

**УДК 004.6  
ББК 32.972  
Л50**

- Л50 Юре Лесковец, Ананд Раджараман, Джекфри Д. Ульман  
 Анализ больших наборов данных. / Пер. с англ. Слинкин А. А. – М.: ДМК Пресс, 2016. – 498 с.: ил.

**ISBN 978-5-97060-190-7**

Эта книга написана ведущими специалистами в области технологий баз данных и веба. Она будет в равной мере полезна студентам и программистам-практикам. Благодаря популярности веба и интернет-торговли появилось много чрезвычайно объемных баз данных, для извлечения информации из которых можно применить методы добычи данных.

В книге описываются алгоритмы, которые реально использовались для решения важнейших задач добычи данных и могут быть с успехом применены даже к очень большим наборам данных. Изложение начинается с рассмотрения технологии MapReduce – важного средства распараллеливания алгоритмов. Излагаются алгоритмы хэширования с учетом близости и потоковой обработки данных, которые поступают слишком быстро для тщательного анализа. В последующих главах рассматривается идея показателя PageRank, нахождение частых предметных наборов и кластеризация. Во второе издание включен дополнительный материал о социальных сетях, машинном обучении и понижении размерности.

Original English language edition published by Cambridge University Press, 132 Avenue of the Americas, New York, NY 10013-2473, USA. Copyright © 2010, 2011, 2012, 2013, 2014 Anand Rajaraman, Jure Leskovec, Jeffrey D. Ullman. Russian-language edition copyright © 2015 by DMK Press. All rights reserved.

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

Материал, изложенный в данной книге, многократно проверен. Но, поскольку вероятность технических ошибок все равно существует, издательство не может гарантировать абсолютную точность и правильность приводимых сведений. В связи с этим издательство не несет ответственности за возможные ошибки, связанные с использованием книги.

ISBN 978-1-118-62986-4 (англ.)

ISBN 978-5-97060-190-7 (рус.)

© 2010, 2011, 2012, 2013, 2014 Anand Rajaraman, Jure Leskovec, Jeffrey D. Ullman  
 © Оформление, издание, ДМК Пресс, 2016



# ОГЛАВЛЕНИЕ

<b>Предисловие .....</b>	<b>17</b>
О чем эта книга.....	17
Требования к читателю .....	18
Упражнения.....	18
Поддержка в вебе.....	18
Автоматизированные домашние задания .....	18
Благодарности .....	19
 <b>ГЛАВА 1.</b>	
<b>Добыча данных .....</b>	<b>20</b>
1.1. Что такое добыча данных? .....	20
1.1.1. Статистическое моделирование .....	20
1.1.2. Машинное обучение .....	21
1.1.3. Вычислительные подходы к моделированию .....	21
1.1.4. Обобщение .....	22
1.1.5. Выделение признаков.....	23
1.2. Статистические пределы добычи данных .....	23
1.2.1. Тотальное владение информацией .....	24
1.2.2. Принцип Бонферрони .....	24
1.2.3. Пример применения принципа Бонферрони .....	25
1.2.4. Упражнения к разделу 1.2 .....	26
1.3. Кое-какие полезные сведения .....	26
1.3.1. Важность слов в документах .....	27
1.3.2. Хэш-функции .....	28
1.3.3. Индексы.....	29
1.3.4. Внешняя память.....	31
1.3.5. Основание натуральных логарифмов .....	31
1.3.6. Степенные зависимости .....	32
1.3.7. Упражнения к разделу 1.3 .....	34
1.4. План книги .....	35
1.5. Резюме .....	37
1.6. Список литературы .....	38

**ГЛАВА 2.**

<b>MapReduce и новый программный стек .....</b>	<b>39</b>
2.1. Распределенные файловые системы .....	40
2.1.1. Физическая организация вычислительных узлов .....	40
2.1.2. Организация больших файловых систем.....	42
2.2. MapReduce .....	42
2.2.1. Задачи-распределители .....	44
2.2.2. Группировка по ключу .....	44
2.2.3. Задачи-редукторы .....	45
2.2.4. Комбинаторы.....	45
2.2.5. Детали выполнения MapReduce .....	46
2.2.6. Обработка отказов узлов .....	48
2.2.7. Упражнения к разделу 2.2 .....	48
2.3. Алгоритмы, в которых используется MapReduce .....	48
2.3.1. Умножение матрицы на вектор с применением MapReduce .....	49
2.3.2. Если вектор $v$ не помещается в оперативной памяти.....	50
2.3.3. Операции реляционной алгебры.....	51
2.3.4. Вычисление выборки с помощью MapReduce .....	53
2.3.5. Вычисление проекции с помощью MapReduce.....	54
2.3.6. Вычисление объединения, пересечения и разности с помощью MapReduce .....	54
2.3.7. Вычисление естественного соединения с помощью MapReduce.....	55
2.3.8. Вычисление группировки и агрегирования с помощью MapReduce ....	56
2.3.9. Умножение матриц .....	56
2.3.10. Умножение матриц за один шаг MapReduce.....	57
2.3.11. Упражнения к разделу 2.3 .....	58
2.4. Обобщения MapReduce .....	59
2.4.1. Системы потоков работ .....	60
2.4.2. Рекурсивные обобщения MapReduce.....	61
2.4.3. Система Pregel .....	64
2.4.4. Упражнения к разделу 2.4 .....	65
2.5. Модель коммуникационной стоимости .....	65
2.5.1. Коммуникационная стоимость для сетей задач.....	65
2.5.2. Физическое время .....	68
2.5.3. Многопутевое соединение.....	68
2.5.4. Упражнения к разделу 2.5 .....	71
2.6. Теория сложности MapReduce .....	73
2.6.1. Размер редукции и коэффициент репликации .....	73
2.6.2. Пример: соединение по сходству.....	74
2.6.3. Графовая модель для проблем MapReduce.....	76
2.6.4. Схема сопоставления .....	78
2.6.5. Когда присутствуют не все входы.....	79
2.6.6. Нижняя граница коэффициента репликации .....	80
2.6.7. Пример: умножение матриц.....	82
2.6.8. Упражнения к разделу 2.6 .....	86
2.7. Резюме .....	87



2.8. Список литературы .....	89
------------------------------	----

## ГЛАВА 3.

### Поиск похожих объектов ..... 92

3.1. Приложения поиска близкого соседям .....	92
3.1.1. Сходство множеств по Жаккарду .....	93
3.1.2. Сходство документов.....	93
3.1.3. Коллаборативная фильтрация как задача о сходстве множеств .....	94
3.1.4. Упражнения к разделу 3.1 .....	96
3.2. Разбиение документов на шинглы .....	96
3.2.1. k-шинглы .....	97
3.2.2. Выбор размера шингла .....	97
3.2.3. Хэширование шинглов .....	98
3.2.4. Шинглы, построенные из слов .....	98
3.2.5. Упражнения к разделу 3.2 .....	99
3.3. Сигнатуры множеств с сохранением сходства .....	100
3.3.1. Матричное представление множеств.....	100
3.3.2. Минхэш .....	101
3.3.3. Минхэш и коэффициент Жаккара.....	102
3.3.4. Минхэш-сигнатуры .....	102
3.3.5. Вычисление минхэш-сигнатур .....	103
3.3.6. Упражнения к разделу 3.3 .....	105
3.4. Хэширование документов с учетом близости.....	107
3.4.1. LSH для минхэш-сигнатур .....	107
3.4.2. Анализ метода разбиения на полосы .....	109
3.4.3. Сочетание разных методов .....	110
3.4.4. Упражнения к разделу 3.4 .....	111
3.5. Метрики.....	111
3.5.1. Определение метрики .....	112
3.5.2. Евклидовы метрики .....	112
3.5.3. Расстояние Жаккара.....	113
3.5.4. Косинусное расстояние .....	114
3.5.5. Редакционное расстояние .....	114
3.5.6. Расстояние Хэмминга.....	115
3.5.7. Упражнения к разделу 3.5 .....	116
3.6. Теория функций, учитывающих близость .....	118
3.6.1. Функции, учитывающие близость .....	119
3.6.2. LSH-семейства для расстояния Жаккара .....	120
3.6.3. Расширение LSH-семейства .....	120
3.6.4. Упражнения к разделу 3.6 .....	122
3.7. LSH-семейства для других метрик .....	123
3.7.1. LSH-семейства для расстояния Хэмминга .....	123
3.7.2. Случайные гиперплоскости и косинусное расстояние.....	124
3.7.3 Эскизы.....	125
3.7.4. LSH-семейства для евклидова расстояния .....	126
3.7.5. Другие примеры LSH-семейств в евклидовых пространствах .....	127

3.7.6. Упражнения к разделу 3.7 .....	128
<b>3.8. Применения хэширования с учетом близости .....</b>	<b>129</b>
3.8.1. Отождествление объектов .....	129
3.8.2. Пример отождествления объектов .....	129
3.8.3. Проверка отождествления записей .....	131
3.8.4. Сравнение отпечатков пальцев .....	132
3.8.5. LSH-семейство для сравнения отпечатков пальцев .....	132
3.8.6. Похожие новости .....	134
3.8.7. Упражнения к разделу 3.8 .....	135
<b>3.9. Методы для высокой степени сходства .....</b>	<b>136</b>
3.9.1. Поиск одинаковых объектов .....	137
3.9.2. Представление множеств в виде строк .....	137
3.9.3. Фильтрация на основе длины строки .....	138
3.9.4. Префиксное индексирование .....	138
3.9.5. Использование информации о позиции .....	140
3.9.6. Использование позиции и длины в индексах .....	141
3.9.7. Упражнения к разделу 3.9 .....	144
<b>3.10. Резюме .....</b>	<b>144</b>
<b>3.11. Список литературы .....</b>	<b>147</b>

**ГЛАВА 4.****Анализ потоков данных.....** **149**

<b>4.1. Потоковая модель данных.....</b>	<b>149</b>
4.1.1. Система управления потоками данных .....	150
4.1.2. Примеры источников потоков данных .....	151
4.1.3. Запросы к потокам.....	152
4.1.4. Проблемы обработки потоков.....	153
<b>4.2. Выборка данных из потока .....</b>	<b>154</b>
4.2.1. Пояснительный пример .....	154
4.2.2. Получение репрезентативной выборки .....	155
4.2.3. Общая постановка задачи о выборке .....	155
4.2.4. Динамическое изменение размера выборки.....	156
4.2.5. Упражнения к разделу 4.2 .....	156
<b>4.3. Фильтрация потоков .....</b>	<b>157</b>
4.3.1. Пояснительный пример .....	157
4.3.2. Фильтр Блума .....	158
4.3.3. Анализ фильтра Блума .....	158
4.3.4. Упражнения к разделу 4.3 .....	160
<b>4.4. Подсчет различных элементов в потоке .....</b>	<b>160</b>
4.4.1. Проблема Count-Distinct .....	160
4.4.2. Алгоритм Флажоле-Мартена .....	161
4.4.3. Комбинирование оценок.....	162
4.4.4. Требования к памяти.....	163
4.4.5. Упражнения к разделу 4.4 .....	163
<b>4.5. Оценивание моментов .....</b>	<b>163</b>
4.5.1. Определение моментов .....	163

## Оглавление



4.5.2. Алгоритм Алона-Матиаса-Сегеди для вторых моментов .....	164
4.5.3. Почему работает алгоритм Алона-Матиаса-Сегеди .....	165
4.5.4. Моменты высших порядков.....	166
4.5.5. Обработка бесконечных потоков.....	166
4.5.6. Упражнения к разделу 4.5 .....	168
<b>4.6. Подсчет единиц в окне .....</b>	<b>169</b>
4.6.1. Стоимость точного подсчета.....	169
4.6.2. Алгоритм Датара-Гиониса-Индыка-Мотвани .....	170
4.6.3. Требования к объему памяти для алгоритма DGIM .....	171
4.6.4. Ответы на вопросы в алгоритме DGIM.....	172
4.6.5. Поддержание условий DGIM .....	172
4.6.6. Уменьшение погрешности .....	174
4.6.7. Обобщения алгоритма подсчета единиц.....	174
4.6.8. Упражнения к разделу 4.6 .....	175
<b>4.7. Затухающие окна .....</b>	<b>176</b>
4.7.1. Задача о самых частых элементах.....	176
4.7.2. Определение затухающего окна .....	176
4.7.3. Нахождение самых популярных элементов .....	177
<b>4.8. Резюме .....</b>	<b>178</b>
<b>4.9. Список литературы .....</b>	<b>180</b>

## ГЛАВА 5.

### Анализ ссылок ..... **182**

<b>5.1. PageRank .....</b>	<b>182</b>
5.1.1. Ранние поисковые системы и спам термов .....	183
5.1.2. Определение PageRank .....	184
5.1.3. Структура веба .....	187
5.1.4. Избегание тупиков.....	189
5.1.5. Паучьи ловушки и телепортация .....	192
5.1.6. Использование PageRank в поисковой системе .....	194
5.1.7. Упражнения к разделу 5.1 .....	194
<b>5.2. Эффективное вычисление PageRank.....</b>	<b>196</b>
5.2.1. Представление матрицы переходов .....	196
5.2.2. Итеративное вычисление PageRank с помощью MapReduce .....	197
5.2.3. Использование комбинаторов для консолидации результатирующего вектора.....	198
5.2.4. Представление блоков матрицы переходов .....	199
5.2.5. Другие эффективные подходы к итеративному вычислению PageRank .....	200
5.2.6. Упражнения к разделу 5.2 .....	201
<b>5.3. Тематический PageRank.....</b>	<b>202</b>
5.3.1. Зачем нужен тематический PageRank .....	202
5.3.2. Смещенное случайное блуждание .....	202
5.3.3. Использование тематического PageRank.....	204
5.3.4. Вывод тем из слов .....	205
5.3.5. Упражнения к разделу 5.3 .....	205



5.4. Ссылочный спам .....	206
5.4.1. Архитектура спам-фермы .....	206
5.4.2. Анализ спам-фермы .....	207
5.4.3. Борьба со ссылочным спамом .....	208
5.4.4. TrustRank .....	208
5.4.5. Спамная масса .....	209
5.4.6. Упражнения к разделу 5.4 .....	210
5.5. Хабы и авторитетные страницы.....	210
5.5.1. Предположения, лежащие в основе HITS .....	211
5.5.2. Формализация хабов и авторитетных страниц.....	211
5.5.3. Упражнения к разделу 5.5 .....	214
5.6. Резюме .....	214
5.7. Список литературы .....	218

## ГЛАВА 6.

### Частые предметные наборы ..... **219**

6.1. Модель корзины покупок .....	219
6.1.1. Определение частого предметного набора.....	220
6.1.2. Применения частых предметных наборов .....	221
6.1.3. Ассоциативные правила .....	223
6.1.4. Поиск ассоциативных правил с высокой достоверностью.....	225
6.1.5. Упражнения к разделу 6.1 .....	225
6.2. Корзины покупок и алгоритм Apriori .....	226
6.2.1. Представление данных о корзинах покупок.....	227
6.2.2. Использование оперативной памяти для подсчета предметных наборов.....	228
6.2.3. Монотонность предметных наборов .....	230
6.2.4. Доминирование подсчета пар.....	230
6.2.5. Алгоритм Apriori .....	231
6.2.6. Применение Apriori для поиска всех частых предметных наборов .....	232
6.2.7. Упражнения к разделу 6.2 .....	235
6.3. Обработка больших наборов данных в оперативной памяти.....	236
6.3.1. Алгоритм Парка-Чена-Ю (PCY) .....	236
6.3.2. Многоэтапный алгоритм .....	238
6.3.3. Многохэшевый алгоритм .....	240
6.3.4. Упражнения к разделу 6.3 .....	242
6.4. Алгоритм с ограниченным числом проходов .....	244
6.4.1. Простой рандомизированный алгоритм .....	244
6.4.2. Предотвращение ошибок в алгоритмах формирования выборки .....	245
6.4.3. Алгоритм SON.....	246
6.4.4. Алгоритм SON и MapReduce .....	247
6.4.5. Алгоритм Тойвонена .....	248
6.4.6. Почему алгоритм Тойвонена работает .....	249
6.4.7. Упражнения к разделу 6.4 .....	249
6.5. Подсчет частых предметных наборов в потоке .....	250
6.5.1. Методы выборки из потока .....	250

## Оглавление

6.5.2. Частые предметные наборы в затухающих окнах .....	251
6.5.3. Гибридные методы .....	253
6.5.4. Упражнения к разделу 6.5 .....	253
6.6. Резюме .....	254
6.7. Список литературы .....	256
<b>ГЛАВА 7.</b>	
<b>Кластеризация.....</b>	<b>258</b>
7.1. Введение в методы кластеризации .....	258
7.1.1. Точки, пространства и расстояния .....	258
7.1.2. Стратегии кластеризации .....	260
7.1.3. Проклятие размерности.....	260
7.1.4. Упражнения к разделу 7.1 .....	262
7.2. Иерархическая кластеризация.....	262
7.2.1. Иерархическая кластеризация в евклидовом пространстве .....	263
7.2.2. Эффективность иерархической кластеризации .....	265
7.2.3. Альтернативные правила управления иерархической кластеризацией .....	266
7.2.4. Иерархическая кластеризация в неевклидовых пространствах.....	268
7.2.5. Упражнения к разделу 7.2 .....	269
7.3. Алгоритм k средних .....	270
7.3.1. Основы алгоритма k средних .....	270
7.3.2. Инициализация кластеров в алгоритме k средних.....	271
7.3.3. Выбор правильного значения k .....	272
7.3.4. Алгоритм Брэдли-Файяда-Рейна .....	273
7.3.5. Обработка данных в алгоритме BFR .....	275
7.3.6. Упражнения к разделу 7.3 .....	277
7.4. Алгоритм CURE .....	278
7.4.1. Этап инициализации в CURE.....	278
7.4.2. Завершение работы алгоритма CURE .....	279
7.4.3. Упражнения к разделу 7.4 .....	280
7.5. Кластеризация в неевклидовых пространствах .....	280
7.5.1. Представление кластеров в алгоритме GRGPF .....	281
7.5.2. Инициализация дерева кластеров .....	281
7.5.3. Добавление точек в алгоритме GRGPF .....	282
7.5.4. Разделение и объединение кластеров .....	283
7.5.5. Упражнения к разделу 7.5 .....	285
7.6. Кластеризация для потоков и параллелизм .....	285
7.6.1. Модель потоковых вычислений.....	285
7.6.2. Алгоритм кластеризации потока .....	286
7.6.3. Инициализация интервалов .....	286
7.6.4. Объединение кластеров .....	287
7.6.5. Ответы на вопросы .....	289
7.6.6. Кластеризация в параллельной среде .....	290
7.6.7. Упражнения к разделу 7.6 .....	290
7.7. Резюме .....	290



7.8. Список литературы .....	294
------------------------------	-----

## ГЛАВА 8.

### Реклама в Интернете.....**295**

8.1. Проблемы онлайновой рекламы .....	295
8.1.1. Возможности рекламы.....	295
8.1.2. Прямое размещение рекламы.....	296
8.1.3. Акцидентные объявления.....	297
8.2. Онлайновые алгоритмы .....	298
8.2.1. Онлайновые и офлайновые алгоритмы .....	298
8.2.2. Жадные алгоритмы .....	299
8.2.3. Коэффициент конкурентоспособности .....	300
8.2.4. Упражнения к разделу 8.2 .....	300
8.3. Задача о паросочетании .....	301
8.3.1. Паросочетания и совершенные паросочетания .....	301
8.3.2. Жадный алгоритм нахождения максимального паросочетания .....	302
8.3.3. Коэффициент конкурентоспособности жадного алгоритма паросочетания .....	303
8.3.4. Упражнения к разделу 8.3 .....	304
8.4. Задача о ключевых словах.....	304
8.4.1. История поисковой рекламы.....	304
8.4.2. Постановка задачи о ключевых словах .....	305
8.4.3. Жадный подход к задаче о ключевых словах .....	306
8.4.4. Алгоритм Balance.....	307
8.4.5. Нижняя граница коэффициента конкурентоспособности в алгоритме Balance .....	308
8.4.6. Алгоритм Balance при большом числе участников аукциона.....	310
8.4.7. Обобщенный алгоритм Balance .....	311
8.4.8. Заключительные замечания по поводу задачи о ключевых словах....	312
8.4.9. Упражнения к разделу 8.4 .....	313
8.5. Реализация алгоритма Adwords .....	313
8.5.1. Сопоставление предложений с поисковыми запросами .....	314
8.5.2. Более сложные задачи сопоставления.....	314
8.5.3. Алгоритм сопоставления документов и ценовых предложений .....	315
8.6. Резюме .....	318
8.7. Список литературы .....	320

## ГЛАВА 9.

### Рекомендательные системы .....

<b>321</b>	
9.1. Модель рекомендательной системы .....	321
9.1.1. Матрица предпочтений.....	322
9.1.2. Длинный хвост .....	323
9.1.3. Применения рекомендательных систем.....	323
9.1.4. Заполнение матрицы предпочтений .....	325
9.2. Рекомендации на основе фильтрации содержимого .....	326



## Оглавление

9.2.1. Профили объектов.....	326
9.2.2. Выявление признаков документа.....	327
9.2.3. Получение признаков объектов из меток .....	328
9.2.4. Представление профиля объекта.....	329
9.2.5. Профили пользователей.....	330
9.2.6. Рекомендование объектов пользователям на основе содержимого....	331
9.2.7. Алгоритм классификации .....	332
9.2.8. Упражнения к разделу 9.2 .....	335
<b>9.3. Коллаборативная фильтрация.....</b>	<b>336</b>
9.3.1. Измерение сходства .....	336
9.3.2. Двойственность сходства .....	339
9.3.3. Кластеризация пользователей и объектов .....	340
9.3.4. Упражнения к разделу 9.3 .....	341
<b>9.4. Понижение размерности .....</b>	<b>342</b>
9.4.1. UV-декомпозиция .....	343
9.4.2. Среднеквадратичная ошибка .....	343
9.4.3. Инкрементное вычисление UV-декомпозиции .....	344
9.4.4. Оптимизация произвольного элемента.....	347
9.4.5. Построение полного алгоритма UV-декомпозиции .....	348
9.4.6. Упражнения к разделу 9.4 .....	351
<b>9.5. Задача NetFlix .....</b>	<b>351</b>
<b>9.6. Резюме .....</b>	<b>353</b>
<b>9.7. Список литературы .....</b>	<b>355</b>

## ГЛАВА 10.

### Анализ графов социальных сетей ..... **356**

<b>10.1. Социальные сети как графы .....</b>	<b>356</b>
10.1.1. Что такое социальная сеть? .....	357
10.1.2. Социальные сети как графы .....	357
10.1.3. Разновидности социальных сетей.....	358
10.1.4. Графы с вершинами нескольких типов .....	360
10.1.5. Упражнения к разделу 10.1 .....	361
<b>10.2. Кластеризация графа социальной сети.....</b>	<b>361</b>
10.2.1. Метрики для графов социальных сетей.....	361
10.2.2. Применение стандартных методов кластеризации .....	362
10.2.3. Промежуточность .....	363
10.2.4. Алгоритм Гирвана-Ньюмана.....	364
10.2.5. Использование промежуточности для нахождения сообществ.....	366
10.2.6. Упражнения к разделу 10.2 .....	368
<b>10.3. Прямое нахождение сообществ .....</b>	<b>368</b>
10.3.1. Нахождение клик .....	368
10.3.2. Полные двудольные графы .....	369
10.3.3. Нахождение полных двудольных подграфов .....	370
10.3.4. Почему должны существовать полные двудольные графы .....	370
10.3.5. Упражнения к разделу 10.3 .....	372



10.4. Разрезание графов .....	373
10.4.1. Какое разрезание считать хорошим? .....	373
10.4.2. Нормализованные разрезы.....	374
10.4.3. Некоторые матрицы, описывающие графы .....	374
10.4.4. Собственные значения матрицы Лапласа .....	375
10.4.5. Другие методы разрезания .....	378
10.4.6. Упражнения к разделу 10.4 .....	379
10.5. Нахождение пересекающихся сообществ .....	379
10.5.1. Природа сообществ.....	379
10.5.2. Оценка максимального правдоподобия .....	380
10.5.3. Модель графа принадлежности .....	382
10.5.4. Как избежать дискретных изменений членства .....	384
10.5.5. Упражнения к разделу 10.5 .....	385
10.6. Simrank .....	386
10.6.1. Случайные блуждания в социальном графе .....	386
10.6.2. Случайное блуждание с перезапуском .....	387
10.6.3. Упражнения к разделу 10.6 .....	389
10.7. Подсчет треугольников .....	390
10.7.1. Зачем подсчитывать треугольники? .....	390
10.7.2. Алгоритм нахождения треугольников.....	390
10.7.3. Оптимальность алгоритма нахождения треугольников.....	392
10.7.4. Нахождение треугольников с помощью MapReduce .....	392
10.7.5. Использование меньшего числа редукторов .....	394
10.7.6. Упражнения к разделу 10.7 .....	395
10.8. Окрестности в графах .....	396
10.8.1. Ориентированные графы и окрестности .....	396
10.8.2. Диаметр графа .....	397
10.8.3. Транзитивное замыкание и достижимость .....	399
10.8.4. Вычисление транзитивного замыкания с помощью MapReduce .....	399
10.8.5. Интеллектуальное транзитивное замыкание .....	402
10.8.6. Транзитивное замыкание посредством сокращения графа .....	403
10.8.7. Аппроксимация размеров окрестностей .....	405
10.8.8. Упражнения к разделу 10.8 .....	407
10.9. Резюме .....	408
10.10. Список литературы .....	411
<b>ГЛАВА 11.</b>	
<b>Понижение размерности .....</b>	<b>414</b>
11.1. Собственные значения и собственные векторы.....	414
11.1.1. Определения .....	415
11.1.2. Вычисление собственных значений и собственных векторов .....	415
11.1.3. Нахождение собственных пары степенным методом .....	417
11.1.4. Матрица собственных векторов .....	420
11.1.5. Упражнения к разделу 11.1 .....	421
11.2. Метод главных компонент .....	422
11.2.1. Иллюстративный пример .....	422

## Оглавление



11.2.2. Использование собственных векторов для понижения размерности .....	425
11.2.3. Матрица расстояний.....	426
11.2.4. Упражнения к разделу 11.2 .....	427
11.3. Сингулярное разложение.....	427
11.3.1. Определение сингулярного разложения .....	428
11.3.2. Интерпретация сингулярного разложения .....	429
11.3.3. Понижение размерности с помощью сингулярного разложения .....	431
11.3.4. Почему обнуление малых сингулярных значений работает.....	432
11.3.5. Запросы с использованием концептов.....	434
11.3.6. Вычисление сингулярного разложения матрицы.....	434
11.3.7. Упражнения к разделу 11.3 .....	435
<b>11.4. CUR-декомпозиция.....</b>	<b>436</b>
11.4.1. Определение CUR-декомпозиции.....	437
11.4.2. Правильный выбор строк и столбцов .....	438
11.4.3. Построение средней матрицы .....	440
11.4.4. Полная CUR-декомпозиция.....	441
11.4.5. Исключение дубликатов строк и столбцов.....	441
11.4.6. Упражнения к разделу 11.4 .....	442
11.5. Резюме .....	442
11.6. Список литературы .....	444

**ГЛАВА 12.****Машинное обучение на больших данных ..... 446**

12.1. Модель машинного обучения .....	447
12.1.1. Обучающие наборы.....	447
12.1.2. Пояснительные примеры .....	447
12.1.3. Подходы к машинному обучению .....	449
12.1.4. Архитектура машинного обучения.....	451
12.1.5. Упражнения к разделу 12.1 .....	454
<b>12.2. Перцептроны .....</b>	<b>454</b>
12.2.1. Обучение перцептрана с нулевым порогом.....	455
12.2.2. Сходимость перцептронов.....	457
12.2.3. Алгоритм Winnow .....	458
12.2.4. Переменный порог.....	459
12.2.5. Многоклассовые перцептроны.....	461
12.2.6. Преобразование обучающего набора .....	462
12.2.7. Проблемы, связанные с перцептранами .....	463
12.2.8. Параллельная реализация перцептронов .....	464
12.2.9. Упражнения к разделу 12.2 .....	466
<b>12.3. Метод опорных векторов.....</b>	<b>466</b>
12.3.1. Механизм метода опорных векторов.....	466
12.3.2. Нормировка гиперплоскости .....	468
12.3.3. Нахождение оптимальных приближенных разделителей.....	470
12.3.4. Нахождение решений в методе опорных векторов с помощью градиентного спуска .....	472



12.3.5. Стохастический градиентный спуск .....	476
12.3.6. Параллельная реализация метода опорных векторов .....	477
12.3.7. Упражнения к разделу 12.3 .....	477
<b>12.4. Обучение по ближайшим соседям.....</b>	<b>478</b>
12.4.1. Инфраструктура для вычисления ближайших соседей .....	478
12.4.2. Обучение по одному ближайшему соседу .....	479
12.4.3. Обучение одномерных функций .....	480
12.4.4. Ядерная регрессия .....	482
12.4.5. Данные в многомерном евклидовом пространстве .....	483
12.4.6. Неевклидовы метрики.....	484
12.4.7. Упражнения к разделу 12.4 .....	485
12.5. Сравнение методов обучения .....	486
12.6. Резюме .....	487
12.7. Список литературы .....	489
<b>Предметный указатель .....</b>	<b>490</b>