

УДК 575.112
ББК 28.071.3
К63

Певзнер П., Компо Ф.

К63 Алгоритмы биоинформатики / пер. с англ. И. Л. Люско. – М.: ДМК Пресс, 2023. – 682 с.: ил.

ISBN 978-5-93700-175-7

Перед вами одно из самых популярных за рубежом руководств по биоинформатике. Книга обеспечивает уникальный баланс между практическими задачами современной биологии и фундаментальными алгоритмическими идеями. Каждая глава начинается с биологического вопроса, а затем неуклонно развивается алгоритмическая сложность, необходимая для ответа на него. Сотни упражнений включены непосредственно в текст и помогают разобраться в непростом материале.

Издание предназначено специалистам в области анализа данных, а также будет полезно ученым, инженерам, студентам и аспирантам, работающим на стыке биологии и информатики.

УДК 575.112
ББК 28.071.3

Все права защищены. Любая часть этой книги не может быть воспроизведена в какой бы то ни было форме и какими бы то ни было средствами без письменного разрешения владельцев авторских прав.

ISBN 978-5-93700-175-7 (рус.)

Copyright © 2015 by Phillip Compeau
and Pavel Pevzner
© Перевод, оформление, издание,
ДМК Пресс, 2022

Содержание

От издательства	16
------------------------------	----

Глава 1. В каком месте генома начинается репликация ДНК?	17
Путешествие в тысячу миль.....	18
Скрытые сообщения в точке начала репликации	20
DnaA-боксы	20
Скрытые сообщения в «Золотом жуке»	21
Подсчет слов	22
Задача поиска часто встречающихся слов	23
Более быстрый подход к задаче частых слов	24
Часто используемые слова в <i>Vibrio cholerae</i>	26
Некоторые скрытые сообщения более примечательны, чем другие	27
Взрыв скрытых сообщений.....	31
Поиск скрытых сообщений в нескольких геномах.....	31
Задача поиска сгустков	32
Самое простое объяснение процесса репликации ДНК	34
Асимметрия репликации	37
Специфическая статистика прямой и обратной полуцепей	41
Неизвестный биологический феномен или статистическая случайность?	41
Дезаминирование	43
Диаграмма смещения	44
Некоторые скрытые сообщения более неуловимы, чем другие.....	47
Последняя попытка найти DnaA-боксы в <i>E. Coli</i>	51
Эпилог. Осложнения в предсказаниях <i>ori</i>	53
Открытые проблемы	55
Множественные точки начала репликации в бактериальном геноме.....	55
Поиск источников репликации у архей	57
Поиск точек начала репликации у дрожжей.....	59
Вычисление вероятностей паттернов в строке	60
Зарядные станции	61
Массив частот	61
Преобразование Patterns в Numbers и наоборот	63
Поиск часто встречающихся слов путем сортировки.....	65

6 СОДЕРЖАНИЕ

Решение задачи поиска сгустков	66
Решение задачи часто встречающихся слов с несовпадениями	68
Генерация окрестности строки	69
Поиск часто встречающихся слов с несовпадениями путем сортировки.....	71
Сопутствующие материалы	72
Оценка «О большого» (Big-O).....	72
Вероятности паттернов в строке	73
Самый красивый эксперимент в биологии.....	78
Направленность цепей ДНК.....	80
Ханойские башни.....	81
Парадокс перекрывающихся слов	83
Библиографические примечания.....	85
Глава 2. Какие сегменты ДНК играют роль молекулярных часов?	87
Есть ли у нас «часовой ген»?	88
Найти мотив сложнее, чем вы думаете.....	89
Идентификация вечернего элемента	89
Игра в прятки с мотивами.....	90
Метод грубой силы поиска мотива	92
Считаем мотивы	93
От мотивов к матрицам профиля и консенсусным строкам	93
На пути к более адекватной функции оценки мотивов	96
Энтропия и motif logo	97
От поиска мотива к поиску медианной строки.....	98
Задача поиска мотива.....	98
Переформулировка задачи поиска мотива.....	99
Задача поиска медианной строки.....	101
Почему мы переформулировали задачу поиска мотива?	103
Жадный алгоритм поиска мотива.....	104
Использование матрицы профиля для бросания костей	104
Анализ жадного алгоритма поиска мотива	106
Поиск мотива и Оливер Кромвель.....	107
Какова вероятность того, что завтра не взойдет солнце?	107
Правило преемственности Лапласа.....	108
Улучшенный алгоритм жадного поиска мотивов	109
Рандомизированный поиск мотива.....	112
Игра в кости для поиска мотивов	112
Почему рандомизированный поиск мотивов работает	114
Почему рандомизированный алгоритм работает так хорошо?	116
Сэмплирование по Гиббсу	119
Сэмплирование по Гиббсу в действии	121
Эпилог. Как туберкулез впадает в спячку, чтобы спрятаться от антибиотиков?	124
Зарядная станция	127
Решение задачи медианной строки	127
Сопутствующие материалы	128
Экспрессия генов	128
ДНК-чибы	128
Игла Бюффона	129

Сложности в поиске мотива	132
Относительная энтропия	132
Библиографические примечания	134
Глава 3. Как мы собираем геномы?	135
Взрывающиеся газеты.....	136
Задача реконструкции строки	139
Сборка генома сложнее, чем вы думаете	139
Реконструкция строк из k-меров	139
Повторы усложняют сборку генома.....	142
Реконструкция строк как прогулка по графу перекрытий	143
От строки к графу.....	143
Геном исчезает	146
Два способа представления графов	147
Гамильтоновы пути и универсальные строки	148
Другой граф для реконструкции строк	150
Склейивание узлов и графы де Брюйна	150
Прогулка по графу де Брюйна.....	152
Эйлеровы пути	152
Другой способ построения графов де Брюйна.....	153
Построение графов де Брюйна из композиции k-меров	155
Графы де Брюйна в сравнении с графиками перекрытия	156
Семь мостов Кенигсберга	157
Теорема Эйлера.....	160
От теоремы Эйлера к алгоритму нахождения эйлеровых циклов	163
Построение эйлеровых циклов	163
От эйлеровых циклов к эйлеровым путям.....	164
Создание универсальных строк	165
Сборка геномов из рид-пар	167
От ридов к рид-парам.....	167
Преобразование рид-пар в длинные виртуальные риды	169
От композиции к спаренной композиции.....	170
Парные графы графы де Брюйна	172
Ловушка парных графов де Брюйна	173
Эпилог. Сборка генома – работа с реальными данными секвенирования.....	176
Разбиваем риды на k-меры	176
Фрагментация генома на контиги	177
Сборка ридов с возможными ошибками	179
Определение кратности ребер в графах де Брюйна	180
Зарядные станции	181
Влияние склейки на матрицу смежности	181
Генерация всех эйлеровых циклов	182
Реконструкция строки, записанной как путь в парном графе де Брюйна.....	184
Максимальные неветвящиеся пути в графике	186
Сопутствующие материалы	187
Краткая история технологий секвенирования ДНК	187
Повторы в геноме человека	189
Графы	190
Игра «Икосиан»	193
Разрешимые и неразрешимые задачи	194

8 СОДЕРЖАНИЕ

От Эйлера до Гамильтона и де Брюйна	195
Семь мостов Калининграда.....	196
Подводные камни сборки двухцепочечной ДНК.....	197
«ЛУЧШАЯ» теорема.....	198
Библиографические примечания.....	199
Глава 4. Как мы секвенируем антибиотики?	201
Открытие антибиотиков	202
Как бактерии производят антибиотики?.....	203
Как пептиды кодируются геномом.....	203
Где в геноме <i>Bacillus brevis</i> закодирован тироцидин?	206
От линейных к циклическим пептидам.....	207
Уклоняясь от центральной догмы молекулярной биологии	208
Секвенирование антибиотиков путем их дробления на части	209
Введение в масс-спектрометрию.....	209
Задача секвенирования циклопептидов	210
Алгоритм грубой силы для секвенирования циклопептидов	212
Алгоритм ветвей и границ для секвенирования циклопептидов.....	214
Масс-спектрометрия и гольф.....	217
От теоретических к реальным спектрам	217
Адаптация секвенирования циклопептидов для спектров с ошибками	218
От 20 до более чем 100 аминокислот.....	222
Спектральная свертка спасает положение	223
Эпилог. От смоделированных спектров – к реальным	227
Зарядные станции	229
Создание теоретического спектра пептида	229
На сколько быстро выполняется алгоритм CyclopeptideSequencing?	231
Сокращение списка пептидов Leaderboard	232
Сопутствующие материалы	233
Гаузе и лысенковщина	233
Открытие кодонов	235
Чувство кворума	236
Молекулярная масса	236
Сленоцистейн и пирролизин	237
Псевдополиномиальный алгоритм для Теоремы магистрали	237
Расщепленные гены.....	238
Библиографические примечания	240
Глава 5. Как мы сравниваем участки ДНК?	241
Взлом неривбосомного кода.....	242
Клуб галстуков РНК.....	242
От сравнения белков к неривбосомному коду.....	242
Что общего между онкогенами и факторами роста?	244
Введение в выравнивание последовательностей.....	245
Выравнивание последовательности как игра	245
Выравнивание последовательностей и самая длинная общая подследовательность	247
Туристическая задача Манхэттена.....	248
Какова наилучшая стратегия осмотра достопримечательностей?	248

Достопримечательности в произвольном ориентированном графе	252
Выравнивание последовательности – это замаскированная туристическая задача Манхэттена	253
Введение в динамическое программирование: задача размена монет	257
Жадный обмен денег	257
Рекурсивный размен денег	258
Размениваем деньги с помощью динамического программирования	259
Новый взгляд на туристическую задачу Манхэттена	261
От Манхэттена к произвольному DAG	266
Выравнивание последовательности как построение графа в стиле Манхэттена	266
Динамическое программирование в произвольном графе DAG	267
Топологические порядки	269
Возвращаясь к графу выравнивания	274
Считаем выравнивания	277
Что не так с моделью LCS?	277
Матрицы счета	279
От глобального к локальному выравниванию	280
Глобальное выравнивание	280
Ограничения глобального выравнивания	281
Бесплатные поездки на такси в графе выравнивания	284
Меняющиеся грани выравнивания последовательности	287
Задача 1. Расстояние редактирования	287
Задача 2. Настройка выравнивания	288
Задача 3. Выравнивание с перекрытием	289
Штрафы за вставки и удаления при выравнивании последовательности	290
Штрафы за аффинные пробелы	290
Строительство графа Манхэттена на трех уровнях	293
Компактное выравнивание последовательности	296
Вычисление счета выравнивания с использованием линейной памяти	296
Задача среднего узла	298
Удивительно быстрый и экономичный алгоритм выравнивания	301
Задача среднего ребра	303
Эпilog. Множественное выравнивание последовательностей	305
Построение трехмерного Манхэттена	305
Жадный алгоритм множественного выравнивания	307
Сопутствующие материалы	310
Светлячки и нерибосомный код	310
Поиск LCS без постройки города	311
Построение топологической сортировки	312
Матрица счета PAM	313
Алгоритмы «разделяй и властвуй»	314
Счет множественных выравниваний	316
Библиографические примечания	318
Глава 6. Есть ли в человеческом геноме «хрупкие» области?	319
О мышах и людях	320
Насколько различаются гены человека и мыши?	321
Синтенные блоки	321

10 СОДЕРЖАНИЕ

Реверсии	322
Точки перестановки.....	324
Модель эволюции хромосом со случайными разрывами.....	325
Сортировка по реверсиям.....	328
Жадный алгоритм сортировки по реверсиям	332
Точки останова.....	334
Что такое точки останова?.....	334
Счет точек останова	335
Сортировка по реверсиям для устранения точек останова	336
Рекомбинации в геномах опухолей.....	338
От монохромосомных к мультихромосомным геномам.....	339
Транслокации, слияния и расщепления.....	339
От генома к графу	340
Двойные разрывы	341
Графы точек останова.....	344
Вычисление дистанции двойного разрыва	347
Горячие точки рекомбинации в геноме человека	350
Модель случайных разрывов соответствует теореме о дистанции двойного разрыва	350
Модель хрупких разрывов	351
Эпилог. Конструирование синтенных блоков	353
Геномные точечные диаграммы и общие k-меры.....	353
Поиск общих k-меров	354
Построение синтенных блоков из общих k-меров	357
Синтенные блоки как связные компоненты в графах	359
Зарядные станции	363
От геномов к графу точек останова	363
Решение задачи сортировки по двойным разрывам	366
Сопутствующие материалы	368
Почему генный состав X-хромосом так консервативен?	368
Открытие геномных рекомбинаций	368
Экспоненциальное распределение.....	369
Сортировка блинов Билла Гейтса и Дэвида Х. Коэна	370
Сортировка линейных перестановок по реверсиям	371
Библиографические примечания	373
Глава 7. Какое животное заразило нас коронавирусом?.....	375
Самая быстрая вспышка	376
Проблемы в отеле «Метрополь».....	376
Эволюция SARS	376
Преобразование матриц расстояний в эволюционные деревья	378
Построение матрицы расстояний из геномов коронавируса	378
Эволюционные деревья в виде графов	379
Построение филогенетии по расстояниям	383
На путях к алгоритму построения филогенетии по расстоянию	386
В поисках соседних листьев	386
Вычисление длины ветвей	388
Аддитивная филогенетия	391

Обрезка дерева.....	391
Прикрепление ветви.....	392
Алгоритм реконструкции филогении по расстоянию.....	393
Построение эволюционного дерева коронавирусов	394
Использование метода наименьших квадратов для построения приблизительных филогений.....	395
Ультраметрические эволюционные деревья.....	397
Алгоритм объединения соседей	402
Преобразование матрицы расстояний в матрицу объединения соседей	402
Анализ коронавирусов с помощью алгоритма объединения соседей	406
Ограничения методов реконструкции эволюционного дерева по расстояниям....	408
Реконструкция эволюционного дерева по признакам	408
Таблицы признаков	408
От анатомических к генетическим признакам	409
Сколько раз эволюция изобретала крылья для насекомых?.....	410
Задача минимального показателя экономии.....	411
Задача максимальной экономии.....	418
Эпилог. Эволюционные деревья в борьбе с преступностью.....	425
Сопутствующие материалы	426
Когда HIV перешел от приматов к человеку?.....	426
Поиск дерева с помощью настройки матрицы расстояний.....	427
Условие четырех точек.....	429
Заразили ли нас атипичной пневмонией летучие мыши?	430
Почему алгоритм объединения соседей работает?.....	432
Вычисление длин ветвей в алгоритме объединения соседей.....	436
Большая панда: медведь или енот?	437
Откуда пришли люди?	439
Библиографические примечания	441
Глава 8. Как дрожжи научились делать вино?	443
Эволюционная история виноделия	444
Как давно мы зависим от алкоголя?	444
Диауксический сдвиг	445
Идентификация генов, ответственных за диауксический сдвиг	445
Две эволюционные гипотезы с разными судьбами	445
Какие гены дрожжей вызывают диауксический сдвиг.....	446
Введение в кластеризацию	447
Анализ экспрессии генов	447
Кластеризация генов дрожжей	451
Принцип правильной кластеризации.....	452
Кластеризация как задача оптимизации.....	454
Самый дальний первый обход.....	456
Самый дальний первый обход	456
Кластеризация k -средних	458
Искажение квадрата ошибки	458
Кластеризация k -средних и центр тяжести	460
Алгоритм Ллойда.....	462
От центров к кластерам и обратно	462
Инициализация алгоритма Ллойда	465

12 СОДЕРЖАНИЕ

Инициализатор k-means++	466
Кластеризация генов, вовлеченных в диауксический сдвиг	466
Ограничения кластеризации k -средних	468
Ограничения кластеризации k -средних	468
От подбрасывания монеты к кластеризации k -средних	470
Подбрасывание монет с неизвестной симметрией.....	470
В чем же состоит вычислительная задача?	473
От подбрасывания монеты к алгоритму Ллойда.....	474
Вернемся к кластеризации.....	476
Принятие мягких решений при подбрасывании монет	477
Максимизация ожиданий: Е-шаг.....	477
Максимизация ожиданий: М-шаг.....	478
Алгоритм максимизации ожидания.....	480
Мягкая кластеризация k -средних.....	480
Применение алгоритма максимизации ожидания к кластеризации	480
От центров к мягким кластерам	480
От мягких кластеров к центрам.....	483
Иерархическая кластеризация	484
Введение в кластеризацию по расстояниям	484
Определение кластеров по структуре дерева	487
Анализ диауксического сдвига с иерархической кластеризацией.....	490
Эпилог. Кластеризация образцов опухоли	493
Сопутствующие материалы	494
Полногеномная дупликация или серия дупликаций?.....	494
Измерение экспрессии генов	495
ДНК-микрочипы	496
Доказательство теоремы о центре тяжести	496
Матрица экспрессии генов и матрица расстояний/сходств	498
Кластеризация и испорченные клики	499
Библиографические примечания	501
Глава 9. Как мы обнаруживаем локацию	
болезнетворных мутаций?	503
Что вызывает синдром Одо?.....	504
Введение во множественное выравнивание последовательностей	505
Объединение Patterns в префиксное дерево	506
Построение префиксного дерева Trie.....	506
Применение префиксного дерева к множественному выравниванию	508
Предварительная обработка генома как альтернатива	511
Суффиксные попытки (suffix tries)	511
Использование суффиксных попыток для сопоставления последовательностей.....	512
Суффиксные деревья (suffix trees)	514
Суффиксные массивы.....	518
Выравнивание паттерна с суффиксным массивом	519
Преобразование Барроуза–Уилера.....	521
Сжатие генома.....	521
Построение преобразования Барроуза–Уилера.....	521
От повторов к сериям	523
Первая попытка инвертирования преобразования Барроуза–Уилера.....	524

Свойство «первый–последний» и инвертирование преобразования	
Барроуза–Уилера527
Свойство «первый–последний»527
Использование свойства «первый–последний» для инвертирования	
преобразования Барроуза–Уилера530
Сопоставление последовательностей с помощью преобразования	
Барроуза–Уилера534
Первая попытка сопоставления паттернов Барроуза–Уилера534
Перемещение по последовательности назад.....	.535
Маппинг «последний–первый»537
Делаем сопоставление паттернов по Барроуз–Уилеру быстрее539
Замена маппинга «последний–первый» оценочными массивами539
Удаление первого столбца матрицы Барроуза–Уилера542
Где находятся совпадающие паттерны?543
Барроуз и Уиллер устанавливают контрольные точки545
Эпилог. Устойчивое к несовпадениям картирование рида547
Сведение приблизительного сопоставления с паттерном к точному.....	.547
BLAST: Сравнение последовательности с базой данных549
Приблизительное сопоставление последовательностей с помощью	
преобразования Барроуза–Уилера550
Сопутствующие материалы552
Построение суффиксного дерева552
Решение задачи самой длинной общей подстроки.....	.555
Построение частичного суффиксного массива.....	.557
Эталонный геном человека558
Рекомбинации, вставки и делеции в геномах человека.....	.558
Алгоритм Ахо–Корасик559
Суффиксные массивы и суффиксные деревья560
Бинарный поиск.....	.565
Библиографические примечания566
Глава 10. Почему биологи до сих пор не разработали	
вакцину от ВИЧ?567
Классификация фенотипа ВИЧ.....	.568
Каким образом ВИЧ ускользает от иммунной системы человека?568
Ограничения метода выравнивания последовательностей570
Азартные игры с якудза572
Две монеты в рукаве у дилера573
Поиск CG-островов575
Скрытые марковские модели576
От подбрасывания монеты к скрытой марковской модели576
Диаграмма НММ577
Переформулировка задачи казино578
Задача декодирования581
Граф Виттерби.....	.581
Алгоритм Виттерби583
Насколько быстр алгоритм Виттерби?584
Поиск наиболее вероятного результата НММ586
Профильтрованные НММ для выравнивания последовательностей588

14 СОДЕРЖАНИЕ

Как НММ связаны с выравниванием последовательностей?	588
Создание профильной НММ	591
Вероятности перехода и эмиссии профильной НММ	594
Классификация белков с помощью профильных НММ	597
Выравнивание белков по профильной НММ	597
Возвращение псевдосчетов	598
Проблема с молчащими состояниями	601
Действительно ли профильные НММ так полезны?	606
Обучение параметров НММ	608
Определение параметров НММ, когда скрытый путь известен	608
Обучение Виттерби	609
Мягкие решения для определения параметров	611
Задача мягкого декодирования	611
Алгоритм «вперед–назад»	612
Обучение Баума–Уэлча	615
Многоликость НММ	617
Эпилог. Природа – мастер, а не изобретатель	618
Сопутствующие материалы	619
Эффект Красной Королевы	619
Гликозилирование	620
Метилирование ДНК	621
Условная вероятность	621
Библиографические примечания	622
Глава 11. Является ли <i>T. rex</i> всего лишь гигантской курицей?	623
Палеонтология встречается с информатикой	624
Какие белки присутствуют в этом образце?	625
Расшифровка идеального спектра	626
От идеального спектра к реальному	629
Секвенирование пептидов	632
Определение пептидов по спектрам	632
Где находятся суффиксные пептиды?	634
Алгоритм секвенирования пептидов	637
Идентификация пептидов	639
Задача идентификации пептидов	639
Идентификация пептидов в неизвестном протеоме <i>T. rex</i>	640
Поиск совпадений пептидов со спектром	640
Идентификация пептидов и теорема о бесконечных обезьянах	642
Частота ложных открытий	642
Статистическая значимость пептид–спектр–совпадений	645
Спектральные словари	647
Пептиды <i>T. rex</i> : постороннее загрязнение или древнее сокровище?	651
Загадка гемоглобина	651
Споры о ДНК динозавров	653
Эпилог. От немодифицированных к модифицированным пептидам. (Часть 1)	654
Посттрансляционные модификации	654
Поиск модификаций как задача выравнивания	655
Построение сетки Манхэттена для спектрального выравнивания	657
Эпилог. От немодифицированных к модифицированным пептидам (Часть 2)	660

Алгоритм спектрального выравнивания	660
Сопутствующие материалы	663
Предсказание генов	663
Поиск всех путей в графе.....	664
Задача антисимметричного пути	665
Преобразование спектров в спектральные векторы.....	666
Теорема о бесконечных обезьянах	670
Вероятностное пространство пептидов в словаре	671
Действительно ли динозавры являются предками птиц?	672
Библиографические примечания	673
Предметный указатель	674