

Стивен Л. Брантон, Дж. Натан Куц

Анализ данных в науке и технике

Data-Driven Science and Engineering

**Machine Learning, Dynamical Systems,
and Control**

Steven L. Brunton, J. Nathan Kutz



Содержание

От издательства	13
Об авторах	14
Предисловие.....	15
Общепотребительные методы оптимизации, уравнения, символы и акронимы	20
Часть I. Понижение размерности и преобразования	31
Глава 1. Сингулярное разложение (SVD)	32
1.1. Общие сведения.....	33
Определение SVD	34
Вычисление SVD	35
Историческая справка	36
Использование в этой книге и предположения о подготовке читателей.....	37
1.2. Аппроксимация матриц	37
Усечение.....	38
Пример: сжатие изображения.....	38
1.3. Математические свойства и манипуляции	41
Интерпретация с привлечением доминирующих корреляций	41
Метод моментальных снимков.....	43
Геометрическая интерпретация	43
Инвариантность SVD относительно унитарных преобразований	45
Левые унитарные преобразования	46
Правые унитарные преобразования	46
1.4. Псевдообращение, метод наименьших квадратов и регрессия.....	47
Одномерная линейная регрессия	49
Полилинейная регрессия	51
Предостережение	53
1.5. Метод главных компонент (PCA)	53
Вычисление	54
Пример: данные с гауссовым шумом.....	55
Пример: данные о раке яичников	57
1.6. Пример: «собственные лица»	58
1.7. Отсечение и выравнивание.....	64
Оптимальный жесткий порог отсечения.....	64
Важность выравнивания данных.....	68

1.8. Рандомизированное сингулярное разложение	71
Рандомизированная линейная алгебра	71
Рандомизированный алгоритм SVD	72
Пример рандомизированного SVD	75
1.9. Тензорные разложения и N -мерные массивы данных	76
Рекомендуемая литература	81

Глава 2. Преобразование Фурье

и вейвлет-преобразование	82
2.1. Ряд Фурье и преобразование Фурье	83
Скалярные произведения функций и векторов	83
Ряд Фурье	84
Преобразование Фурье	89
2.2. Дискретное преобразование Фурье (ДПФ) и быстрое преобразование Фурье (БПФ)	92
Дискретное преобразование Фурье	93
Быстрое преобразование Фурье	95
Пример БПФ: фильтрация шума	96
Пример БПФ: спектральные производные	98
2.3. Преобразование дифференциальных уравнений в частных производных	100
Уравнение теплопроводности	101
Одностороннее волновое уравнение	103
Уравнение Бюргера	105
2.4. Преобразование Габора и спектрограмма	107
Дискретное преобразование Габора	108
Пример: сигнал с квадратичной частотной модуляцией	108
Пример: «Патетическая соната» Бетховена	110
Принцип неопределенности	112
2.5. Вейвлеты и многомасштабный анализ	113
Дискретное вейвлет-преобразование	115
2.6. Двумерные преобразования и обработка сигналов	116
Двумерное преобразование Фурье для изображений	116
Двумерное вейвлет-преобразование изображений	119
Рекомендуемая литература	122

Глава 3. Разреженность и сжатие измерений

3.1. Разреженность и сжатие	124
Пример: сжатие изображения	125
Почему сигналы допускают сжатие: просторность пространства изображений	127
3.2. Сжатое измерение	128
Заявление об отказе от ответственности	132
Другие формулировки	133
3.3. Примеры сжатых измерений	133
Норма ℓ_1 и разреженные решения недоопределенной системы	134

Восстановление звукового сигнала по разреженным измерениям	135
3.4. Геометрия сжатия	137
Свойство ограниченной изометрии (RIP)	139
Некогерентность и матрицы измерений.....	139
Плохие измерения	140
3.5. Разреженная регрессия.....	140
Отбрасывание выбросов и робастность	141
Отбор признаков и LASSO-регрессия	142
3.6. Разреженное представление.....	146
3.7. Робастный метод главных компонент (RPCA).....	151
3.8. Разреженное размещение датчиков.....	153
Разреженное размещение датчиков для реконструкции	154
Разреженная классификация	158
Рекомендуемая литература	159

Часть II. МАШИННОЕ ОБУЧЕНИЕ И АНАЛИЗ

ДАННЫХ

Глава 4. Регрессия и выбор модели.....

4.1. Классическая аппроксимация кривой	163
Методы наименьших квадратов	163
Линия наименьших квадратов.....	166
Линеаризация данных.....	167
4.2. Нелинейная регрессия и градиентный спуск.....	169
Градиентный спуск	170
Метод переменных направлений	175
4.3. Регрессия и уравнение $Ax = b$: переопределенные и недоопределенные системы.....	176
Переопределенные системы	176
Недоопределенные системы	180
4.4. Оптимизация как краеугольный камень регрессии	183
4.5. Парето-фронт и Lex Parsimoniae	188
Переобучение.....	190
4.6. Выбор модели: перекрестная проверка	191
<i>k</i> -групповая перекрестная проверка.....	195
Перекрестная проверка с контролем по <i>p</i> точкам	197
4.7. Выбор модели: информационный критерий	197
Информационные критерии: AIC и BIC	200
Вычисление AIC и BIC.....	201
Рекомендуемая литература	202

Глава 5. Кластеризация и классификация

5.1. Выделение признаков и добыча данных	204
5.2. Обучение с учителем и без учителя.....	210
5.3. Обучение без учителя: кластеризация методом <i>k</i> средних	214
5.4. Иерархическая кластеризация без учителя: дендрограмма.....	219

5.5. Смесовые модели и EM-алгоритм.....	223
5.6. Обучение с учителем и линейные дискриминанты.....	227
5.7. Метод опорных векторов (SVM)	233
Линейный SVM	233
Нелинейный SVM.....	235
Ядерные методы в сочетании с SVM.....	236
5.8. Решающие деревья и случайные леса	238
Случайные леса.....	243
5.9. 10 лучших алгоритмов по версии Data Mining 2008.....	244
Алгоритм k средних	244
EM-алгоритм (смесовые модели)	245
Метод опорных векторов (SVM).....	245
CART (Classification and Regression Tree – дерево классификации и регрессии)	245
Метод k ближайших соседей (kNN).....	246
Наивная байесовская классификация.....	246
AdaBoost (ансамблевое обучение с усилением)	246
C4.5 (ансамблевое обучение решающих деревьев).....	247
Алгоритм Apriori	247
PageRank	247
Рекомендуемая литература	248

Глава 6. Нейронные сети и глубокое обучение

6.1. Нейронные сети: однослойные сети	250
Однослойная сеть.....	252
6.2. Многослойные сети и функции активации	255
6.3. Алгоритм обратного распространения.....	260
6.4. Алгоритм стохастического градиентного спуска	264
6.5. Глубокие сверточные нейронные сети.....	267
Сверточные слои	268
Пулинговые слои	269
Полносвязные слои	269
Прореживание	270
6.6. Нейронные сети для динамических систем.....	272
6.7. Разнообразие нейронных сетей	277
Перцептрон	277
Сети прямого распространения (FF)	277
Рекуррентная нейронная сеть (RNN)	279
Автокодировщик (AE).....	279
Марковская цепь (MC)	280
Сеть Хопфилда (HN).....	280
Машина Больцмана (BM)	280
Ограниченная машина Больцмана (RBM)	281
Сеть глубокого доверия (DBN).....	281
Глубокая сверточная нейронная сеть (DCNN).....	281
Антисверточная сеть (DN).....	281
Глубокая сверточная сеть обратной графики (DCIGN).....	282

Порождающая состязательная сеть (GAN)	282
Машина неустойчивых состояний (LSM).....	282
Машина экстремального обучения (ELM)	283
Сеть с эхо-состояниями (ESN)	283
Глубокая остаточная сеть (DRN).....	283
Сеть Кохонена (KN)	284
Нейронная машина Тьюринга (NTM).....	284
Рекомендуемая литература	284

Часть III. ДИНАМИЧЕСКИЕ СИСТЕМЫ И УПРАВЛЕНИЕ.....

285

Глава 7. Динамические системы, управляемые данными.....

286

7.1. Обзор, мотивация и проблемы	287
Динамические системы	287
Цели и проблемы современной теории динамических систем	291
7.2. Разложение по динамическим модам (DMD)	294
Алгоритм DMD	295
Пример и код	300
Расширения, приложения и ограничения.....	300
7.3. Разреженная идентификация нелинейной динамики (SINDy)	308
Нахождение дифференциальных уравнений в частных производных.....	314
Обобщение SINDy на рациональные нелинейности	316
Применение информационного критерия для выбора модели	319
7.4. Теория оператора Купмана	320
Математическая теория оператора Купмана.....	320
Разложение по модам Купмана и конечные представления.....	324
Примеры погружений Купмана	326
Аналитическое разложение собственных функций в ряд	329
История и недавние достижения.....	331
7.5. Управляемый данными анализ Купмана	332
Расширенный DMD	332
Аппроксимация собственных функций Купмана на основе данных	334
Управляемый данными анализ Купмана и запаздывающие координаты	336
Нейронные сети для погружений Купмана.....	340
Рекомендуемая литература	342

Глава 8. Теория линейного управления

344

Типы управления	345
8.1. Управление с замкнутым контуром обратной связи.....	346
Примеры преимуществ управления с обратной связью.....	348
8.2. Линейные стационарные системы	351
Линеаризация нелинейной динамики	351

Неуправляемая линейная система	352
Управляемая линейная система	354
Системы с дискретным временем	355
Пример: обратный маятник	356
8.3. Управляемость и наблюдаемость	357
Управляемость	357
Наблюдаемость	359
Критерий управляемости РВН	360
Теорема Кэли–Гамильтона и достижимость	361
Грамианы и степень управляемости и наблюдаемости	362
Стабилизируемость и распознаваемость	364
8.4. Оптимальное управление полным состоянием: линейно-квадратичный регулятор (ЛКР)	364
Вывод уравнения Риккати оптимального управления	366
8.5. Оптимальное оценивание полного состояния: фильтр Калмана	369
8.6. Оптимальное управление с использованием датчиков: линейно-квадратичное гауссово управление (ЛКГ)	372
8.7. Практический пример: обратный маятник на тележке	374
Управление маятником на тележке с обратной связью	376
Оценка полного состояния системы маятник–тележка	379
Управление с обратной связью системой маятник–тележка с использованием датчиков	382
8.8. Робастное управление и методы анализа в частотной области	384
Методы в частотной области	384
Качество управления и передаточная функция контура: чувствительность и дополнительная чувствительность	389
Обращение динамики	392
Робастное управление	393
Рекомендуемая литература	396

Глава 9. Сбалансированные модели, пригодные

для управления	397
9.1. Упрощение модели и идентификация системы	397
9.2. Сбалансированное упрощение модели	399
Цель упрощения модели	399
Замена переменных в системах управления	401
Балансирующие преобразования	403
Сбалансирование усечения	407
Вычисление сбалансированных реализаций	408
Пример сбалансированного упрощения модели	413
9.3. Идентификация системы	415
Алгоритм реализации собственной системы	416
Идентификация наблюдателей с помощью фильтра Калмана	419
Комбинация ERA и OKID	423
Рекомендуемая литература	425

Глава 10. Управление на основе данных	426
10.1. Идентификация нелинейной системы для управления.....	427
DMD с управлением.....	428
Нелинейное управление с помощью оператора Купмана.....	430
SINDy с управлением.....	432
Пример управления на основе прогнозирующих моделей (MPC).....	432
10.2. Управление с машинным обучением.....	436
Обучение с подкреплением.....	438
Управление с итеративным обучением.....	439
Генетические алгоритмы.....	439
Генетическое программирование.....	441
Пример: применение генетического алгоритма для настройки ПИД-регулятора.....	443
10.3. Адаптивное управление с поиском экстремума.....	448
Простой пример управления с поиском экстремума.....	452
Пример управления с поиском экстремума в сложной ситуации.....	455
Приложения управления с поиском экстремума.....	456
Рекомендуемая литература.....	458
Часть IV. МОДЕЛИ ПОНИЖЕННОГО ПОРЯДКА	460
Глава 11. Модели пониженного порядка (ROM)	461
11.1. POD для дифференциальных уравнений в частных производных.....	462
Разложение по модам Фурье.....	465
Специальные функции и теория Штурма–Лиувилля.....	466
Понижение размерности.....	467
11.2. Элементы оптимального базиса: собственное ортогональное разложение.....	468
Проекция Галеркина на POD-моды.....	470
Пример: гармонический осциллятор.....	471
11.3. POD и динамика солитонов.....	475
Упрощение солитона ($N = 1$).....	477
Упрощение солитона ($N = 2$).....	479
11.4. POD в непрерывной формулировке.....	480
Квадратурные правила для R: правило трапеций.....	482
Квадратурные правила более высокого порядка.....	483
POD-моды и квадратурные формулы.....	485
11.5. POD с симметриями: повороты и сдвиги.....	486
Сдвиг: распространение волн.....	486
Поворот: спиральные волны.....	488
Рекомендуемая литература.....	492
Глава 12. Интерполяция для ROM	494
12.1. Неполное POD.....	494
Разреженные измерения и реконструкция.....	496

Моды гармонического осциллятора	497
12.2. Ошибка и сходимость неполного POD	501
Случайная выборка и сходимость	501
Неполные измерения и качество реконструкции.....	503
12.3. Неполные измерения: минимизация числа обусловленности	504
Замены числа обусловленности.....	510
12.4. Неполные измерения: максимальная дисперсия.....	512
12.5. POD и дискретный эмпирический метод интерполяции (DEIM) ...	517
POD и DEIM	518
DEIM	519
12.6. Реализация алгоритма DEIM.....	521
Алгоритм QDEIM	523
12.7. Машинное обучение ROM.....	524
Выбор POD-моды	525
Пример: обтекание цилиндра	527
Рекомендуемая литература	529
Глоссарий.....	531
Список литературы.....	538
Предметный указатель	539

От издательства

Отзывы и пожелания

Мы всегда рады отзывам наших читателей. Расскажите нам, что вы думаете об этой книге – что понравилось или, может быть, не понравилось. Отзывы важны для нас, чтобы выпускать книги, которые будут для вас максимально полезны.

Вы можете написать отзыв на нашем сайте www.dmkpress.com, зайдя на страницу книги и оставив комментарий в разделе «Отзывы и рецензии». Также можно послать письмо главному редактору по адресу dmkpress@gmail.com; при этом укажите название книги в теме письма.

Если вы являетесь экспертом в какой-либо области и заинтересованы в написании новой книги, заполните форму на нашем сайте по адресу http://dmkpress.com/authors/publish_book/ или напишите в издательство по адресу dmkpress@gmail.com.

Скачивание исходного кода примеров

Скачать файлы с дополнительной информацией для книг издательства «ДМК Пресс» можно на сайте www.dmkpress.com на странице с описанием соответствующей книги.

Список опечаток

Хотя мы приняли все возможные меры для того, чтобы обеспечить высокое качество наших текстов, ошибки все равно случаются. Если вы найдете ошибку в одной из наших книг, мы будем очень благодарны, если вы сообщите о ней главному редактору по адресу dmkpress@gmail.com. Сделав это, вы избавите других читателей от недопонимания и поможете нам улучшить последующие издания этой книги.

Нарушение авторских прав

Пиратство в интернете по-прежнему остается насущной проблемой. Издательства «ДМК Пресс» и Springer очень серьезно относятся к вопросам защиты авторских прав и лицензирования. Если вы столкнетесь в интернете с незаконной публикацией какой-либо из наших книг, пожалуйста, пришлите нам ссылку на интернет-ресурс, чтобы мы могли применить санкции.

Ссылку на подозрительные материалы можно прислать по адресу электронной почты dmkpress@gmail.com.

Мы высоко ценим любую помощь по защите наших авторов, благодаря которой мы можем предоставлять вам качественные материалы.

Об авторах

Стивен Л. Брантон – доцент факультета общего машиностроения в Вашингтонском университете. Также является внештатным сотрудником отделения прикладной математики и науки о данных Института eScience. Область его научных интересов охватывает применение науки о данных и машинного обучения к динамическим системам и управлению в области гидрогазодинамики, биомеханики, оптики, энергетических систем и производства. Является автором двух учебников, лауреатом премии армии и ВВС для молодых ученых, получил право преподавания в инженерном колледже Вашингтонского университета и удостоен премии для молодых преподавателей.

Дж. Натан Куц – профессор прикладной математики Вашингтонского университета, был деканом факультета до 2015 года. Также является внештатным профессором отделения электротехники и физики и старшим научным сотрудником отделения науки о данных в Институте eScience. Область научных интересов охватывает сложные системы и анализ данных, конкретно применение методов машинного обучения и динамических систем и управление в разнообразных приложениях. Автор двух учебников, лауреат премии Боинг за отличное преподавание прикладной математики, а также премии CAREER Национального научного фонда США.

Предисловие

Эта книга посвящена растущей области знаний на пересечении методов обработки больших данных, прикладной оптимизации и классических дисциплин инженерной математики и математической физики. Мы готовили данный материал на протяжении ряда лет, в основном для лекций, читаемых студентам старших курсов и аспирантам технических и физических факультетов.

Обычно такие студенты имеют подготовку в области линейной алгебры, дифференциальных уравнений и научных расчетов, а инженеры также знакомы с теорией управления и (или) дифференциальными уравнениями в частных производных. Однако в большинстве учебных программ научно-технических вузов методы обработки данных и (или) оптимизации освещаются слабо или не включены вовсе. С другой стороны, студенты, обучающиеся по специальностям «информатика» и «статистика», плохо знакомы с динамическими системами и теорией управления. Нашей целью было написать введение в прикладную науку о данных для обеих групп. Включенные в книгу методы отбирались по трем критериям: (1) релевантность, (2) простота и (3) общность. Мы стремились представить широкий круг тем от вводного материала до методов, реально применяемых в исследованиях.

Открытие на основе анализа данных революционизировало наши подходы к моделированию, прогнозированию поведения и управлению сложными системами. Самые насущные научно-технические задачи нашего времени не поддаются эмпирическим моделям и выводам, основанным на первопринципах. Все чаще исследователи обращаются к подходам на основе анализа данных при изучении широкого спектра сложных систем, как то: турбулентность, науки о мозге, климатология, эпидемиология, финансы, робототехника, автономные системы. Такие системы обычно являются нелинейными, динамическими, многомасштабными в пространстве и во времени, многомерными и имеют доминирующие паттерны, которые необходимо охарактеризовать и смоделировать, чтобы в конечном итоге обеспечить сбор данных, прогнозирование, оценку и управление. Благодаря современным математическим методам вкупе с невиданной ранее доступностью данных и располагаемыми вычислительными ресурсами мы теперь можем подступить к неприступным до недавнего времени проблемам. Упомянем лишь малую толику новых методов: надежное восстановление изображения по разреженным и зашумленным измерениям случайных пикселей, управление турбулентностью с помощью машинного обучения, оптимальное размещение датчиков и приводов, идентификация допускающих интерпретацию нелинейных динамических систем на основе одних лишь данных и модели пониженного порядка, позволяющие ускорить изучение и оптимизацию систем со сложной многомасштабной физикой.

Движущим началом современной науки о данных является доступность больших и постоянно увеличивающихся объемов данных вследствие заме-

чательных инноваций в области разработки дешевых датчиков, возросших на порядки вычислительных мощностей и практически неограниченной емкости устройств хранения и скорости передачи. Такое изобилие данных открывает перед учеными и инженерами во всех областях новые возможности для изобретений на основе анализа данных; часто в этой связи говорят о четвертой парадигме научного открытия [245]. Эта четвертая парадигма представляет собой естественную кульминацию первых трех: эмпирического эксперимента, аналитического вывода и вычислительного исследования. Интеграция всех трех методик создает новаторскую платформу для новых открытий на основе данных. Этот процесс научного открытия не нов и по сути дела повторяет усилия титанов научной революции: Иоганна Кеплера (1571–1630) и сэра Исаака Ньютона (1642–1727). Оба сыграли ключевую роль в разработке теоретических принципов небесной механики на базе сочетания эмпирических подходов, основанных на анализе данных, и аналитических вычислений. Наука о данных не заменяет математическую физику и технику, но дополняет ее с учетом достижений XXI века, что больше напоминает возрождение, нежели революцию.

Наука о данных сама по себе не нова, ее предложил больше 50 лет назад Джон Тьюки, предвидевший появление науки, в центре внимания которой будет обучение на данных, или *анализ данных* [152]. С тех пор в науке о данных преобладают два разных подхода [78]. Сообщество *машинного обучения* состоит в основном из специалистов по информатике и интересуется в первую очередь разработкой быстрых, масштабируемых и качественных алгоритмов прогнозирования. Сообщество же *статистического обучения*, которое вовсе необязательно во всем противопоставлять первому, больше сосредоточено на факультетах математической статистики и занимается выводом интерпретируемых моделей. Обе методологии могут похвастаться значительными успехами и закладывают математические и вычислительные основания методов науки о данных. Целью ученых и инженеров должно стать использование этих методов для выведения из результатов наблюдений и обсчета моделей (чаще нелинейных), которые правильно улавливают динамику системы и количественно и качественно обобщаются на ненаблюдавшиеся области фазового, параметрического или прикладного пространства. А в этой книге нашей целью будет применение статистических методов и методов машинного обучения к решению технических задач.

РАССМАТРИВАЕМЫЕ ВОПРОСЫ

В книге обсуждается целый ряд ключевых тем. Во-первых, во многих сложных системах присутствуют доминирующие *низкоразмерные паттерны* данных, несмотря на быстрое увеличение разрешающей способности измерений и вычислений. Базовая структура открывает возможность эффективного размещения датчиков и компактного представления для моделирования и управления. Выделение паттернов тесно связано со второй темой: отысканием *преобразований координат*, позволяющих упростить систему.

Действительно, богатая история математической физики вращается вокруг преобразований координат (например, спектральные разложения, преобразование Фурье, обобщенные функции и т. д.), хотя эти методы в большинстве своем были ограничены простой идеализированной геометрией и линейной динамикой. Умение выводить преобразования *на основе данных* открывает возможность обобщить их на новые задачи с более сложной геометрией и граничными условиями. На протяжении всей книги мы будем интересоваться *динамическими системами и управлением*, т. е. применением методов, основанных на анализе данных, к моделированию и управлению систем, изменяющихся во времени. Красной нитью проходит тема *управляемой данными прикладной оптимизации*, поскольку едва ли не каждый рассматриваемый вопрос так или иначе связан с оптимизацией (например, нахождение *оптимальных* низкоразмерных паттернов, *оптимальное* расположение датчиков, *оптимизация* в машинном обучении, *оптимальное* управление и т. д.). И еще одна, даже более фундаментальная тема – большинство данных организовано в массивы для анализа, а широкое развитие численных инструментов линейной алгебры, начиная с 1960-х годов, лежит в основе математических методов матричных разложений и стратегий решения, встречающихся в этой книге.

БЛАГОДАРНОСТИ

Мы в долгу перед многими талантливыми студентами, сотрудниками и коллегами, которые делились ценными замечаниями и предложениями и оказывали нам поддержку. Особенно мы благодарны Джошуа Проктору (Joshua Proctor), который стоял у истоков этой книги и помогал при планировании ее структуры и организации. Мы также извлекли много полезного из бесед с Бингом Брантоном (Bing Brunton), Игорем Мезичем (Igor Mezić), Берндом Ноаком (Bernd Noack) и Сэмом Тайрой (Sam Taira). Эта книга была бы невозможна без помощи наших сотрудников и коллег, исследования которых отражены в тексте.

На протяжении работы над книгой и чтения соответствующих курсов мы получали чрезвычайно ценные отзывы и комментарии от наших замечательных студентов и постдоков: Трэвиса Эшкама (Travis Askham), Майкла Ау-Юнга (Michael Au-Yeung), Цзе Бая (Zhe Bai), Идо Брайта (Ido Bright), Кэтлин Чемпион (Kathleen Champion), Эмили Кларк (Emily Clark), Чарльза Делаханта (Charles Delahunt), Даниэля Дылевски (Daniel Dylewski), Бена Эрричсона (Ben Erichson), Чарли Фислера (Charlie Fiesler), Синь Фу (Xing Fu), Чена Гонга (Chen Gong), Тарена Гормана (Taren Gorman), Джейкоба Гросека (Jacob Grosek), Сета Хирша (Seth Hirsh), Микала Джонсона (Mikala Johnson), Юрики Кайзер (Eurika Kaiser), Мейсона Кема (Mason Kamb), Джеймса Кьюнерта (James Kunert), Бетани Луш (Bethany Lusch), Педро Майа (Pedro Maia), Критики Манохара (Krithika Manohar), Найла Мэнгана (Niall Mangan), Арианы Мендибль (Ariana Mendible), Томаса Морена (Thomas Mohren), Меган Моррисон (Megan Morrison), Маркуса Куэйда (Markus Quade), Сэма Руди (Sam Rudy), Сюзанны Саргсян

(Susanna Sargsyan), Изабель Шерл (Isabel Scherl), Эли Шлизермана (Eli Shlizerman), Джорджа Степанянца (George Stepaniants), Бена Строма (Ben Strom), Чань Суна (Chang Sun), Роя Тэйлора (Roy Taylor), Мегханы Велагар (Meghana Velagar), Джейка Вехолта (Jake Weholt) и Мэтта Уильямса (Matt Williams). Наши студенты подвигли нас на написание этой книги, благодаря им мы каждый день приходим на работу с радостью и волнением.

Мы также благодарны руководителю издательской группы в Cambridge University Press Лоран Коулз (Lauren Cowles), на которую могли положиться на протяжении всего процесса работы.

ОНЛАЙНОВЫЕ МАТЕРИАЛЫ

Мы с самого начала предполагали, что к книге будут прилагаться обширные дополнительные материалы: код, данные, видео, домашние задания и рекомендуемые способы построения курса. Все эти материалы можно найти на сайте databookuw.com.

Код в сети полнее, чем в книге, в частности включен код генерации рисунков, пригодных для публикации. Визуализация данных была поставлена на первое место среди методов науки о данных в опросе «Состояние науки о данных и машинного обучения», проведенном на Kaggle 2017. Поэтому мы настоятельно рекомендуем читателям скачать код с сайта и в полной мере воспользоваться командами построения графиков и диаграмм.

Мы также записали и разместили на YouTube лекции по большинству тем, включенных в книгу. Есть дополнительные видео для студентов, желающих восполнить пробелы в подготовке по научным расчетам и основам прикладной математики. Этот текст задуман одновременно как справочное пособие и источник материалов к нескольким курсам, рассчитанным на студентов разного уровня. Большинство глав самостоятельны, на их основе можно разработать *курсы молодого бойца*, рассчитанные примерно на 10 часов каждый.

КАК ПОЛЬЗОВАТЬСЯ ЭТОЙ КНИГОЙ

Книга рассчитана на начинающих аспирантов и продвинутых студентов старших курсов научных и технических факультетов. Поэтому методы машинного обучения излагаются с азов, но при этом мы предполагаем, что студенты умеют моделировать физические системы с помощью дифференциальных уравнений и решать их с помощью таких программ, как **ode45**. Рассматриваются как начальные вопросы, так и актуальные исследовательские методы. Наша цель – представить цельный взгляд и математический инструментарий для решения научно-технических задач. Но книга может быть также полезна студентам, изучающим информатику и статистику, которые зачастую мало знают о динамических системах и теории управления. На основе представленного материала можно разработать несколько курсов,

программы некоторых из них имеются на сайте книги и включают домашние задания, наборы данных и код.

Прежде всего мы хотели, чтобы книга была интересной, чтобы она вдохновляла, открывала глаза и вооружала знаниями молодых ученых и инженеров. Мы пытались по возможности не слишком усложнять, не жертвуя при этом глубиной и широтой охвата, без которых не может быть никакой исследовательской работы. Многие главы можно было бы развернуть в целые книги, и такие книги есть. Однако мы также стремились к полноте в той мере, в какой этого можно ожидать от книги, посвященной столь обширной и быстро развивающейся области. Мы надеемся, что книга придется вам по вкусу, что вы овладеете всеми описанными в ней методами и измените мир с помощью прикладной науки о данных!

Общеупотребительные методы оптимизации, уравнения, символы и акронимы

НАИБОЛЕЕ РАСПРОСТРАНЕННЫЕ СТРАТЕГИИ ОПТИМИЗАЦИИ

Метод наименьших квадратов (обсуждается в главах 1 и 4) минимизирует сумму квадратов разностей (ошибок) между фактическими данными и предсказаниями модели. В случае линейного метода наименьших квадратов, когда данные аппроксимируются линейной функцией, имеется решение в замкнутой форме, которое можно найти, приравняв к нулю производную ошибки по каждому неизвестному. Этот подход широко используется в технике и прикладных науках для аппроксимации полиномиальными функциями. Применение нелинейного метода наименьших квадратов обычно требует итеративного уточнения путем аппроксимации нелинейного решения линейным на каждой итерации.

Градиентный спуск (обсуждается в главах 4 и 6) – основной метод выпуклой оптимизации в многомерных системах. Для минимизации ошибки вычисляется градиент аппроксимирующей функции. Решение обновляется итеративно путем *спуска с горы* в пространстве решений. Одномерным вариантом градиентного спуска является метод Ньютона–Рафсона. В многомерном пространстве метод часто находит только локальный минимум. Важнейшими инновациями в приложениях больших данных являются стохастический градиентный спуск и алгоритм обратного распространения, благодаря чему оптимизация сводится к самому вычислению градиента.

Чередующийся градиентный спуск (Alternating Descent Method – ADM) (обсуждается в главе 4) позволяет избежать вычисления градиента за счет того, что на каждом шаге производится оптимизация по одной неизвестной. Таким образом, все неизвестные переменные считаются постоянными, за исключением одной, по которой производится линейный поиск (выпуклая оптимизация). Эта переменная обновляется, после чего фиксируется, и то же самое повторяется для другой переменной. На одном шаге итерации пере-

бираются все неизвестные, а сами итерации продолжают до тех пор, пока не будет достигнута желаемая точность.

Расширенный метод Лагранжа (Augmented Lagrange Method – ALM) (обсуждается в главах 3 и 8) – класс алгоритмов для решения задач условной оптимизации. Они похожи на методы штрафования тем, что заменяют задачу оптимизации с ограничениями последовательностью задач без ограничений и прибавляют к целевой функции штрафной член, который играет роль множителя Лагранжа. Расширенный метод Лагранжа – не то же самое, что метод множителей Лагранжа.

Линейное программирование и симплекс-метод – безотказные алгоритмы выпуклой оптимизации. В линейном программировании целевая функция линейно зависит от неизвестных, а ограничениями являются линейные равенства и неравенства. Вычислив область допустимых решений – выпуклый политоп, – алгоритм линейного программирования находит в полиэдре точку, в которой функция принимает наименьшее (или наибольшее) значение, если таковая существует. Симплекс-метод – это конкретная итеративная процедура линейного программирования, которая по заданному опорному допустимому решению пытается найти другое опорное решение, для которого целевая функция принимает меньшее значение, и тем самым производит оптимизацию.

НАИБОЛЕЕ УПОТРЕБИТЕЛЬНЫЕ УРАВНЕНИЯ И СИМВОЛЫ

Линейная алгебра

Линейная система уравнений

$$\mathbf{Ax} = \mathbf{b}. \quad (0.1)$$

Матрица $\mathbf{A} \in \mathbb{R}^{p \times n}$ и вектор $\mathbf{b} \in \mathbb{R}^p$ обычно известны, а вектор $\mathbf{x} \in \mathbb{R}^n$ неизвестен.

Уравнение для собственных значений

$$\mathbf{AT} = \mathbf{T}\mathbf{\Lambda}. \quad (0.2)$$

Столбец ξ_k матрицы \mathbf{T} является собственным вектором матрицы $\mathbf{A} \in \mathbb{C}^{n \times n}$, соответствующим собственному значению λ_k : $\mathbf{A}\xi_k = \lambda_k \xi_k$. Матрица $\mathbf{\Lambda}$ – диагональная матрица, содержащая эти собственные значения, в простейшем случае все n собственных значений различны.

Замена координат

$$\mathbf{x} = \mathbf{\Psi}\mathbf{a}. \quad (0.3)$$

Вектор $\mathbf{x} \in \mathbb{R}^n$ можно записать как $\mathbf{a} \in \mathbb{R}^n$ в системе координат, определенной столбцами матрицы $\mathbf{\Psi} \in \mathbb{R}^{n \times n}$.

Уравнение измерений

$$\mathbf{y} = \mathbf{C}\mathbf{x}. \tag{0.4}$$

Вектор $\mathbf{y} \in \mathbb{R}^p$ является измерением состояния $\mathbf{x} \in \mathbb{R}^n$ в результате применения матрицы измерений $\mathbf{C} \in \mathbb{R}^{p \times n}$.

Сингулярное разложение

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* \approx \tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}\tilde{\mathbf{V}}^*. \tag{0.5}$$

Матрицу $\mathbf{X} \in \mathbb{C}^{n \times m}$ можно разложить в произведение трех матриц $\mathbf{U} \in \mathbb{C}^{n \times n}$, $\mathbf{\Sigma} \in \mathbb{C}^{n \times m}$ и $\mathbf{V} \in \mathbb{C}^{m \times m}$. Матрицы \mathbf{U} и \mathbf{V} унитарные, т. е. $\mathbf{U}\mathbf{U}^* = \mathbf{U}^*\mathbf{U} = \mathbf{I}^{n \times n}$ и $\mathbf{V}\mathbf{V}^* = \mathbf{V}^*\mathbf{V} = \mathbf{I}^{m \times m}$, где $*$ обозначает операцию комплексного сопряжения и транспонирования. Столбцы \mathbf{U} (соответственно \mathbf{V}) ортогональны и называются левыми (соответственно правыми) *сингулярными векторами*. На главной диагонали диагональной матрицы $\mathbf{\Sigma}$ находятся убывающие неотрицательные элементы, называемые *сингулярными значениями*.

Часто \mathbf{X} аппроксимируется матрицей низкого ранга $\tilde{\mathbf{X}} = \tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}\tilde{\mathbf{V}}^*$, где $\tilde{\mathbf{U}}$ и $\tilde{\mathbf{V}}$ содержат первые $r \ll n$ столбцов \mathbf{U} и \mathbf{V} соответственно, а $\tilde{\mathbf{\Sigma}}$ – левый верхний блок $\mathbf{\Sigma}$ размера $r \times r$. В контексте пространственных мод, моделей пониженного порядка и размещения датчиков матрица $\tilde{\mathbf{U}}$ часто обозначается буквой $\mathbf{\Psi}$.

Регрессия и оптимизация

Оптимизация переопределенных и недоопределенных линейных систем

$$\operatorname{argmin}_{\mathbf{x}} (\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2 + \lambda g(\mathbf{x})) \text{ или} \tag{0.6a}$$

$$\operatorname{argmin}_{\mathbf{x}} g(\mathbf{x}) \text{ при условии } \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2 \leq \varepsilon. \tag{0.6b}$$

Здесь $g(\mathbf{x})$ – штраф регрессии (со штрафным параметром λ для переопределенных систем). Для переопределенных и недоопределенных систем линейных уравнений $\mathbf{A}\mathbf{x} = \mathbf{b}$, когда решений либо не существует, либо бесконечно много, для нахождения решения нужно задать ограничение или штраф; эта процедура называется *регуляризацией*.

Оптимизация переопределенных и недоопределенных линейных систем

$$\operatorname{argmin}_{\mathbf{x}} (f(\mathbf{A}, \mathbf{x}, \mathbf{b}) + \lambda g(\mathbf{x})) \text{ или} \tag{0.7a}$$

$$\operatorname{argmin}_{\mathbf{x}} g(\mathbf{x}) \text{ при условии } f(\mathbf{A}, \mathbf{x}, \mathbf{b}) \leq \varepsilon. \tag{0.7b}$$

Это обобщение линейной системы на нелинейную систему $f(\cdot)$ с регуляризацией $g(\cdot)$. Такие переопределенные и недоопределенные системы часто решаются методами градиентного спуска.

Композиционная оптимизация для нейронных сетей

$$\operatorname{argmin}_{\mathbf{A}_j} (f_M(\mathbf{A}_M, \dots, f_2(\mathbf{A}_2, (f_1(\mathbf{A}_1, \mathbf{x})) \dots)) + \lambda g(\mathbf{A}_j)). \tag{0.8}$$

Здесь \mathbf{A}_k – матрицы весов связей между k -м и $(k + 1)$ -м слоями нейронной сети. Обычно это сильно недоопределенная система, которая регуляризи-

руется прибавлением $g(\mathbf{A}_j)$. Композиция и регуляризация весьма важны как для порождения выразительных представлений данных, так и для предотвращения переобучения.

Динамические системы и системы пониженного порядка

*Нелинейное обыкновенное дифференциальное уравнение
(динамическая система)*

$$\frac{d}{dt} \mathbf{x}(t) = \mathbf{f}(\mathbf{x}(t), t; \boldsymbol{\beta}). \quad (0.9)$$

Вектор $\mathbf{x}(t) \in \mathbb{R}^n$ описывает состояние системы, изменяющейся во времени t , $\boldsymbol{\beta}$ – вектор параметров, а \mathbf{f} – векторное поле. В общем случае \mathbf{f} – липшицева функция, что гарантирует существование и единственность решения.

Система с линейной зависимостью выхода от входа

$$\frac{d}{dt} \mathbf{x} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u} \quad (0.10a)$$

$$\mathbf{y} = \mathbf{C}\mathbf{x} + \mathbf{D}\mathbf{u}. \quad (0.10b)$$

Состояние системы представлено вектором $\mathbf{x} \in \mathbb{R}^n$, входы (приводы) – вектором $\mathbf{u} \in \mathbb{R}^q$, а выходы (датчики) – вектором $\mathbf{y} \in \mathbb{R}^p$. Матрицы \mathbf{A} , \mathbf{B} , \mathbf{C} , \mathbf{D} определяют динамику, управляющее воздействие, стратегию работы датчиков и эффект сквозного управления соответственно.

*Нелинейное отображение
(динамические системы с дискретным временем)*

$$\mathbf{x}_{k+1} = \mathbf{F}(\mathbf{x}_k). \quad (0.11)$$

Состояние системы на k -й итерации представлено вектором $\mathbf{x}_k \in \mathbb{R}^n$, а \mathbf{F} – потенциально нелинейное отображение. Часто это отображение описывает продвижение итераций во времени, т. е. $\mathbf{x}_k = \mathbf{x}(k\Delta t)$; в таком случае потоковое отображение обозначается $\mathbf{F}_{\Delta t}$.

Операторное уравнение Купмана (с дискретным временем)

$$\mathcal{K}g = g \circ \mathbf{F}_t \Rightarrow \mathcal{K}_t \varphi = \lambda \varphi. \quad (0.12)$$

Линейный оператор Купмана \mathcal{K}_t экстраполирует функции измерения состояния $g(\mathbf{x})$ с помощью потока \mathbf{F}_t . Собственные значения и собственные векторы \mathcal{K}_t обозначаются λ и $\varphi(\mathbf{x})$ соответственно. Оператор \mathcal{K}_t применяется к гильбертову пространству измерений.

Нелинейные дифференциальные уравнения в частных производных (УрЧП)

$$\mathbf{u}_t = \mathbf{N}(\mathbf{u}, \mathbf{u}_x, \mathbf{u}_{xx}, \dots, x, t; \boldsymbol{\beta}). \quad (0.13)$$

Состояние УрЧП описывается вектором \mathbf{u} , \mathbf{N} – нелинейный оператор эволюции, нижние индексы обозначают взятие частных производных, а x и t –

пространственная и временная переменные соответственно. УрЧП параметризуется значениями, собранными в векторе $\mathbf{\beta}$. Состояние УрЧП \mathbf{u} может быть непрерывной функцией $u(x, t)$, а может быть дискретизировано в нескольких точках пространства, $\mathbf{u}(t) = [u(x_1, t) \ u(x_2, t) \ \dots \ u(x_n, t)]^T \in \mathbb{R}^n$.

Разложение Галеркина

Непрерывное разложение Галеркина имеет вид:

$$u(x, t) \approx \sum_{k=1}^r a_k(t) \psi_k(x). \tag{0.14}$$

Функции $a_k(t)$ – коэффициенты, отражающие временную динамику, а $\psi_k(x)$ – пространственные моды. Для многомерного дискретизированного состояния разложение Галеркина принимает вид $\mathbf{u}(t) \approx \sum_{k=1}^r a_k(t) \boldsymbol{\Psi}_k$. Пространственные моды $\boldsymbol{\Psi}_k \in \mathbb{R}^n$ могут быть столбцами матрицы $\boldsymbol{\Psi} = \tilde{\mathbf{U}}$.

СПИСОК ОБОЗНАЧЕНИЙ

Размерности

- K количество ненулевых элементов K -разреженного вектора \mathbf{s}
- m количество снимков данных (т. е. столбцов \mathbf{X})
- n размерность состояния $\mathbf{x} \in \mathbb{R}^n$
- p размерность измерения, или выходной переменной $\mathbf{y} \in \mathbb{R}^p$
- q размерность выходной переменной $\mathbf{u} \in \mathbb{R}^q$
- r ранг усеченного сингулярного разложения или иной низкоранговой аппроксимации

Скаляры

- s частота в лапласовой области
- t время
- δ скорость обучения в методе градиентного спуска
- Δt временной шаг
- x пространственная переменная
- Δx пространственный шаг
- σ сингулярное значение
- λ собственное значение
- λ параметр разреженности при разреженной оптимизации (раздел 7.3)
- λ множитель Лагранжа (разделы 3.7, 8.4 и 11.4)
- τ порог

Векторы

- \mathbf{a} вектор амплитуд мод \mathbf{x} в базисе $\boldsymbol{\Psi}$, $\mathbf{a} \in \mathbb{R}^r$
- \mathbf{b} вектор измерений в линейной системе $\mathbf{Ax} = \mathbf{b}$
- \mathbf{b} вектор амплитуд мод в разложении по динамическим модам (раздел 7.2)
- \mathbf{Q} вектор, содержащий функцию потенциала в алгоритме PDE-FIND

- r** вектор невязок
s разреженный вектор $\mathbf{s} \in \mathbb{R}^n$
u регулируемая переменная (главы 8, 9, 10)
u вектор состояния УрЧП (главы 11 и 12)
w экзогенные входы
 \mathbf{w}_d возмущения системы
 \mathbf{w}_n шум измерений
 \mathbf{w}_r опорная траектория
x состояние системы $\mathbf{x} \in \mathbb{R}^n$
 \mathbf{x}_k снимок данных в момент t_k
 \mathbf{x}_j пример данных $j \in Z := \{1, 2, \dots, m\}$ (главы 5 и 6)
 $\tilde{\mathbf{x}}$ упрощенное состояние $\tilde{\mathbf{x}} \in \mathbb{R}^r$, т. е. $\mathbf{x} \approx \tilde{\mathbf{U}}\tilde{\mathbf{x}}$
 $\hat{\mathbf{x}}$ оценка состояния системы
y вектор измерений $\mathbf{y} \in \mathbb{R}^p$
 y_j метка данных $j \in Z := \{1, 2, \dots, m\}$ (главы 5 и 6)
 $\hat{\mathbf{y}}$ оценка измерения выхода
z преобразованное состояние $\mathbf{x} = \mathbf{Tz}$ (главы 8 и 9)
ε вектор ошибок
β бифуркационные параметры
ξ собственный вектор оператора Купмана (разделы 7.4 и 7.5)
ξ разреженный вектор коэффициентов (раздел 7.3)
φ мода в разложении по динамическим модам
ψ мода собственного ортогонального разложения (POD)
Υ вектор измерений УрЧП в алгоритме PDE-FIND

Матрицы

- A** матрица системы уравнений, или динамики
 $\hat{\mathbf{A}}$ редуцированная динамика в r -мерном подпространстве POD
 \mathbf{A}_x матричное представление линейной динамики с состоянием \mathbf{x}
 \mathbf{A}_y матричное представление линейной динамики с наблюдаемыми переменными \mathbf{y}
(A, B, C, V) матрицы системы с непрерывным пространством состояний
($\hat{\mathbf{A}}, \hat{\mathbf{B}}, \hat{\mathbf{C}}, \hat{\mathbf{V}}$) матрицы системы с дискретным пространством состояний
($\tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \tilde{\mathbf{C}}, \tilde{\mathbf{V}}$) матрицы пространства состояний системы в новых координатах $\mathbf{z} = \mathbf{T}^{-1}\mathbf{x}$
($\tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \tilde{\mathbf{C}}, \tilde{\mathbf{V}}$) матрицы пространства состояний упрощенной системы ранга r
B матрица входных данных с приводов
C матрица линейных измерений состояний
C матрица управляемости
F дискретное преобразование Фурье
G матричное представление линейной динамики состояний и входов $[\mathbf{x}^T \mathbf{u}^T]^T$
H матрица Ганкеля
H' матрица Ганкеля с временным сдвигом
I единичная матрица
K матричная форма оператора Купмана (глава 7)
K коэффициент усиления системы управления с замкнутым контуром (глава 8)

- K_f коэффициент усиления фильтра Калмана
 K_r коэффициент усиления линейно-квадратичного регулятора (ЛКР)
 L низкоранговая часть матрицы X (глава 3)
 O матрица наблюдаемости
 P унитарная матрица, применяемая к столбцам X
 Q весовая матрица стоимости отклонений от нулевого состояния в ЛКР (раздел 8.4)
 Q ортогональная матрица в QR-разложении
 R весовая матрица стоимости управляющих воздействий в ЛКР (раздел 8.4)
 R верхнетреугольная матрица в QR-разложении
 S разреженная часть матрицы X (глава 3)
 T матрица собственных векторов (глава 8)
 T замена координат (главы 8 и 9)
 U левые сингулярные векторы X , $U \in \mathbb{R}^{n \times n}$
 \hat{U} левые сингулярные векторы экономичного сингулярного разложения X , $\hat{U} \in \mathbb{R}^{n \times n}$
 \tilde{U} левые сингулярные векторы (POD-моды) усеченного сингулярного разложения X , $\tilde{U} \in \mathbb{R}^{n \times r}$
 V правые сингулярные векторы X , $V \in \mathbb{R}^{m \times m}$
 \tilde{V} правые сингулярные векторы (POD-моды) усеченного сингулярного разложения X , $\tilde{V} \in \mathbb{R}^{m \times r}$
 Σ матрица сингулярных значений X , $\Sigma \in \mathbb{R}^{n \times m}$
 $\hat{\Sigma}$ матрица сингулярных значений экономичного сингулярного разложения X , $\hat{\Sigma} \in \mathbb{R}^{m \times m}$
 $\tilde{\Sigma}$ матрица сингулярных значений усеченного сингулярного разложения X , $\tilde{\Sigma} \in \mathbb{R}^{r \times r}$
 W собственные векторы \tilde{A}
 W_c грамиан управляемости
 W_o грамиан наблюдаемости
 X матрица данных, $X \in \mathbb{R}^{n \times m}$
 X' матрица данных с временным сдвигом, $X' \in \mathbb{R}^{n \times m}$
 Y проекция матрицы X на ортогональный базис в рандомизированном сингулярном разложении (раздел 1.8)
 Y матрица данных наблюдаемых величин, $Y = g(X)$, $Y \in \mathbb{R}^{p \times m}$ (глава 7)
 Y' матрица данных наблюдаемых величин со сдвигом, $Y' = g(X')$, $Y' \in \mathbb{R}^{p \times m}$ (глава 7)
 Z эскиз матрицы для рандомизированного сингулярного разложения, $Z \in \mathbb{R}^{n \times r}$ (раздел 1.8)
 Θ матрица измерений, умноженная на разреживающий базис, $\Theta = C\Psi$ (глава 3)
 Θ матрица функций-кандидатов для SINDy (раздел 7.3)
 Γ матрица производных функций-кандидатов для SINDy (раздел 7.3)
 Ξ матрица коэффициентов функций-кандидатов для SINDy (раздел 7.3)
 Ξ матрица нелинейных снимков для DEIM (раздел 12.5)
 Λ диагональная матрица собственных значений
 Y матрица входных снимков, $Y \in \mathbb{R}^{q \times m}$
 Φ матрица DMD-мод, $\Phi \triangleq X'V'\Sigma^{-1}W$

Ψ ортонормированный базис (например, моды Фурье или POD-моды)

Тензоры

$(\mathcal{A}, \mathcal{B}, \mathcal{M})$ тензоры N -мерных массивов размера $I_1 \times I_2 \times \dots \times I_N$

Нормы

$\|\cdot\|_0$ псевдонорма ℓ_0 вектора \mathbf{x} : количество ненулевых элементов \mathbf{x}

$\|\cdot\|_1$ норма ℓ_1 вектора \mathbf{x} : $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$

$\|\cdot\|_2$ норма ℓ_2 вектора \mathbf{x} : $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n (x_i^2)}$

$\|\cdot\|_2$ норма ℓ_2 матрицы \mathbf{X} : $\|\mathbf{X}\|_2 = \max_x \frac{\|\mathbf{X}\mathbf{x}\|_2}{\|\mathbf{x}\|_2}$

$\|\cdot\|_F$ норма Фробениуса матрицы \mathbf{X} : $\|\mathbf{X}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^m |X_{ij}|^2}$

$\|\cdot\|_*$ ядерная норма матрицы \mathbf{X} : $\|\mathbf{X}\|_* = \text{trace}(\sqrt{\mathbf{X}^* \mathbf{X}}) = \sum_{i=1}^m \sigma_i$ (для $m \leq n$)

$\langle \cdot, \cdot \rangle$ скалярное произведение. Для функций $\langle f(x), g(x) \rangle = \int_{-\infty}^{\infty} f(x)g^*(x)dx$

$\langle \cdot, \cdot \rangle$ скалярное произведение. Для векторов $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^* \mathbf{v}$

Операторы, функции и отображения

\mathcal{F} преобразование Фурье

\mathbf{F} отображение динамической системы с дискретным временем

\mathbf{F}_t дискретное потоковое отображение динамической системы

\mathbf{f} динамическая система с непрерывным временем

G преобразование Габора

\mathbf{G} передаточная функция, отображающая входы на выходы (глава 8)

g скалярная функция измерения \mathbf{x}

\mathbf{g} скалярная функция измерения \mathbf{x}

J функция стоимости для регулирования

ℓ функция потерь в методе опорных векторов (глава 5)

\mathcal{K} оператор Купмана (с непрерывным временем)

\mathcal{K}_t оператор Купмана, ассоциированный с потоковым отображением

\mathcal{L} преобразование Лапласа

\mathbf{L} передаточная функция контура (глава 8)

\mathbf{L} линейное дифференциальное уравнение в частных производных (главы 11 и 12)

\mathbf{N} нелинейное дифференциальное уравнение в частных производных

\mathcal{O} порядок величины

\mathbf{S} функция чувствительности (глава 8)

\mathbf{T} дополнительная функция чувствительности (глава 8)

\mathcal{W} вейвлет-преобразование

μ несогласованность между матрицей измерений \mathbf{C} и базисом Ψ

κ число обусловленности

ϕ собственная функция Купмана

∇ оператор градиента

$*$ оператор свертки

Наиболее употребительные акронимы

БПФ	быстрое преобразование Фурье
ГО	глубокое обучение
ОДУ	обыкновенное дифференциальное уравнение
СНС	сверточная нейронная сеть
УрЧП	дифференциальное уравнение в частных производных
DMD	разложение по динамическим модам (dynamic mode decomposition)
РСА	метод главных компонент (principal components analysis)
POD	собственное ортогональное разложение (proper orthogonal decomposition)
ROM	модель пониженного порядка (reduced order model)
SVD	сингулярное разложение (singular value decomposition)

Прочие акронимы

ДПФ	дискретное преобразование Фурье
ИНС	искусственная нейронная сеть
ЛДА	линейный дискриминантный анализ
НУШ	нелинейное уравнение Шрёдингера
ОПФ	оконное преобразование Фурье (short time Fourier transform)
ПИД	пропорционально-интегрально-дифференцирующий регулятор
РНС	рекуррентная нейронная сеть
СГС	стохастический градиентный спуск
ADM	метод переменных направлений (alternating directions method)
AIC	информационный критерий Акаике (Akaike information criterion)
ALM	расширенный метод множителей Лагранжа (augmented Lagrange multiplier)
ARMA	авторегрессионное скользящее среднее (autoregressive moving average)
ARMAX	авторегрессионное скользящее среднее с экзогенным входом (autoregressive moving average with exogenous input)
BIC	байесовский информационный критерий (Bayesian information criterion)
BPOD	сбалансированное собственное ортогональное разложение (balanced proper orthogonal decomposition)
ССА	канонический корреляционный анализ (canonical correlation analysis)
CFD	вычислительная гидродинамика (computational fluid dynamics)
CoSaMP	согласованное преследование со сжатой выборкой (compressive sampling matching pursuit)
CWT	непрерывное вейвлет-преобразование (continuous wavelet transform)
DCT	дискретное косинусное преобразование (discrete cosine transform)
DEIM	дискретный эмпирический метод интерполяции (discrete empirical interpolation method)
DMDc	разложение по динамическим модам с управлением (dynamic mode decomposition with control)

DMDc	разложение по динамическим модам с управлением (dynamic mode decomposition with control)
DNS	прямое численное моделирование (direct numerical simulation)
DWT	дискретное вейвлет-преобразование
ECOG	электрокортикография (electrocorticography)
eDMD	расширенное DMD (extended DMD)
EIM	эмпирический метод интерполяции (empirical interpolation method)
EM	математическое ожидание-максимизация (expectation maximization)
EOF	эмпирические ортогональные функции (empirical orthogonal functions)
ERA	алгоритм реализации собственной системы (eigensystem realization algorithm)
ESC	управление с поиском экстремума (extremum-seeking control)
GMM	модель гауссовой смеси (Gaussian mixture model)
HAVOK	ганкелево альтернативное представление оператора Купмана (Hankel alternative view of Koopman)
ICA	метод независимых компонент (independent component analysis)
JL	Джонсона–Линденштраусса (Johnson–Lindenstrauss)
KL	Кульбака–Лейблера (Kullback–Leibler)
KLT	преобразование Карунена–Лоэва (Karhunen–Loève transform)
LAD	наименьшее абсолютное отклонение (least absolute deviations)
LASSO	оператор наименьшего абсолютного сжатия и выборки (least absolute shrinkage and selection operator)
LQE	линейно-квадратичная модель оценки (linear quadratic estimator)
LQG	линейно-квадратичный гауссов регулятор (linear quadratic Gaussian controller)
LQR	линейно-квадратичный регулятор
LTI	линейная стационарная система (linear time invariant system)
MIMO	с несколькими входами и несколькими выходами (multiple input, multiple output)
MLC	управление на основе машинного обучения (machine learning control)
MPE	оценка отсутствующей точки (missing point estimation)
mrDMD	многомасштабное разложение по динамическим модам (multi-resolution dynamic mode decomposition)
NARMAX	нелинейная авторегрессионная модель с экзогенными входами (nonlinear autoregressive model with exogenous inputs)
OKID	идентификация наблюдателей с помощью фильтра Калмана (observer Kalman filter identification)
PВН	критерий Попова–Белевича–Хаутуса (Popov–Belevitch–Hautus test)
PCP	преследование главных компонент (principal component pursuit)
PDE-FIND	функциональная идентификация нелинейной динамики с уравнениями в частных производных (partial differential equation functional identification of nonlinear dynamics)

PDF	функция распределения вероятностей (probability distribution function)
PIV	анемометрия по изображениям частиц (particle image velocimetry)
RIP	свойство ограниченной изометрии (restricted isometry property)
RKHS	гильбертово пространство с воспроизводящим ядром (reproducing kernel Hilbert space)
RPCA	робастный метод главных компонент (robust principal components analysis)
rSVD	рандомизированное SVD (randomized SVD)
SINDy	разреженная идентификация нелинейных систем (sparse identification of nonlinear dynamics)
SISO	с одним входом и одним выходом (single input, single output)
SRC	разреженное представление для классификации (sparse representation for classification)
SSA	анализ сингулярного спектра (singular spectrum analysis)
STLS	последовательный метод наименьших квадратов с порогом (sequential thresholded least-squares)
SVM	метод опорных векторов (support vector machine)
TICA	метод независимых компонент с задержкой (time-lagged independent component analysis)
VAC	вариационный подход к конформационной динамике (variational approach of conformation dynamics)

Часть I



**ПОНИЖЕНИЕ
РАЗМЕРНОСТИ
И ПРЕОБРАЗОВАНИЯ**

Глава 1

Сингулярное разложение (SVD)

Сингулярное разложение (SVD) – одно из самых важных разложений матриц, появившихся в компьютерную эру, оно лежит в основе почти всех методов обработки данных, рассматриваемых в этой книге. SVD – это численно устойчивое разложение матрицы, которое применимо для самых разных целей и при этом гарантированно существует. Мы будем использовать SVD для получения низкоранговых аппроксимаций матриц и для вычисления псевдообращения неквадратных матриц, т. е. нахождения решения системы уравнений вида $\mathbf{Ax} = \mathbf{b}$. Еще одно важное применение SVD находит как алгоритм, лежащий в основе метода главных компонент (principal component analysis – PCA), идея которого – выделение статистически значимых факторов из многомерных данных. Комбинация SVD+PCA применяется к решению широкого круга научно-технических задач.

В каком-то смысле SVD является обобщением быстрого преобразования Фурье (БПФ) – темы следующей главы. Учебники прикладной математики часто начинаются с описания БПФ, поскольку это основа многих классических аналитических и численных результатов. Однако БПФ работает в идеализированной постановке, тогда как SVD – более общая техника, основанная на анализе предъявленных данных. Поскольку эта книга посвящена данным, мы начнем с сингулярного разложения, которое можно рассматривать как базис, созданный *специально* под конкретные данные, в отличие от БПФ, предоставляющего *общий* базис.

Во многих предметных областях сложные системы порождают данные, которые естественно организуются в виде больших матриц или, более общо, массивов. Например, временной ряд, получающийся в результате эксперимента или имитационного моделирования, можно представить в виде матрицы, каждый столбец которой содержит все измерения в один момент времени. Если данные в каждый момент времени многомерные, как при построении трехмерной модели погоды высокой разрешающей способности, то их можно *разгладить*, преобразовав в длинный вектор-столбец, и собрать такие столбцы в большую матрицу. Аналогично значения пикселей полутонового изображения можно сохранить в виде матрицы или преобразовать в длинные векторы-столбцы, которые в совокупности образуют матрицу, представляющую кадры фильма. Примечательно, что данные, порождаемые

такими системами, обычно имеют низкий ранг. Это означает, что можно выделить несколько доминирующих паттернов, которые объясняют многомерные данные. SVD как раз и является численно устойчивым и эффективным методом обнаружения таких паттернов в данных.

1.1. ОБЩИЕ СВЕДЕНИЯ

Сейчас мы вкратце опишем SVD и разовьем интуицию, которая поможет продемонстрировать SVD на ряде мотивирующих примеров. SVD лежит в основе многих других методов, описанных в этой книге, в т. ч. методов классификации в главе 5, разложения по динамическим модам (DMD) в главе 7 и собственного ортогонального разложения (POD) в главе 11. Подробно математические свойства обсуждаются в последующих разделах.

Высокая размерность – самая главная трудность при обработке данных, порождаемых сложными системами. Такие наборы данных могут включать аудио, изображения и видео. Данные могут также генерироваться физической системой, например записи электрической активности мозга или измерения скорости течения жидкости в модели либо в эксперименте. Замечено, что во многих естественно возникающих системах в данных присутствуют преобладающие паттерны, характеризуемые аттрактором или многообразием низкой размерности [252, 251].

Рассмотрим, например, типичное изображение – оно содержит много измерений (пикселей) и, стало быть, является точкой в многомерном векторном пространстве. Однако большинство изображений отлично сжимаются, т. е. существенную информацию можно представить подпространством гораздо более низкой размерности. Сжимаемость изображений будет подробно обсуждаться на страницах этой книги. Сложные гидрогазодинамические системы, например атмосфера Земли или турбулентная спутная струя за кормой автомобиля, также дают убедительные примеры низкоразмерной структуры, скрывающейся в пространстве состояний высокой размерности. Хотя в точных моделях поведения жидкости или газа количество степеней свободы исчисляется миллионами и миллиардами, часто в потоке можно выделить доминирующие когерентные структуры, например периодические завихрения за кормой автомобиля или ураганы в атмосфере.

SVD предлагает систематический способ нахождения низкоразмерной аппроксимации данных высокой размерности в терминах доминирующих паттернов. Эту технику можно назвать *управляемыми данными*, поскольку для обнаружения паттернов нужны только данные, без привлечения экспертных знаний или интуиции. Метод SVD обладает численной устойчивостью и дает иерархическое представление данных в новой системе координат, определяемой доминирующими корреляциями внутри данных. Кроме того, гарантируется, что сингулярное разложение, в отличие от спектрального, существует для любой матрицы.

SVD имеет много полезных приложений, помимо понижения размерности данных. Этот метод используется для вычисления псевдообратной матрицы для неквадратных матриц и тем самым позволяет найти решения недоопре-

деленных или переопределенных матричных уравнений вида $\mathbf{Ax} = \mathbf{b}$. Также мы будем использовать SVD для очистки наборов данных от шумов. Не менее важно, что SVD позволяет охарактеризовать входную и выходную геометрию линейного отображения векторных пространств. Все эти приложения будут рассмотрены в данной главе и дадут нам возможность развить интуитивное понимание матриц и данных высокой размерности.

Определение SVD

В общем случае нас интересует анализ большого набора данных $\mathbf{X} \in \mathbb{C}^{n \times m}$:

$$\mathbf{X} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_m \\ | & | & & | \end{bmatrix}. \quad (1.1)$$

Столбцы $\mathbf{x}_k \in \mathbb{C}^n$ могут представлять измерения, полученные в процессе моделирования или эксперимента. Например, это могут быть изображения, вытянутые в векторы-столбцы длиной, равной количеству пикселей в изображении. Векторы-столбцы могут также представлять состояние нестационарной физической системы, например скорости течения жидкости в выбранных точках, множество измерений электрической активности мозга или состояние метеорологической модели с разрешением 1 км.

Индекс k – это метка, обозначающая k -й набор измерений. Во многих примерах из этой книги \mathbf{X} будет содержать *временные ряды* данных, а $\mathbf{x}_k = \mathbf{x}(k\Delta t)$. Зачастую *размерность состояния* n очень велика – порядка миллионов или миллиардов степеней свободы. Столбцы иногда называют (моментальными) *снимками*, а m обозначает количество снимков в \mathbf{X} . Для многих систем $n \gg m$, поэтому получается *высокая и тонкая* матрица, в отличие от *низкой и толстой* в случае, когда $n \ll m$.

SVD – это однозначно определенное разложение, существующее для любой комплексной матрицы $\mathbf{X} \in \mathbb{C}^{n \times m}$:

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*, \quad (1.2)$$

где $\mathbf{U} \in \mathbb{C}^{n \times n}$ и $\mathbf{V} \in \mathbb{C}^{m \times m}$ – *унитарные* матрицы¹ с ортонормированными столбцами, а $\mathbf{\Sigma} \in \mathbb{R}^{n \times m}$ – матрица, диагональные элементы которой вещественны и неотрицательны, а все остальные равны нулю. Здесь символ * обозначает операцию комплексного сопряжения и транспонирования². Ниже мы увидим, что унитарность \mathbf{U} и \mathbf{V} важна и часто используется.

Если $n \geq m$, то на диагонали матрицы $\mathbf{\Sigma}$ расположено не более m ненулевых элементов, и ее можно записать в виде $\mathbf{\Sigma} = \begin{bmatrix} \hat{\mathbf{\Sigma}} \\ \mathbf{0} \end{bmatrix}$. Поэтому \mathbf{X} можно *точно* представить, воспользовавшись *экономной* формой SVD:

¹ Квадратная матрица \mathbf{U} называется унитарной, если $\mathbf{U}\mathbf{U}^* = \mathbf{U}^*\mathbf{U} = \mathbf{I}$.

² Для вещественных матриц эта операция совпадает с обычным транспонированием: $\mathbf{X}^* = \mathbf{X}^T$.

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* = \left[\hat{\mathbf{U}} \quad \hat{\mathbf{U}}^\perp \right] \begin{bmatrix} \hat{\mathbf{\Sigma}} \\ \mathbf{0} \end{bmatrix} \mathbf{V}^* = \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\mathbf{V}^*. \quad (1.3)$$

Полное и экономное SVD показаны на рис. 1.1. На столбцы матрицы $\hat{\mathbf{U}}^\perp$ натянуто векторное пространство, дополнительное и ортогональное к натянутому на столбцы $\hat{\mathbf{U}}$. Столбцы \mathbf{U} называются *левыми сингулярными векторами* \mathbf{X} , а столбцы \mathbf{V} – *правыми сингулярными векторами*. Диагональные элементы $\hat{\mathbf{\Sigma}} \in \mathbb{C}^{m \times m}$ называются *сингулярными значениями*, они расположены в порядке убывания. Ранг \mathbf{X} равен количеству ненулевых сингулярных значений.

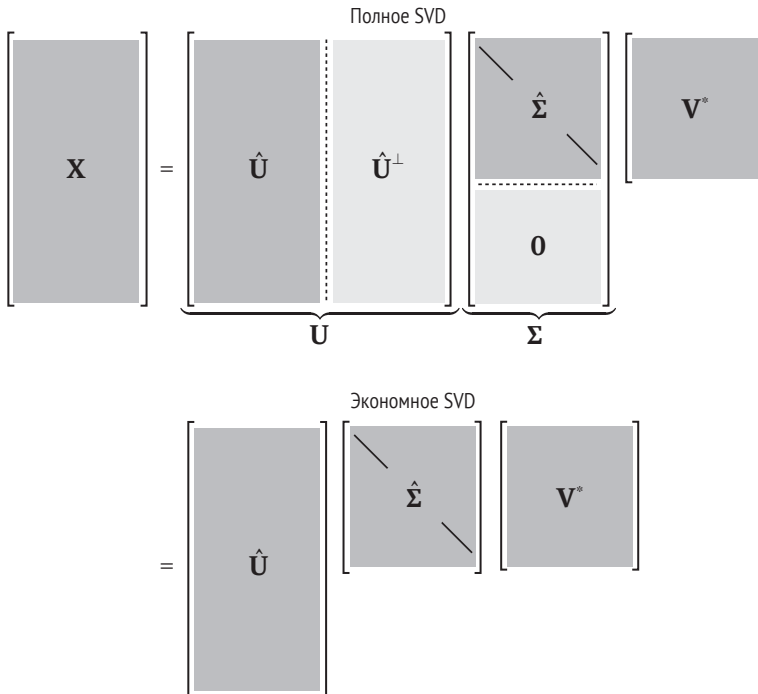


Рис. 1.1 ❖ Схема матриц в полном и экономном SVD

Вычисление SVD

SVD – краеугольный камень вычислительных методов в науке и технике, а численная реализация SVD важна и познавательна с математической точки зрения. Впрочем, большинство стандартных численных реализаций хорошо отработаны, к ним существует простой интерфейс из многих современных языков программирования, что позволяет нам абстрагироваться от деталей вычисления SVD. Как правило, мы будем просто использовать SVD как часть решения более крупной задачи и примем существование эффективных и устойчивых численных алгоритмов как данность. В последующих разделах мы продемон-

стрируем использование SVD на разных языках программирования, а также обсудим большинство стандартных стратегий вычисления и их ограничений. На тему вычисления SVD есть немало важных результатов [212, 106, 211, 292, 238]. Более полное обсуждение вычислительных аспектов можно найти в работе [214]. Для вычисления SVD очень больших матриц все чаще используют рандомизированные алгоритмы, этот вопрос обсуждается в разделе 1.8.

MATLAB. В MATLAB SVD вычисляется прямолинейно:

```
>>X = randn(5,3); % создать случайную матрицу 5×3
>>[U,S,V] = svd(X); % сингулярное разложение
```

Для неквадратных матриц X экономное SVD эффективнее:

```
>>[Uhat,Shat,V] = svd(X,'econ'); % экономное SVD
```

Python

```
>>> import numpy as np
>>> X = np.random.rand(5, 3) % create random data matrix
>>> U, S, V = np.linalg.svd(X,full_matrices=True) % full SVD
>>> Uhat, Shat, Vhat = np.linalg.svd(X, full_matrices=False)
% economy SVD
```

R

```
> X <- replicate(3, rnorm(5))
> s <- svd(X)
> U <- s$u
> S <- diag(s$d)
> V <- s$v
```

Mathematica

```
In:= X=RandomReal[{0,1},{5,3}]
In:= {U,S,V} = SingularValueDecomposition[X]
```

Другие языки

SVD реализовано и на других языках, например Fortran и C++. На самом деле большинство реализаций SVD основаны на библиотеке LAPACK (Linear Algebra Package) [13], написанной на Fortran. В LAPACK подпрограмма вычисления SVD называется **DGESVD**, а она уже обернута функциями C++ в библиотеках **Armadillo** и **Eigen**.

Историческая справка

SVD имеет давнюю и богатую историю, от ранних работ, в которых был заложен теоретический фундамент, до современных исследований по численной устойчивости и эффективности. Отличный исторический обзор приведен в работе Stewart [502], где описан общий контекст и многие важные детали. В этом обзоре много внимания уделено ранним теоретическим работам Бельтрами и Джордана (1873), Сильвестра (1889), Шмидта (1907) и Вейля (1912). Обсуждаются также более поздние работы, включая основополагаю-

щие труды по вычислительной стороне проблемы Голуба с сотрудниками [212, 211]. Кроме того, в современных учебниках [524, 17, 316] имеются прекрасно написанные главы, посвященные SVD.

Использование в этой книге и предположения о подготовке читателей

SVD – основа многих методов понижения размерности. К ним относятся метод главных компонент (PCA) в статистике [418, 256, 257], преобразование Карунена–Лоэва (KLT) [280, 340], эмпирические ортогональные функции (EOF) в изучении климата [344], собственное ортогональное разложение (POD) в гидродинамике [251] и канонический корреляционный анализ (CCA) [131]. Хотя многие из этих методов разрабатывались независимо для разных областей знания, различаются они только способами сбора и предварительной обработки данных. В работе Gerbrands [204] приведено превосходное обсуждение связи между SVD, KLT и PCA.

SVD широко используется для идентификации систем и в теории управления для получения моделей пониженного порядка, сбалансированных в том смысле, что состояния иерархически упорядочены в терминах доступности наблюдению и управляемости с помощью приводов [388].

В этой главе предполагается, что читатель знаком с линейной алгеброй и имеет некоторый опыт применения вычислительных методов. Если необходимо освежить познания, то существует много отличных книг по численной линейной алгебре, в т. ч. с обсуждением SVD, например [524, 17, 316].

1.2. АППРОКСИМАЦИЯ МАТРИЦ

Пожалуй, самым полезным и определяющим свойством SVD является тот факт, что оно дает *оптимальную* низкоранговую аппроксимацию матрицы \mathbf{X} . На самом деле SVD позволяет построить *иерархию* низкоранговых аппроксимаций, поскольку для получения аппроксимации ранга r нужно просто оставить первые r сингулярных значений и векторов, а остальные отбросить.

Шмидт обобщил SVD на пространства функций и доказал теорему, устанавливающую, что усеченное SVD является оптимальной низкоранговой аппроксимацией исходной матрицы \mathbf{X} [476]. Теорема Шмидта была заново открыта в работе Eckart and Young [170], поэтому иногда ее называют теоремой Эккарта–Янга.

Теорема 1 (Eckart–Young [170]). *Оптимальную в смысле наименьших квадратов аппроксимацию \mathbf{X} ранга r дает усеченное SVD $\tilde{\mathbf{X}}$ ранга r :*

$$\operatorname{argmin}_{\tilde{\mathbf{X}} \text{ такое, что } \operatorname{rank}(\tilde{\mathbf{X}})=r} \|\mathbf{X} - \tilde{\mathbf{X}}\|_F = \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^* \tag{1.4}$$

Здесь $\tilde{\mathbf{U}}$ и $\tilde{\mathbf{V}}$ обозначают матрицы, образованные первыми r столбцами \mathbf{U} и \mathbf{V} , а $\tilde{\Sigma}$ содержит левый верхний блок Σ размера $r \times r$. $\|\cdot\|_F$ – норма Фробениуса.

Согласно этой нотации базис усеченного SVD (и аппроксимирующая матрица $\tilde{\mathbf{X}}$) обозначается $\tilde{\mathbf{X}} = \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^*$. Поскольку Σ – диагональная матрица, SVD-аппроксимация ранга r может быть представлена суммой r матриц ранга 1:

$$\tilde{\mathbf{X}} = \sum_{k=1}^r \sigma_k \mathbf{u}_k \mathbf{v}_k^* = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^* + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^* + \dots + \sigma_r \mathbf{u}_r \mathbf{v}_r^*. \quad (1.5)$$

Это так называемое *диадическое* суммирование. Для любого ранга r не существует лучшей аппроксимации \mathbf{X} в смысле нормы ℓ_2 , чем усеченная SVD-аппроксимация $\tilde{\mathbf{X}}$. Таким образом, данные высокой размерности хорошо описываются несколькими доминирующими паттернами, определяемыми столбцами $\tilde{\mathbf{U}}$ и $\tilde{\mathbf{V}}$.

Это важное свойство SVD, к которому мы будем много раз возвращаться. Есть многочисленные примеры наборов данных, содержащих измерения высокой размерности, которые приводят к большой матрице \mathbf{X} . Однако в данных часто присутствуют доминирующие паттерны низкой размерности, и базис усеченного SVD $\tilde{\mathbf{U}}$ определяет преобразование координат из пространства измерений высокой размерности в пространство паттернов низкой размерности. В результате уменьшается размер и размерность больших наборов данных, а значит, открывается возможность для визуализации и анализа. Наконец, многие системы, рассматриваемые в этой книге, *динамические* (см. главу 7), а базис SVD дает иерархию мод, характеризующую наблюдаемый аттрактор, на который мы можем спроецировать динамическую систему низкой размерности для получения моделей пониженного порядка (см. главу 12).

Усечение

Усеченное SVD показано на рис. 1.2, где $\tilde{\mathbf{U}}$, $\tilde{\Sigma}$ и $\tilde{\mathbf{V}}$ обозначают усеченные матрицы. Если \mathbf{X} – матрица неполного ранга, то некоторые сингулярные значения в $\tilde{\Sigma}$ могут быть равны нулю, и тогда усеченное SVD остается точным. Однако если r меньше числа ненулевых сингулярных значений (т. е. ранга \mathbf{X}), то усеченное SVD является всего лишь аппроксимацией \mathbf{X} :

$$\mathbf{X} \approx \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^*. \quad (1.6)$$

Есть много способов выбора ранга усеченной матрицы r , они обсуждаются в разделе 1.7. Если мы требуем, чтобы усеченная матрица содержала все ненулевые сингулярные значения, то равенство $\mathbf{X} = \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^*$ точное.

Пример: сжатие изображения

Проиллюстрируем идею аппроксимации матриц на простом примере: сжатие изображения. Тема, которая красной нитью проходит через всю книгу, – наличие в больших наборах данных паттернов, благодаря чему возможны низкоранговые представления. Естественные изображения дают простой

и интуитивно понятный пример такой внутренне присущей *сжимаемости*. Полутонное изображение можно рассматривать как вещественную матрицу $\mathbf{X} \in \mathbb{R}^{n \times m}$, где n и m – числа пикселей по вертикали и по горизонтали соответственно¹. В зависимости от базиса представления (пространство пикселей, частотная область в смысле преобразования Фурье, преобразованные с помощью SVD координаты) изображение может иметь очень компактные аппроксимации.

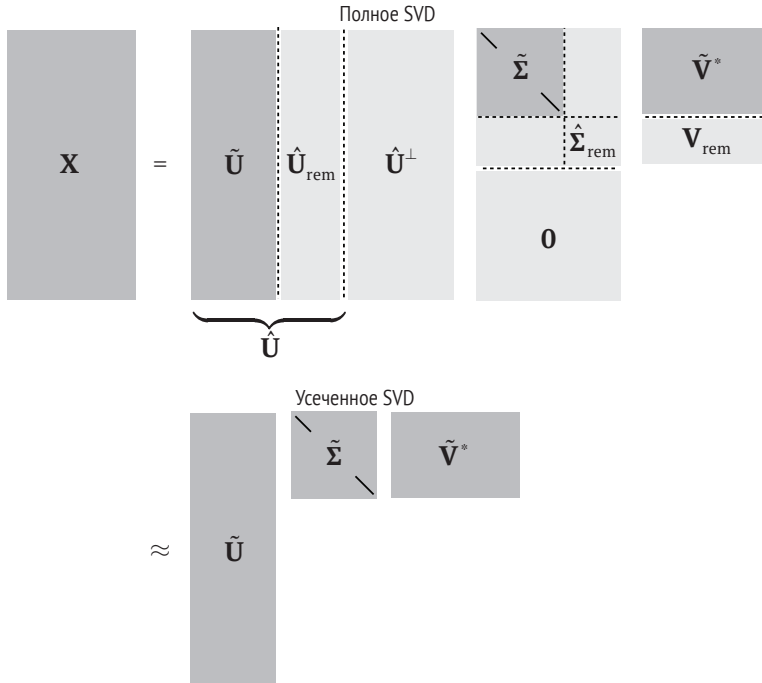


Рис. 1.2 ❖ Схема усеченного SVD.
Нижний индекс «rem» обозначает остаток $\hat{\mathbf{U}}$, $\hat{\mathbf{\Sigma}}$ или \mathbf{V} после усечения

Рассмотрим изображение собаки по кличке Мордекай на снегу (рис. 1.3). Его размер 2000×1500 пикселей. Мы можем вычислить SVD этого изображения и нанести на график сингулярные значения (рис. 1.4). На рис. 1.3 показаны приближенные матрицы $\tilde{\mathbf{X}}$ для разных значений r . При $r = 100$ реконструированное изображение вполне точное, а сингулярные значения отражают почти 80 % неоднородности изображения. Усечение SVD приводит к сжатию исходного изображения, поскольку в $\tilde{\mathbf{U}}$, $\tilde{\mathbf{\Sigma}}$ и $\tilde{\mathbf{V}}$ нужно хранить только первые 100 столбцов \mathbf{U} и \mathbf{V} плюс первые 100 диагональных элементов $\mathbf{\Sigma}$.

¹ Размер изображения часто задают, указывая сначала размер по горизонтали, а затем по вертикали, т. е. $\mathbf{X}^T \in \mathbb{R}^{m \times n}$, но мы будем придерживаться противоположного соглашения, совпадающего с общепринятым порядком обозначения размера матрицы.

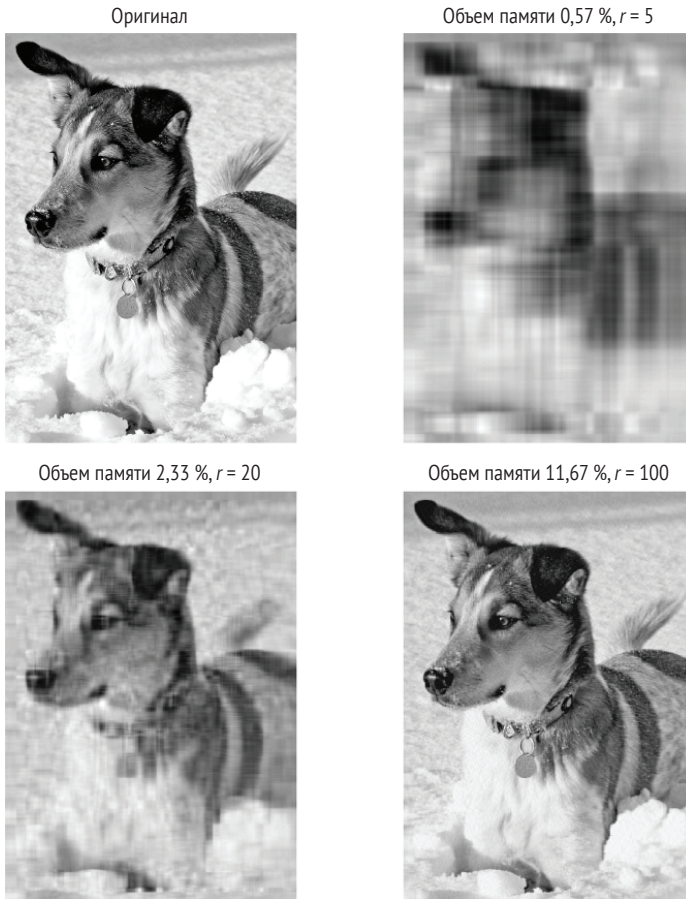


Рис. 1.3 ❖ Сжатие изображения собаки Мордекая на снегу путем усечения SVD с различными значениями ранга r . Разрешение исходного изображения 2000×1500

Сначала загрузим изображение:

```
A=imread('..\\DATA\\dog.jpg');
X=double(rgb2gray(A)); % преобразовать RGB в полутоновое, 256 бит->double.
nx = size(X,1); ny = size(X,2);
imagesc(X), axis off, colormap gray
```

и вычислим SVD:

```
[U,S,V] = svd(X);
```

Затем вычислим приближенную матрицу, используя усеченные SVD с различными рангами ($r = 5, 20, 100$):

```
for r=[5 20 100]; % Truncation value
    Xapprox = U(:,1:r)*S(1:r,1:r)*V(:,1:r)'; % Approx. image
```

```
figure, imagesc(Xapprox), axis off
title(['r=', num2str(r, '%d'), '']);
end
```

Наконец, построим графики сингулярных значений и суммарной энергии, изображенные на рис. 1.4.

```
subplot(1,2,1), semilogy(diag(S), 'k')
subplot(1,2,2), plot(cumsum(diag(S))/sum(diag(S)), 'k')
```

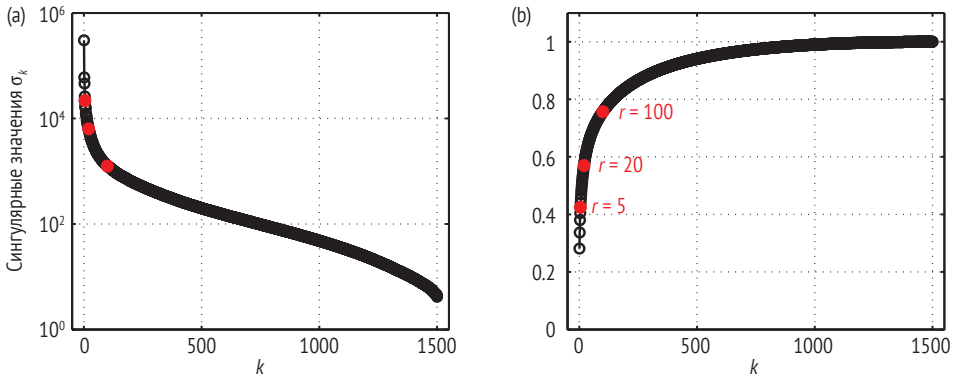


Рис. 1.4 ❖ (a) Сингулярные значения σ_k .
(b) Суммарная энергия в первых k модах

1.3. МАТЕМАТИЧЕСКИЕ СВОЙСТВА И МАНИПУЛЯЦИИ

Опишем важные математические свойства SVD, в т. ч. геометрические интерпретации унитарных матриц \mathbf{U} и \mathbf{V} , а также обсудим SVD в терминах доминирующих корреляций, присутствующих в данных \mathbf{X} . Связь между SVD и корреляциями данных мы объясним в разделе 1.5, посвященном методу главных компонент.

Интерпретация с привлечением доминирующих корреляций

Сингулярное разложение тесно связано с задачей о собственных значениях, в которой фигурируют корреляционные матрицы $\mathbf{X}\mathbf{X}^*$ и $\mathbf{X}^*\mathbf{X}$, показанные на рис. 1.5 для конкретного изображения, а на рис. 1.6 и 1.7 для матриц общего вида. Подставив (1.3) в $\mathbf{X}\mathbf{X}^*$ и $\mathbf{X}^*\mathbf{X}$, получаем:

$$\mathbf{X}\mathbf{X}^* = \mathbf{U} \begin{bmatrix} \hat{\Sigma} \\ \mathbf{0} \end{bmatrix} \mathbf{V}^* \mathbf{V} \begin{bmatrix} \hat{\Sigma} & \mathbf{0} \end{bmatrix} \mathbf{U}^* = \mathbf{U} \begin{bmatrix} \hat{\Sigma}^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{U}^*; \quad (1.7a)$$

$$\mathbf{X}^* \mathbf{X} = \mathbf{V} \begin{bmatrix} \hat{\Sigma} & \mathbf{0} \end{bmatrix} \mathbf{U}^* \mathbf{U} \begin{bmatrix} \hat{\Sigma} \\ \mathbf{0} \end{bmatrix} \mathbf{V}^* = \mathbf{V} \hat{\Sigma}^2 \mathbf{V}^*. \quad (1.7b)$$

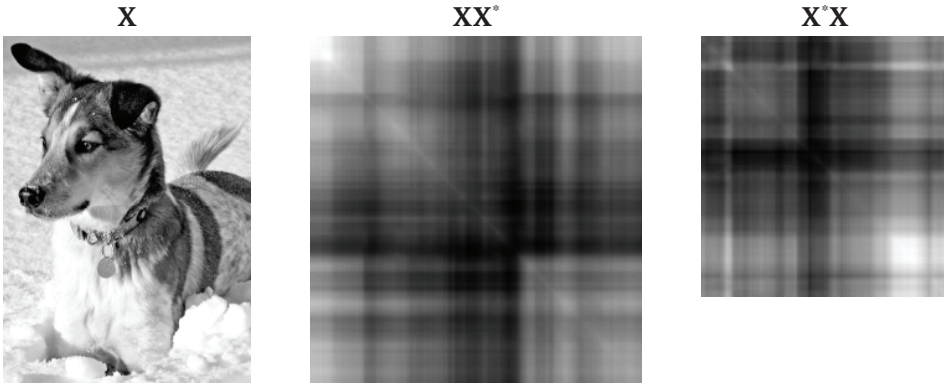


Рис. 1.5 ❖ Корреляционные матрицы $\mathbf{X}\mathbf{X}^*$ и $\mathbf{X}^*\mathbf{X}$ для матрицы \mathbf{X} , полученной из изображения собаки. Заметим, что обе корреляционные матрицы симметричны

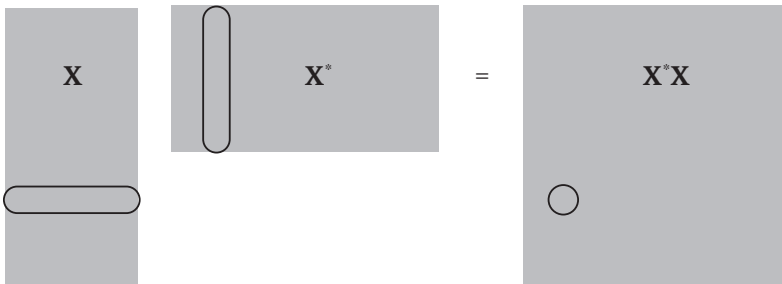


Рис. 1.6 ❖ Корреляционная матрица $\mathbf{X}\mathbf{X}^*$ получена взятием скалярных произведений строк \mathbf{X}

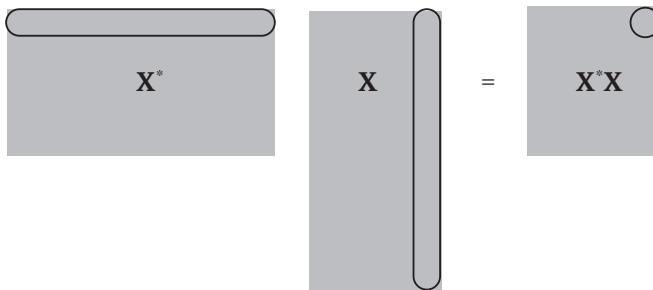


Рис. 1.7 ❖ Корреляционная матрица $\mathbf{X}^*\mathbf{X}$ получена взятием скалярных произведений столбцов \mathbf{X}

Учитывая, что \mathbf{U} и \mathbf{V} унитарны, \mathbf{U} , $\hat{\Sigma}$ и \mathbf{V} являются решениями следующих задач на нахождение собственных значений:

$$\mathbf{X}\mathbf{X}^*\mathbf{U} = \mathbf{U} \begin{bmatrix} \hat{\Sigma}^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}; \quad (1.8a)$$

$$\mathbf{X}^*\mathbf{X}\mathbf{V} = \mathbf{V}\hat{\Sigma}^2. \quad (1.8b)$$

Иными словами, любое ненулевое сингулярное значение \mathbf{X} является положительным квадратным корнем из какого-то собственного значения матриц $\mathbf{X}^*\mathbf{X}$ и $\mathbf{X}\mathbf{X}^*$, которые имеют одинаковые множества собственных значений. Отсюда следует, что если \mathbf{X} самосопряженная (т. е. $\mathbf{X} = \mathbf{X}^*$), то сингулярные значения \mathbf{X} равны абсолютным величинам собственных значений \mathbf{X} .

В результате мы получаем интуитивно понятную интерпретацию SVD: столбцы \mathbf{U} являются собственными векторами корреляционной матрицы $\mathbf{X}\mathbf{X}^*$, а столбцы \mathbf{V} – собственными векторами $\mathbf{X}^*\mathbf{X}$. Мы упорядочиваем сингулярные значения в порядке убывания абсолютной величины, поэтому порядок столбцов \mathbf{U} отражает, какую часть корреляции между столбцами \mathbf{X} они улавливают; аналогично \mathbf{V} улавливает корреляцию между строками \mathbf{X} .

Метод моментальных снимков

На практике часто невозможно построить матрицу $\mathbf{X}\mathbf{X}^*$, поскольку размерность состояния n слишком велика; что уж говорить о нахождении собственных значений. Если \mathbf{X} состоит из миллиона элементов, то число элементов $\mathbf{X}\mathbf{X}^*$ равно триллиону. В 1987 году Сирович (Sirovich) заметил, что можно не вычислять эту большую матрицу, а найти первые m столбцов \mathbf{U} с помощью метода, получившего название «метод (моментальных) снимков» [490].

Вместо того чтобы вычислять спектральное разложение $\mathbf{X}\mathbf{X}^*$ для получения левых сингулярных векторов \mathbf{U} , мы вычисляем только спектральное разложение гораздо меньшей и простой для работы матрицы $\mathbf{X}^*\mathbf{X}$. Затем из (1.8b) находим \mathbf{V} и $\hat{\Sigma}$. Если у $\hat{\Sigma}$ имеются нулевые сингулярные значения, то мы оставляем только ее часть $\tilde{\Sigma}$, соответствующую r ненулевым значениям, и соответствующие столбцы $\tilde{\mathbf{V}}$ матрицы \mathbf{V} . Зная эти матрицы, можно следующим образом аппроксимировать матрицу $\tilde{\mathbf{U}}$, состоящую из первых r столбцов \mathbf{U} :

$$\tilde{\mathbf{U}} = \mathbf{X}\tilde{\mathbf{V}}\tilde{\Sigma}^{-1}. \quad (1.9)$$

Геометрическая интерпретация

Столбцы матрицы \mathbf{U} образуют ортонормированный базис пространства столбцов \mathbf{X} . Аналогично столбцы \mathbf{V} образуют ортонормированный базис пространства строк \mathbf{X} . Если столбцы \mathbf{X} содержат пространственные измерения в разные моменты времени, то \mathbf{U} кодирует пространственные паттерны, а \mathbf{V} – временные паттерны.

Особенно полезным делает SVD тот факт, что \mathbf{U} и \mathbf{V} – унитарные матрицы, так что $\mathbf{U}\mathbf{U}^* = \mathbf{U}^*\mathbf{U} = \mathbf{I}_{n \times n}$ и $\mathbf{V}\mathbf{V}^* = \mathbf{V}^*\mathbf{V} = \mathbf{I}_{m \times m}$. Это означает, что для решения системы уравнений с матрицей \mathbf{U} или \mathbf{V} нужно просто умножить обе части на транспонированную матрицу. Сложность этой операции составляет $O(n^2)$, в отличие от традиционных методов обращения матрицы общего вида, имеющих сложность $O(n^3)$. Как отмечено в предыдущем разделе и в работе [57], SVD тесно связано со спектральными свойствами компактных самосопряженных операторов $\mathbf{X}\mathbf{X}^*$ и $\mathbf{X}^*\mathbf{X}$.

Сингулярное разложение \mathbf{X} можно геометрически интерпретировать, рассмотрев отображение гиперсферы $S^{n-1} \triangleq \{\mathbf{x} \mid \|\mathbf{x}\|_2 = 1\} \subset \mathbb{R}^n$ в эллипсоид, $\{\mathbf{y} \mid \mathbf{y} = \mathbf{X}\mathbf{x} \text{ для } \mathbf{x} \in S^{n-1}\} \subset \mathbb{R}^m$ посредством умножения на \mathbf{X} . Графически это показано на рис. 1.8 для сферы в \mathbb{R}^3 и отображения посредством умножения на \mathbf{X} с тремя ненулевыми сингулярными значениями. Поскольку умножение на матрицу – линейное отображение, достаточно знать, как оно ведет себя на единичной сфере, чтобы вычислить образ любого вектора.

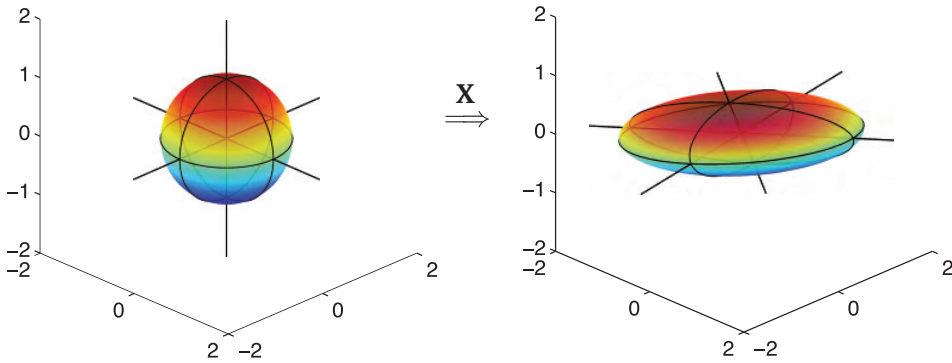


Рис. 1.8 ❖ Геометрическая иллюстрация SVD как отображения сферы в \mathbb{R}^n в эллипсоид в \mathbb{R}^m

Для частного случая на рис. 1.8 мы строим матрицу \mathbf{X} из трех матриц поворота \mathbf{R}_x , \mathbf{R}_y и \mathbf{R}_z и четвертой матрицы, описывающей растяжение вдоль главных осей:

$$\mathbf{X} = \underbrace{\begin{bmatrix} \cos(\theta_3) & -\sin(\theta_3) & 0 \\ \sin(\theta_3) & \cos(\theta_3) & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{R}_z} \underbrace{\begin{bmatrix} \cos(\theta_2) & 0 & \sin(\theta_2) \\ 0 & 1 & 0 \\ -\sin(\theta_2) & 0 & \cos(\theta_2) \end{bmatrix}}_{\mathbf{R}_y} \times \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos(\theta_1) & -\sin(\theta_1) \\ 0 & \sin(\theta_1) & \cos(\theta_1) \end{bmatrix}}_{\mathbf{R}_y} \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{bmatrix}.$$

В этом случае $\theta_1 = \pi/15$, $\theta_2 = -\pi/9$ и $\theta_3 = -\pi/20$, а $\sigma_1 = 3$, $\sigma_2 = 1$, $\sigma_3 = 0.5$. Матрицы поворота не коммутируют, поэтому порядок поворотов имеет значение. Если одно из сингулярных значений равно нулю, то соответствующая размерность исчезает и эллипсоид становится фигурой в подпространстве меньшей размерности. Произведение $\mathbf{R}_x \mathbf{R}_y \mathbf{R}_z$ является унитарной матрицей \mathbf{U} в сингулярном разложении \mathbf{X} . Матрица \mathbf{V} в данном случае единичная.

Листинг 1.1 ❖ Построение матриц поворота

```
theta = [pi/15; -pi/9; -pi/20];
Sigma = diag([3; 1; 0.5]);           % масштабировать x, y, z
Rx = [1 0 0;                        % поворот вокруг оси x
      0 cos(theta(1)) -sin(theta(1));
      0 sin(theta(1)) cos(theta(1))];
Ry = [cos(theta(2)) 0 sin(theta(2)); % поворот вокруг оси y
      0 1 0;
      -sin(theta(2)) 0 cos(theta(2))];
Rz = [cos(theta(3)) -sin(theta(3)) 0; % поворот вокруг оси z
      sin(theta(3)) cos(theta(3)) 0;
      0 0 1];
X = Rz*Ry*Rx*Sigma;                 % повернуть и масштабировать
```

Листинг 1.2 ❖ Изображение сферы

```
[x,y,z] = sphere(25);
h1=surf(x,y,z);
```

Листинг 1.3 ❖ Отображение сферы путем умножения на \mathbf{X} и изображение получившегося эллипсоида

```
xR = 0*x; yR = 0*y; zR = 0*z;
for i=1:size(x,1)
    for j=1:size(x,2)
        vecR = X*[x(i,j); y(i,j); z(i,j)];
        xR(i,j) = vecR(1);
        yR(i,j) = vecR(2);
        zR(i,j) = vecR(3);
    end
end
h2=surf(xR,yR,zR,z); % координата z сферы задает цвет
```

Инвариантность SVD относительно унитарных преобразований

У SVD есть полезное свойство: если умножить матрицу данных \mathbf{X} слева или справа на унитарную матрицу, то члены сингулярного разложения не изменятся, за исключением левой или правой унитарной матрицы \mathbf{U} или \mathbf{V} соответственно. У этого факта есть важные следствия, поскольку дискретное преобразование Фурье (ДПФ, см. главу 2) \mathcal{F} является унитарным преобразо-

ванием, а это значит, что SVD матрицы $\hat{\mathbf{X}} = \mathcal{F}\mathbf{X}$ совпадает с SVD матрицы \mathbf{X} с тем отличием, что моды $\hat{\mathbf{U}}$ будут дискретными преобразованиями Фурье мод \mathbf{U} : $\hat{\mathbf{U}} = \mathcal{F}\mathbf{U}$. Кроме того, инвариантность SVD относительно унитарных преобразований позволяет использовать сжатые измерения для реконструкции мод SVD, разреженных в некотором базисе преобразования (см. главу 3).

Инвариантность SVD относительно унитарных преобразований геометрически очевидна, поскольку унитарное преобразование лишь поворачивает векторы в пространстве, но не изменяет их скалярные произведения и структуру корреляций. Будем обозначать левое унитарное преобразование \mathbf{C} , так что $\mathbf{Y} = \mathbf{C}\mathbf{X}$, а правое унитарное преобразование \mathbf{P}^* , так что $\mathbf{Y} = \mathbf{X}\mathbf{P}^*$. SVD матрицы \mathbf{X} будем обозначать $\mathbf{U}_X \Sigma_X \mathbf{V}_X^*$, а SVD матрицы \mathbf{Y} будет равно $\mathbf{U}_Y \Sigma_Y \mathbf{V}_Y^*$.

Левые унитарные преобразования

Сначала рассмотрим левое унитарное преобразование \mathbf{X} : $\mathbf{Y} = \mathbf{C}\mathbf{X}$. Вычисление корреляционной матрицы $\mathbf{Y}^* \mathbf{Y}$ дает

$$\mathbf{Y}^* \mathbf{Y} = \mathbf{X}^* \mathbf{C}^* \mathbf{C} \mathbf{X} = \mathbf{X}^* \mathbf{X}. \quad (1.10)$$

Спроецированные данные имеют такой же спектральный состав, т. е. \mathbf{V}_X и Σ_X не изменяются. Применив метод снимков для реконструкции \mathbf{U}_Y , находим

$$\mathbf{U}_Y = \mathbf{Y} \mathbf{V}_X \Sigma_X^{-1} = \mathbf{C} \mathbf{X} \mathbf{V}_X \Sigma_X^{-1} = \mathbf{C} \mathbf{U}_X. \quad (1.11)$$

Таким образом, $\mathbf{U}_Y = \mathbf{C} \mathbf{U}_X$, $\Sigma_Y = \Sigma_X$ и $\mathbf{V}_Y = \mathbf{V}_X$. Тогда SVD матрицы \mathbf{Y} равно:

$$\mathbf{Y} = \mathbf{C} \mathbf{X} = \mathbf{C} \mathbf{U}_X \Sigma_X \mathbf{V}_X^*. \quad (1.12)$$

Правые унитарные преобразования

Для правого унитарного преобразования $\mathbf{Y} = \mathbf{X}\mathbf{P}^*$ корреляционная матрица $\mathbf{Y}^* \mathbf{Y}$ равна

$$\mathbf{Y}^* \mathbf{Y} = \mathbf{P} \mathbf{X}^* \mathbf{X} \mathbf{P}^* = \mathbf{P} \mathbf{V}_X \Sigma_X^2 \mathbf{V}_X^* \mathbf{P}^* \quad (1.13)$$

и имеет такое спектральное разложение:

$$\mathbf{Y}^* \mathbf{Y} \mathbf{P} \mathbf{V}_X = \mathbf{P} \mathbf{V}_X \Sigma_X^2. \quad (1.14)$$

Таким образом, $\mathbf{V}_Y = \mathbf{P} \mathbf{V}_X$ и $\Sigma_Y = \Sigma_X$. Мы можем воспользоваться методом снимков для реконструкции \mathbf{U}_Y :

$$\mathbf{U}_Y = \mathbf{Y} \mathbf{P} \mathbf{V}_X \Sigma_X^{-1} = \mathbf{X} \mathbf{V}_X \Sigma_X^{-1} = \mathbf{U}_X. \quad (1.15)$$

Следовательно, $\mathbf{U}_Y = \mathbf{U}_X$, и SVD матрицы \mathbf{Y} можно записать в виде:

$$\mathbf{Y} = \mathbf{X} \mathbf{P}^* = \mathbf{U}_X \Sigma_X \mathbf{V}_X^* \mathbf{P}^*. \quad (1.16)$$

1.4. ПСЕВДООБРАЩЕНИЕ, МЕТОД НАИМЕНЬШИХ КВАДРАТОВ И РЕГРЕССИЯ

Многие физические системы можно представить линейной системой уравнений

$$Ax = b, \tag{1.17}$$

где матрица ограничений A и вектор b известны, а вектор x неизвестен. Если A – квадратная обратимая матрица (т. е. определитель A не равен нулю), то существует единственное решение x для любого b . Но если A сингулярная или прямоугольная, то может существовать одно, ни одного или бесконечно много решений в зависимости от конкретного b и пространств столбцов и строк A .

Сначала рассмотрим *недоопределенную систему*, т. е. случай, когда $A \in \mathbb{C}^{n \times m}$ и $n \ll m$ (A – низкая и толстая матрица) – уравнений меньше, чем неизвестных. Такая система, скорее всего, будет иметь полный столбцовый ранг, поскольку число столбцов много больше, чем требуется для линейно независимого базиса¹. В общем случае, если низкая и толстая матрица A имеет полный столбцовый ранг, то для каждого b существует бесконечно много решений x . Такая система называется *недоопределенной*, потому что элементов b недостаточно, чтобы определить вектор x высокой размерности.

Точно так же рассмотрим *переопределенную систему*, когда $n \gg m$ (высокая и тощая матрица), т. е. уравнений больше, чем неизвестных. Эта матрица не может иметь полного столбцового ранга, поэтому гарантируется, что существуют такие векторы b , для которых нет ни одного решения x . На самом деле решение x существует, только если b принадлежит пространству столбцов A , т. е. $b \in \text{col}(A)$.

Технически могут существовать b , при которых имеется бесконечно много решений x для высокой и тощей матрицы A , равно как и такие b , при которых не существует ни одного решения для низкой и толстой матрицы. Пространство решений системы (1.17) определяется четырьмя фундаментальными подпространствами $A = \tilde{U}\tilde{\Sigma}\tilde{V}^*$, где ранг r выбран так, чтобы все ненулевые сингулярные значения были включены:

- пространство столбцов, $\text{col}(A)$, натянутое на столбцы A ; оно называется также *областью значений*. Пространство столбцов A совпадает с пространством столбцов \tilde{U} ;
- ортогональное дополнение к $\text{col}(A)$ обозначается $\text{ker}(A^*)$ и совпадает с пространством столбцов матрицы \tilde{U}^\perp на рис. 1.1;
- пространство строк, $\text{row}(A)$, натянутое на строки A и совпадающее с линейной оболочкой столбцов \tilde{V} . Имеет место равенство $\text{row}(A) = \text{col}(A^*)$;

¹ Легко построить вырожденные примеры низкой и толстой матрицы неполного столбцового ранга, например $A = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$.

- ядерное пространство, $\ker(\mathbf{A})$, являющееся ортогональным дополнением к $\text{row}(\mathbf{A})$ и называемое также *нуль-пространством*, или *ядром*. Ядро – это подпространство, состоящее из векторов, которые \mathbf{A} отображает в нуль, т. е. $\mathbf{A}\mathbf{x} = \mathbf{0}$; оно совпадает с $\text{col}(\hat{\mathbf{V}}^\perp)$.

Точнее, если $\mathbf{b} \in \text{col}(\mathbf{A})$ и $\dim(\ker(\mathbf{A})) \neq 0$, то существует бесконечно много решений \mathbf{x} . Заметим, что условие $\dim(\ker(\mathbf{A})) \neq 0$ гарантированно выполняется для низкой и толстой матрицы. Аналогично, если $\mathbf{b} \notin \text{col}(\mathbf{A})$, то решений не существует, и система уравнений (1.17) называется *несовместной*.

Описанные выше фундаментальные подпространства обладают следующими свойствами:

$$\text{col}(\mathbf{A}) \oplus \ker(\mathbf{A}^*) = \mathbb{R}^n \quad (1.18a)$$

$$\text{col}(\mathbf{A}^*) \oplus \ker(\mathbf{A}) = \mathbb{R}^n. \quad (1.18b)$$

Замечание 1. *Имеется обширная литература по теории случайных матриц, в которой перечисленные выше утверждения почти всегда верны, т. е. верны с высокой вероятностью. Например, крайне маловероятно, что система $\mathbf{A}\mathbf{x} = \mathbf{b}$ имеет решение для случайно выбранных матрицы $\mathbf{A} \in \mathbb{R}^{n \times m}$ и вектора $\mathbf{b} \in \mathbb{R}^n$, где $n \gg m$, т. к. мало шансов, что \mathbf{b} принадлежит пространству столбцов \mathbf{A} . Эти свойства случайных матриц будут играть важную роль в теории сжатых измерений (см. главу 3).*

В переопределенном случае, когда решений не существует, нам часто хотелось бы найти вектор \mathbf{x} , который минимизирует сумму квадратов ошибок $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$, он называется решением с *наименьшей среднеквадратической ошибкой*. Заметим, что такое решение минимизирует также величину $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$. В недоопределенном случае, когда существует бесконечно много решений, часто требуется найти решение \mathbf{x} , для которого $\mathbf{A}\mathbf{x} = \mathbf{b}$, а норма $\|\mathbf{x}\|_2$ минимальна; оно называется решением с *минимальной нормой*.

SVD – общепризнанный метод решения этих важных задач оптимизации. Прежде всего если подставить точное усеченное SVD $\mathbf{A} = \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^*$ вместо \mathbf{A} , то каждую из матриц $\tilde{\mathbf{U}}$, $\tilde{\Sigma}$ и $\tilde{\mathbf{V}}^*$ можно будет поочередно «обратить», что даст левую псевдообратную матрицу Мура–Пенроуза [425, 426, 453, 572] \mathbf{A}^\dagger :

$$\mathbf{A}^\dagger \triangleq \tilde{\mathbf{V}}\tilde{\Sigma}^{-1}\tilde{\mathbf{U}}^* \Rightarrow \mathbf{A}^\dagger\mathbf{A} = \mathbf{I}^{m \times m}. \quad (1.19)$$

Ее можно использовать для нахождения решений системы (1.17) с наименьшей среднеквадратической ошибкой и с минимальной нормой:

$$\mathbf{A}^\dagger\mathbf{A}\tilde{\mathbf{x}} = \mathbf{A}^\dagger\mathbf{b} \Rightarrow \tilde{\mathbf{x}} = \tilde{\mathbf{V}}\tilde{\Sigma}^{-1}\tilde{\mathbf{U}}^*\mathbf{b}. \quad (1.20)$$

Подставляя решение $\tilde{\mathbf{x}}$ обратно в (1.17), получаем:

$$\mathbf{A}\tilde{\mathbf{x}} = \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^*\tilde{\mathbf{V}}\tilde{\Sigma}^{-1}\tilde{\mathbf{U}}^*\mathbf{b} \quad (1.21a)$$

$$\mathbf{A}\tilde{\mathbf{x}} = \tilde{\mathbf{U}}\tilde{\mathbf{U}}^*\mathbf{b}. \quad (1.21b)$$

Заметим, что $\tilde{U}\tilde{U}^*$ – необязательно единичная матрица, но является проекцией на пространство столбцов \tilde{U} . Поэтому \tilde{x} будет точным решением системы (1.17), только если \mathbf{b} принадлежит пространству столбцов \tilde{U} , а значит, пространству столбцов \mathbf{A} .

Вычислить псевдообратную матрицу \mathbf{A}^\dagger можно эффективно, если предварительно произвести затратное вычисление SVD. Обращение унитарных матриц \tilde{U} и \tilde{V} сводится к умножению на транспонированные матрицы, для чего требуется $O(n^2)$ операций. Обращение диагональной матрицы $\tilde{\Sigma}$ еще эффективнее и требует всего $O(n)$ операций. С другой стороны, обращение плотной квадратной матрицы потребовало бы $O(n^3)$ операций.

Одномерная линейная регрессия

Регрессия – важный статистический инструмент установления связи между величинами на основе имеющихся данных [360]. Рассмотрим набор данных, изображенный на рис. 1.9. Точки, обозначенные красными крестиками, получены прибавлением гауссова белого шума к черной прямой, как показано в листинге 1.4. Мы предполагаем, что между данными имеется линейная связь, как в (1.17), и используем псевдообратную матрицу для нахождения решения с наименьшей среднеквадратической ошибкой – синей штриховой прямой с угловым коэффициентом x , – как показано в листинге 1.5.

$$\begin{bmatrix} | \\ | \\ | \end{bmatrix} \mathbf{b} = \begin{bmatrix} | \\ | \\ | \end{bmatrix} \mathbf{a} x = \tilde{U}\tilde{\Sigma}\tilde{V}^* x. \tag{1.22a}$$

$$\Rightarrow x = \tilde{V}\tilde{\Sigma}^{-1}\tilde{U}^* \mathbf{b}. \tag{1.22b}$$

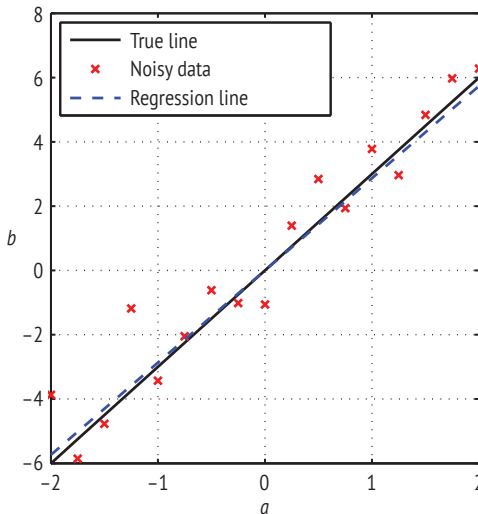


Рис. 1.9 ❖ Линейная регрессия с зашумленными данными

В (1.22b) $\tilde{\Sigma} = \|\mathbf{a}\|_2$, $\tilde{\mathbf{V}} = 1$, $\tilde{\mathbf{U}} = \mathbf{a}/\|\mathbf{a}\|_2$. Умножение на левую псевдообратную матрицу дает

$$x = \frac{\mathbf{a}^* \mathbf{b}}{\|\mathbf{a}\|_2^2}. \quad (1.23)$$

Это имеет физический смысл, если интерпретировать x как значение, которое дает наилучшее отображение нашего вектора \mathbf{a} в вектор \mathbf{b} . Такое наилучшее значение x единственно и получается в результате вычисления скалярного произведения \mathbf{b} с нормированным вектором в направлении \mathbf{a} . И мы добавляем второй нормировочный коэффициент $\|\mathbf{a}\|_2$, потому что \mathbf{a} в формуле (1.22a) не нормирован.

Отметим, что если в (1.22) взять векторы-строки вместо векторов-столбцов, то будут происходить странные вещи. Кроме того, если величина шума велика по сравнению с угловым коэффициентом x , то в точности псевдообратной матрицы произойдет фазовый переход, связанный с жесткой пороговой обработкой, описанной ниже.

Листинг 1.4 ❖ Генерирование зашумленных данных для рис. 1.9

```
x = 3; % истинный угловой коэффициент
a = [-2:.25:2]';
b = a*x + 1*randn(size(a)); % добавить шум
plot(a,x*a,'k') % истинная связь
hold on, plot(a,b,'rx') % зашумленные измерения
```

Листинг 1.5 ❖ Вычисление аппроксимации методом наименьших квадратов для рис. 1.9

```
[U,S,V] = svd(a,'econ');
xtilde = V*inv(S)*U'*b; % аппроксимации методом наименьших квадратов
plot(a,xtilde*a,'b--') % нарисовать график
```

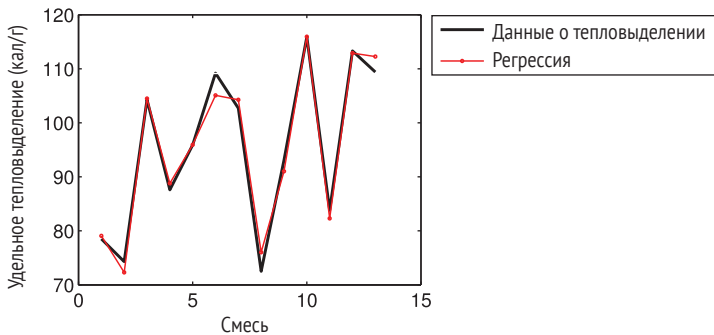


Рис. 1.10 ❖ Данные о тепловыделении для четырехкомпонентных цементных смесей

В статистике описанная выше процедура называется *линейной регрессией*. В MATLAB есть команда **regress**, а также команда **pinv**, которую тоже можно использовать.

Листинг 1.6 ❖ Альтернативные варианты метода наименьших квадратов в MATLAB

```
xtilde1 = V*inv(S)*U'*b
xtilde2 = pinv(a)*b
xtilde3 = regress(b,a)
```

Полилинейная регрессия

Пример 1: данные о тепловыделении цементных смесей

Начнем с простого встроенного в MATLAB набора данных, который описывает тепловыделение различных цементных смесей, составленных из четырех компонент. В этой задаче мы решаем систему (1.17) с $A \in \mathbb{R}^{13 \times 4}$, поскольку компонент четыре, а измерения производятся для 13 разных смесей. Наша цель – найти веса x , описывающие связь между четырьмя ингредиентами и тепловыделением. В листинге 1.7 показано, как найти решение с минимальной ошибкой. Исследуются варианты с использованием функций **regress** и **pinv**.

Листинг 1.7 ❖ Полилинейная регрессия для данных о тепловыделении цементных смесей

```
load hald; % загрузить набор данных Portland Cement
A = ingredients;
b = heat;

[U,S,V] = svd(A,'econ');
x = V*inv(S)*U'*b; % решить систему Ax=b методом SVD

plot(b,'k'); hold on % построить график
plot(A*x,'r-o',); % аппроксимация графика

x = regress(b,A); % вариант 1 (regress)
x = pinv(A)*b; % вариант 2 (pinv)
```

Пример 2: данные о недвижимости в Бостоне

В этом примере мы исследуем более крупный набор данных, чтобы определить, какие факторы лучше всего предсказывают цены на бостонском рынке недвижимости [234]. Этот набор можно скачать из репозитория машинного обучения UCI Machine Learning Repository [24].

Существует 13 атрибутов, коррелирующих с ценой дома, например: индекс преступности на душу населения и ставка налога на имущество. Мы построили регрессионную модель влияния этих факторов на цену и на рис. 1.11 на-

рисовали график, на котором лучшее предсказание цены сопоставляется с истинной стоимостью дома. На рис. 1.12 показаны коэффициенты регрессии. Хотя стоимость дома предсказана неточно, тренды хорошо согласуются. Часто бывает, что простая линейная аппроксимация плохо улавливает выбросы с наибольшими значениями; такую ситуацию мы наблюдаем и в этом примере.

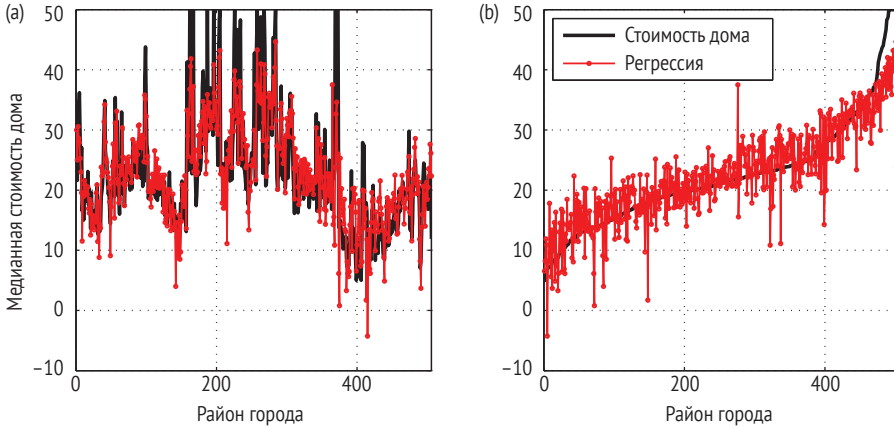


Рис. 1.11 ❖ Полилинейная регрессия цен домов с различными факторами: (а) неотсортированные данные; (б) данные отсортированы по стоимости дома

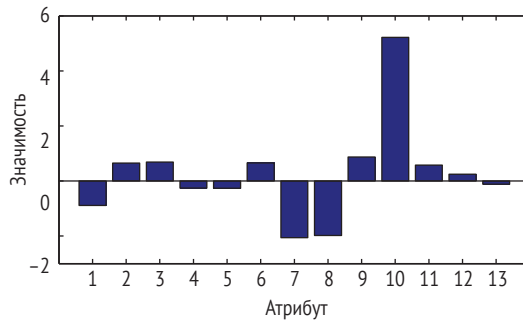


Рис. 1.12 ❖ Значимость различных атрибутов в регрессии

Этот набор данных содержит цены и атрибуты для 506 домов, так что матрица атрибутов имеет размер 506×13 . Важно дополнить эту матрицу столбцом, состоящим из единиц, чтобы учесть возможность ненулевого постоянного смещения в формуле регрессии. В случае одномерной регрессии это соответствует свободному члену.

Листинг 1.8 ❖ Полилинейная регрессия для данных о стоимости домов в Бостоне

```
load housing.data
```

```
b = housing(:,14);           % стоимости домов в тысячах долларов
A = housing(:,1:13);        % прочие факторы
```

```
A = [A ones(size(A,1),1)]; % дополнить единичными свободными членами
x = regress(b,A);
plot(b,'k-o');
hold on, plot(A*x,'r-o');

[b sortind] = sort(housing(:,14)); % отсортированные стоимости
plot(b,'k-o')
hold on, plot(A(sortind,:)*x,'r-o')
```

Предостережение

Вообще говоря, матрица U , столбцами которой являются левые сингулярные векторы X , – унитарная квадратная матрица. Поэтому $U^*U = UU^* = I_{n \times n}$. Однако чтобы вычислить псевдообратную матрицу для X , мы должны вычислить $X^\dagger = \tilde{V}\tilde{\Sigma}^{-1}\tilde{U}^*$, потому что только $\tilde{\Sigma}$ обратима (если все сингулярные значения отличны от нуля), тогда как Σ в общем случае необратима (она даже в общем случае не является квадратной).

До сих пор мы предполагали, что $X = \tilde{U}\tilde{\Sigma}\tilde{V}^*$ – точное SVD, так что ранг r включает все ненулевые сингулярные значения. Это гарантирует, что матрица $\tilde{\Sigma}$ обратима.

Трудности начинаются при работе с усеченным базисом, образованным левыми сингулярными векторами \tilde{U} . По-прежнему верно, что $\tilde{U}^*\tilde{U} = I_{r \times r}$, где r – ранг X . Однако $\tilde{U}\tilde{U}^* \neq I_{n \times n}$, что легко проверить численно на простом примере. Предположение о том, что $\tilde{U}\tilde{U}^*$ равно единичной матрице, – одно из самых типичных заблуждений при использовании SVD¹.

```
>> tol = 1.e-16;
>> [U,S,V] = svd(X,'econ')
>> r = max(find(diag(S)>max(S(:))*tol));
>> invX = V(:,1:r)*S(1:r,1:r)*U(:,1:r)'; % только приближенно
```

1.5. МЕТОД ГЛАВНЫХ КОМПОНЕНТ (PCA)

Метод главных компонент (PCA) – одно из основных применений SVD, он позволяет, ориентируясь только на данные, построить систему координат для представления коррелированных данных высокой размерности. При этом используются корреляционные матрицы, описанные в разделе 1.3. Важно, что в методе PCA, перед тем как применять SVD, производится предварительная обработка данных: вычитание среднего и приведение к единичной дисперсии. Геометрия результирующей системы координат определяется главными компонентами (principal component – PC), которые не коррелируют между собой (ортогональны), но имеют максимальную корреляцию с результатами измерений. Теория была разработана в 1901 году Пирсоном

¹ Его не избежали и авторы, по ошибке включив это неверное тождество в ранний вариант работы [96].

[418] и независимо в 1930-х годах Хотеллингом [256, 257]. Работа Jolliffe [268] содержит хорошее справочное пособие.

Обычно в одном эксперименте производится несколько измерений, которые образуют вектор-строку. Эти измерения могут быть свойствами наблюдаемой величины, например демографическими признаками одного человека. Выполняется несколько экспериментов, и все векторы-строки собираются в большую матрицу \mathbf{X} . Если говорить о демографии, то экспериментальные данные могут быть собраны путем опроса. Заметим, что такое соглашение – \mathbf{X} состоит из строк признаков – отличается от соглашения, принятого в других частях этой главы, где отдельные «снимки» признаков расположены по столбцам. Но мы решили в этом разделе не вступать в противоречие с литературой по PCA. Матрица по-прежнему имеет размер $n \times m$ и может содержать больше строк, чем столбцов, или наоборот.

Вычисление

Теперь вычислим среднее по строкам $\bar{\mathbf{x}}$ (т. е. среднее всех строк) и вычтем его из \mathbf{X} . Среднее $\bar{\mathbf{x}}$ определяется по формуле

$$\bar{x}_j = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_{ij}, \quad (1.24)$$

а средняя матрица – по формуле

$$\bar{\mathbf{X}} = \begin{bmatrix} 1 \\ \vdots \\ \bar{\mathbf{x}} \\ 1 \end{bmatrix}. \quad (1.25)$$

Вычитая $\bar{\mathbf{X}}$ из \mathbf{X} , получаем матрицу \mathbf{V} с нулевым средним:

$$\mathbf{V} = \mathbf{X} - \bar{\mathbf{X}} \quad (1.26)$$

Ковариационная матрица строк \mathbf{V} определяется формулой

$$\mathbf{C} = \frac{1}{n-1} \mathbf{V}^* \mathbf{V}. \quad (1.27)$$

Первая главная компонента \mathbf{u}_1 определяется так:

$$\mathbf{u}_1 = \operatorname{argmax}_{\|\mathbf{u}_1\|=1} \mathbf{u}_1^* \mathbf{V}^* \mathbf{V} \mathbf{u}_1. \quad (1.28)$$

Это собственный вектор $\mathbf{V}^* \mathbf{V}$, соответствующий наибольшему собственному значению. Теперь понятно, что \mathbf{u}_1 – левый сингулярный вектор \mathbf{V} , соответствующий наибольшему сингулярному значению.

Главные компоненты можно получить, вычислив спектральное разложение \mathbf{C} :

$$\mathbf{C} \mathbf{V} = \mathbf{V} \mathbf{D}, \quad (1.29)$$

которое гарантированно существует, потому что матрица \mathbf{C} эрмитова.

Команда `pca`

В MATLAB имеются дополнительные команды `pca` и `princomp` (основанная на `pca`) для метода главных компонент:

```
>> [V,score,s2] = pca(X);
```

Матрица \mathbf{V} эквивалентна матрице \mathbf{V} из сингулярного разложения \mathbf{X} с точностью до изменения знака столбцов. Вектор $\mathbf{s2}$ содержит собственные значения ковариационной матрицы \mathbf{X} , которые называются также дисперсиями главных компонент; это квадраты сингулярных значений. Переменная `score` содержит координаты каждой строки \mathbf{V} (данные после вычитания среднего) в направлениях главных компонент. Вообще говоря, мы часто предпочитаем использовать команду `svd` в сочетании с различными шагами постобработки, описанными ранее в этом разделе.

Пример: данные с гауссовым шумом

Рассмотрим зашумленное облако данных на рис. 1.13 (а), построенном программой в листинге 1.9. Данные генерируются путем выборки 10 000 векторов из двумерного нормального распределения с нулевым средним и единичной дисперсией. Эти векторы затем масштабируются в направлениях x и y с коэффициентами в табл. 1.1 и поворачиваются на угол $\pi/3$. Наконец, все облако данных подвергается параллельному переносу, так что его центр располагается в точке $\mathbf{x}_C = [2 \ 1]^T$.

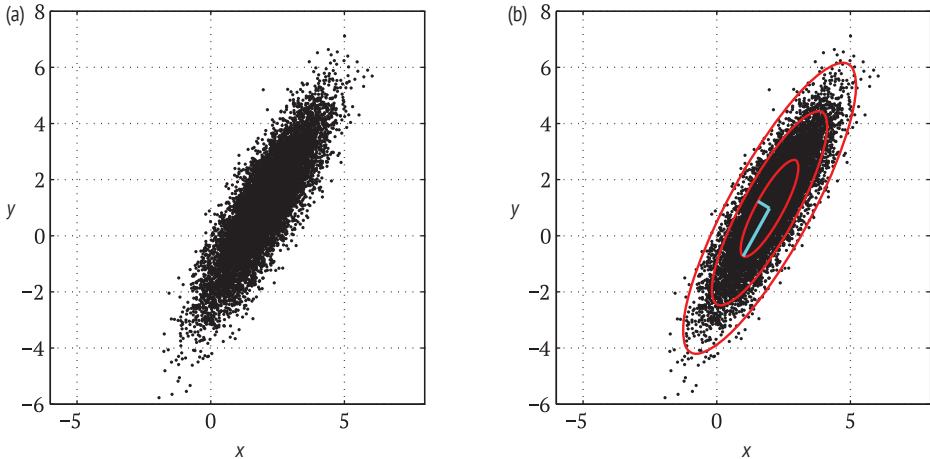


Рис. 1.13 ❖ Главные компоненты улавливают дисперсию гауссовых данных после вычитания среднего (а). На рисунке (б) показаны эллипсы для первых трех стандартных отклонений (красным цветом) и два левых сингулярных вектора, умноженных на сингулярные значения ($\sigma_1 \mathbf{u}_1 + \mathbf{x}_C$ и $\sigma_2 \mathbf{u}_2 + \mathbf{x}_C$, голубым цветом)