

МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ РОССИЙСКОЙ ФЕДЕРАЦИИ
МОСКОВСКИЙ ФИЗИКО-ТЕХНИЧЕСКИЙ ИНСТИТУТ (государственный университет)
ФИЗТЕХ-ШКОЛА ПРИКЛАДНОЙ МАТЕМАТИКИ И ИНФОРМАТИКИ
КАФЕДРА «ИНТЕЛЛЕКТУАЛЬНЫЕ СИСТЕМЫ»

Панченко Святослав Константинович

Геометрическая алгебра для решения задачи декодирования сигналов

03.03.01 — Прикладные физика и математика

ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА МАГИСТРА

Научный руководитель:
д.ф.-м.н. Стрижов Вадим Викторович

Москва
2020

Содержание

Введение	4
1 Постановка задачи восстановления плотности	6
2 Решение задачи восстановления плотности	8
2.1 Распределение Кента как способ описания распределения углов	8
2.2 Выбор плотности распределения компоненты смеси	10
2.3 Алгоритм нахождения оптимальных параметров смеси	11
2.4 Обновление параметров распределения Кента	12
2.4.1 Стохастическая модификация ЕМ-алгоритма	12
2.4.2 Моментные оценки параметров распределения Кента	13
2.4.3 Аналитические формулы для моментных оценок	13
2.5 Определение числа компонент в модели смеси распределений	14
2.6 Инициализация параметров смеси в алгоритме	15
2.7 Окончательный вид алгоритма поиска оптимальных параметров	16
3 Вычислительный эксперимент	17
3.1 Экспериментальные данные	17
3.2 Восстановление плотности распределения пространственных конфигураций пары ALA-C _{ar}	17
3.2.1 Иллюстрации для $r = 7\text{\AA}$	19
3.2.2 Иллюстрации для $r = 11\text{\AA}$	20
3.2.3 Иллюстрации для $r = 15\text{\AA}$	21
3.3 Установление соответствия с другими моделями	22
Заключение	23
Список литературы	25

[ДАЛЕЕ СЛЕДУЕТ ШАБЛОН, АКТУАЛЬНЫЙ ТЕКСТ РАБОТЫ БУДЕТ
ПРЕДОСТАВЛЕН 31-ОГО МАЯ]

Аннотация

Рассмотрена задача восстановления плотности распределения трёхмерного случайного вектора, две компоненты которого представляют собой пару сферических углов. Требуется, чтобы полученные плотности были интерпретируемы с точки зрения эксперта, согласовывались с ранее полученными результатами, а модель восстановления учитывала периодичность углов. Предлагается рассматривать значения пары углов как реализации случайного вектора, областью значений которого является сфера в трёхмерном пространстве. Искомая плотность в таком подходе моделируется в виде смеси, в каждой компоненте которой углы распределены в соответствии с распределением Кента. Параметры смеси находятся с помощью модификации алгоритма Stochastic Expectation-Maximization. Проведено восстановление плотностей распределения пространственных ориентаций различных пар молекул. Демонстрируется, что результаты восстановления согласуются с мнением эксперта и результатами других моделей.

Ключевые слова: *трёхмерная структура белка, восстановление плотности распределения, модель смеси распределений, распределение Кента, алгоритм Stochastic Expectation-Maximization.*

Введение

Актуальность темы. Трёхмерная структура белковой молекулы — это ключ к пониманию её биологических функций и свойств. Однако, определение строения белковой цепи, обычно с помощью рентгеновской кристаллографии или спектроскопии ядерного магнитного резонанса, весьма дорого и трудоёмко. Поэтому число экспериментально определённых белковых структур существенно меньше числа идентифицированных белковых цепочек. Отсюда возникает задача предсказания по последовательности образующих молекулу компонент её трёхмерной структуры. С проблемой можно ознакомиться в работах [1, 2].

Решение данной задачи — одна из самых важных целей биоинформатики и теоретической химии. Оно позволит значительно улучшить существующие генеративные и предсказательные модели в области исследования молекулярных последовательностей. Полученные при помощи этих моделей данные активно используются в медицине и биотехнологии.

Существующие решения. Один из фундаментальных подходов к решению данной задачи — построение потенциала, функции, минимумы которой соответствуют энергетически устойчивым конфигурациям молекул, образующих химическую связь. Имеются два основных подхода к построению таких потенциалов:

Физический подход: в этом подходе потенциал строится на основе законов молекулярной химии, описывающих взаимодействия молекул. Такие потенциалы, как правило, точны, но их вычисление весьма трудоёмко. Они лучше подходят для описания одной конкретной цепочки и мало применимы для анализа произвольных последовательностей. Примеры использования подобных подходов рассмотрены в исследованиях [3, 4].

Статистический подход: в этом подходе предполагается стохастическая модель порождения данных: для каждой пары потенциально взаимодействующих молекул на основе известных и изученных молекулярных структур строятся функции плотности совместной вероятности параметров химической связи, определяющих взаимную пространственную ориентацию этих молекул. С помощью полученных плотностей формируются статистические потенциалы, описывающие структуру неизученной молекулярной цепи на основе вероятностного обобщения известных структур. Многочисленные исследования в этой области представлены в статьях [5–10].

Второй, статистический, подход существенно опирается на оценку совместных плотностей распределения величин, характеризующих молекулярную связь. Для определённого набора хорошо изученных молекул существуют базы данных, содержащие информацию о параметрах химических связей, которую эти молекулы формировали между собой в исследованных структурах. Одна пара таких молекул характеризуется десятками тысяч зарегистрированных конфигураций с различными параметрами. Такое множество конфигураций объясняется тем, что рассматриваемая пара молекул входила в состав громадного количества различных молекулярных последовательностей, каждый элемент которых мог повлиять на значения этих параметров.