factors in each run. Figure 6.3 indicates the experiment result. Figure 6.3 shows that Experiment 2, 4, 6 and 7 almost hit the above observation with 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. 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, the LP model with only SFG spectral information, 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. Even combining IR and SFG spectral information, the constructed LP model cannot help to re-construct the target composition.

The return slack variable and composition of Experiment 4, 6 and 7 are the same as Experiment 2 in each run. This indicates Raman spectral information are dominating the solution in the LP model.

When each amino acid's candidates are expanded from 0° to 180° on θ , the same experiment set is applied 100 times with randomly generated scaling factors in each run. The result in Figure ?? illustrates that none of the experiment in the set meets

the observation with high frequency.

When the return compositions of Experiment 2 and 6 are further analyzed, there are few observation need to be noticed. To facilitate the explanation, one random run is picked as an explicit example. Figure ?? is the target composition. Figure ?? is the return composition of Experiment 2. Figure ?? is the return composition of Experiment 6. In Figure ??, the slack variable still equals $1 - \text{generated scaling factor}$. For each amino acid, the selected candidate in the return composition may not be the exact one as the target composition selected. However, this selected candidate is always either the correct one in the target composition, or the correct one's θ complimentary. In Figure ??, for each amino acid, there are two selected candidates in the return composition. These two selected candidates are the correct one and its θ complimentary. When the percentages of these two selected candidates are added, it is equalled to the percentage returned for the amino acid in Figure ?. $0.0776716 + 0.0252641 = 0.102936$. Between these two selected candidates, the correct one's percentage is always bigger than its θ complimentary. $0.0776716 > 0.0252641$. In conclusion, Experiment 2 achieves in telling the slack variable, the scaling factor, and the ratio between the

The LP model with only Raman spectral information Experiment 2, is able to tell us for each amino acid, candidate with which θ or this θ 's complementary would exist in our target spectra. Because Raman spectra for candidate with θ on one degree is the same as its supplementary. This LP model can not distinguish between candidate with θ and the one with its complementary. However, with the help of SFG, we may be able to know which one between the above two dominates one amino acid's the total fraction, like what we have learnt from Chapter 5 about E6. Therefore we exam E6 here, and it displays which one of the two takes the major composition in the target spectra. With this information, we can decide between the θ and its complementary. With the information coming from E2 and E6, we therefore, obtain the right composition of the target spectrum.

Here goes the example, in one run of the experiment group. The composition for the target spectrum is Array 6.2. In this target spectrum, we have 0.14799 of $\theta = 40^\circ$ Methionine, 0.74202 of $\theta = 50^\circ$ Leucine, 0.08989 of $\theta = 150^\circ$ Ile, 0.01135 of $\theta = 40^\circ$ Ala, 0.00715 of $\theta = 0^\circ$ Thr, 0.0016 of $\theta = 20^\circ$ Val.

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0.14799 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.74202 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.08989 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.01135 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.00715 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0.0016 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \quad (6.2)$$

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0.102936 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.516118 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.0625238 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.0078945 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.00497324 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0.00111289 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.304441 \end{bmatrix} \quad (6.3)$$

The result returned by E2 is shown in Array 6.3. You may notice same as Array ?? to Array ??, Array 6.3 contains one more value than Array 6.2, 0.304441, which is the slack variable. We already know that the scaling factor is 0.695560510845 (generated randomly, but recorded). When we add the slack variable and the scaling factor, the total comes up to 1.0.

In Array 6.3, we get 0.102936 of $\theta = 40^\circ$ Methionine, 0.516118 of $\theta = 50^\circ$ Leucine, 0.0625238 of $\theta = 30^\circ$ Ile, 0.0078945 of $\theta = 40^\circ$ Ala, 0.00497324 of $\theta = 0^\circ$ Thr, 0.00111289 of $\theta = 20^\circ$ Val. From Array 6.4, we can also deduce the value for the scaling factor.

$$\begin{aligned} \frac{0.102936}{0.14799} &= \frac{0.516118}{0.74202} = \frac{0.0625238}{0.08989} = \frac{0.0078945}{0.01135} \\ &= \frac{0.00497324}{0.00715} = \frac{0.00111289}{0.0016} = 0.695560 \end{aligned} \quad (6.4)$$

At first glance, we may guess that this LP model actually return the correct composition. However, not all amino acids' composition is correct. For Ile, it should be 0.0625238 of $\theta = 150^\circ$, but the result returned 0.0625238 of $\theta = 30^\circ$, which is the

complimentary of 150° . It is because these two degrees' Raman spectra are identical, there is no way for current LP model to distinguish these two.

With this information, we know the only thing we need to make sure is: is the θ returned by LP model the exact one in target spectrum or its complementary? (The above conclusion can also be applied to the experiments in Chapter 5 without the scaling factor. The LP model with only Raman spectra information can help us to see which θ of candidate and its complimentary should be for each amino acid. I did not observe this before.)

To answer this question, we need the help of SFG data. Because only SFG can tell us the difference between one angle and its complementary, as their spectra are symmetry, not identical around the axis of wevenumber.

In this experiment group, the result returned E6 is Array 6.5 (Second example???)

$$\left[\begin{array}{cccccccccccccccc} 0 & 0 & 0 & 0 & 0.0776716 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.0252641 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.3894440 & 0 & 0 & 0 & 0 & 0 & 0 & 0.1266740 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.0153456 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.0471782 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.00595697 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.00193762 & 0 & 0 & 0 & 0 \\ 0.0037526 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.00122061 & 0 \\ 0 & 0 & 0.000839749 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.000273144 & 0 & 0 \end{array} \right] \quad (6.5)$$

0.304441(anyrelationforthefraction???)

The value of the slack variable is the same as what returned by E2. However, the returned composition is totally different than what returned by E2. The interesting thing is that for each amino acid, the return existing candidates are complementary on θ . The total percentage of these two candidates, take Methionine as an instance, $0.0776716 + 0.0252641$ makes 0.1029357 which is the same as what is returned by E2. This is the same for every amino acid. What's more, the composition returned by the E6 indicates which θ dominates the composition for one amino acid. For Methionine, $\theta = 40^\circ$ takes major part; for Leucine, $\theta = 50^\circ$ does; for Ile, $\theta = 30^\circ$ does; for Ala, $\theta = 40^\circ$ does; for Thr, $\theta = 0^\circ$ does; for Val, $\theta = 20^\circ$ does; And those candidates are the correct components for target spectra.

IR+SFG, can you do anything with it???

6.2 Results

6.3 Discussion

6.4 Conclusions

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!

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

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

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

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

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