

*Посвящается памяти А. Гарри Клофа*

# Содержание

<b>Вступительное слово</b> .....	11
<b>Предисловие ко второму изданию</b> .....	12
<b>Предисловие к первому изданию</b> .....	17
<b>Обозначения</b> .....	20
<b>От издательства</b> .....	25
<b>Глава 1. Введение</b> .....	26
1.1. Обучение с подкреплением.....	26
1.2. Примеры .....	30
1.3. Элементы обучения с подкреплением .....	31
1.4. Ограничения и круг вопросов .....	33
1.5. Развернутый пример: игра в крестики-нолики.....	34
1.6. Резюме .....	39
1.7. История ранних этапов обучения с подкреплением .....	39
Библиографические замечания.....	49
<b>Часть I. ТАБЛИЧНЫЕ МЕТОДЫ РЕШЕНИЯ</b> .....	50
<b>Глава 2. Многорукие бандиты</b> .....	51
2.1. Задача о $k$ -руком бандите.....	51
2.2. Методы ценности действий .....	53
2.3. 10-рукий испытательный стенд.....	54
2.4. Инкрементная реализация.....	57
2.5. Нестационарная задача .....	59
2.6. Оптимистические начальные значения.....	60
2.7. Выбор действия, дающего верхнюю доверительную границу .....	62
2.8. Градиентные алгоритмы бандита.....	64
2.9. Ассоциативный поиск (контекстуальные бандиты).....	68
2.10. Резюме .....	69
Библиографические и исторические замечания.....	71
<b>Глава 3. Конечные марковские процессы принятия решений</b> .....	74
3.1. Интерфейс между агентом и окружающей средой.....	74
3.2. Цели и вознаграждения .....	80
3.3. Доход и эпизоды .....	82
3.4. Унифицированная нотация для эпизодических и непрерывных задач .....	84
3.5. Стратегии и функции ценности .....	86
3.6. Оптимальные стратегии и оптимальные функции ценности .....	91
3.7. Оптимальность и аппроксимация .....	96
3.8. Резюме .....	97
Библиографические и исторические замечания.....	99

<b>Глава 4. Динамическое программирование</b> .....	102
4.1. Оценивание стратегии (предсказание).....	103
4.2. Улучшение стратегии.....	107
4.3. Итерация по стратегиям.....	109
4.4. Итерация по ценности.....	112
4.5. Асинхронное динамическое программирование.....	115
4.6. Обобщенная итерация по стратегиям.....	116
4.7. Эффективность динамического программирования.....	117
4.8. Резюме.....	118
Библиографические и исторические замечания.....	119
<b>Глава 5. Методы Монте-Карло</b> .....	122
5.1. Предсказание методами Монте-Карло.....	123
5.2. Оценивание ценности действий методом Монте-Карло.....	127
5.3. Управление методом Монте-Карло.....	129
5.4. Управление методом Монте-Карло без исследовательских стартов.....	132
5.5. Предсказание с разделенной стратегией посредством выборки по значимости.....	135
5.6. Инкрементная реализация.....	142
5.7. Управление методом Монте-Карло с разделенной стратегией.....	143
5.8. *Выборка по значимости с учетом обесценивания.....	146
5.9. *Приведенная выборка по значимости.....	147
5.10. Резюме.....	149
Библиографические и исторические замечания.....	150
<b>Глава 6. Обучение на основе временных различий</b> .....	152
6.1. Предсказание TD-методами.....	152
6.2. Преимущества TD-методов предсказания.....	157
6.3. Оптимальность TD(0).....	159
6.4. Sarsa: TD-управление с единой стратегией.....	162
6.5. Q-обучение: TD-управление с разделенной стратегией.....	165
6.6. Expected Sarsa.....	167
6.7. Смещение максимизации и двойное обучение.....	169
6.8. Игры, послесостояния и другие специальные случаи.....	171
6.9. Резюме.....	173
Библиографические и исторические замечания.....	174
<b>Глава 7. <math>n</math>-шаговый бутстрэппинг</b> .....	176
7.1. $n$ -шаговое TD-предсказание.....	176
7.2. $n$ -шаговый алгоритм Sarsa.....	181
7.3. $n$ -шаговое обучение с разделенной стратегией.....	184
7.4. *Приведенные методы с переменным управлением.....	186
7.5. Обучение с разделенной стратегией без выборки по значимости: $n$ -шаговый алгоритм обновления по дереву.....	188
7.6. *Унифицированный алгоритм: $n$ -шаговый Q( $\sigma$ ).....	190
7.7. Резюме.....	193
Библиографические и исторические замечания.....	194
<b>Глава 8. Планирование и обучение табличными методами</b> .....	195
8.1. Модели и планирование.....	195
8.2. Дуна: объединение планирования, исполнения и обучения.....	198

8.3. Когда модель неверна .....	203
8.4. Приоритетный проход .....	206
8.5. Сравнение выборочного и полного обновлений .....	210
8.6. Траекторная выборка .....	213
8.7. Динамическое программирование в реальном времени .....	216
8.8. Планирование в момент принятия решений .....	220
8.9. Эвристический поиск .....	221
8.10. Разыгрывающие алгоритмы .....	224
8.11. Поиск по дереву методом Монте-Карло .....	226
8.12. Резюме главы .....	229
8.13. Резюме части I: оси .....	230
Библиографические и исторические замечания .....	233

## **Часть II. ПРИБЛИЖЕННЫЕ МЕТОДЫ РЕШЕНИЯ .....** 236

### **Глава 9. Предсказание с единой стратегией и аппроксимацией .....** 238

9.1. Аппроксимация функции ценности .....	239
9.2. Целевая функция предсказания ( $\sqrt{E}$ ) .....	240
9.3. Стохастические градиентные и полуградиентные методы .....	242
9.4. Линейные методы .....	246
9.5. Конструирование признаков для линейных методов .....	252
9.5.1. Полиномы .....	252
9.5.2. Базис Фурье .....	254
9.5.3. Грубое кодирование .....	257
9.5.4. Плиточное кодирование .....	260
9.5.5. Радиально-базисные функции .....	265
9.6. Выбор размера шага вручную .....	266
9.7. Нелинейная аппроксимация функций: искусственные нейронные сети .....	267
9.8. Алгоритм TD наименьших квадратов .....	272
9.9. Аппроксимация функций с запоминанием .....	274
9.10. Аппроксимация с помощью ядерных функций .....	276
9.11. Более глубокий взгляд на обучение с единой стратегией: заинтересованность и значимость .....	278
9.12. Резюме .....	280
Библиографические и исторические замечания .....	281

### **Глава 10. Управление с единой стратегией и аппроксимацией .....** 288

10.1. Эпизодическое полуградиентное управление .....	288
10.2. Полуградиентный $n$ -шаговый Sarsa .....	292
10.3. Среднее вознаграждение: новая постановка непрерывных задач .....	294
10.4. Возражения против постановки с обесцениванием .....	299
10.5. Дифференциальный полуградиентный $n$ -шаговый Sarsa .....	301
10.6. Резюме .....	302
Библиографические и исторические замечания .....	303

### **Глава 11. \*Методы с разделенной стратегией и аппроксимацией .....** 304

11.1. Полуградиентные методы .....	305
11.2. Примеры расходимости в случае с разделенной стратегией .....	307
11.3. Смертельная триада .....	312

11.4. Геометрия линейной аппроксимации функций ценности .....	314
11.5. Градиентный спуск по беллмановской ошибке .....	318
11.6. Беллмановская ошибка необучаема .....	322
11.7. Градиентные TD-методы .....	327
11.8. Эмфатические TD-методы .....	330
11.9. Уменьшение дисперсии.....	332
11.10. Резюме .....	334
Библиографические и исторические замечания.....	335

## **Глава 12. Следы приемлемости .....**

12.1. $\lambda$ -доход .....	338
12.2. TD( $\lambda$ ) .....	342
12.3. $n$ -шаговые усеченные $\lambda$ -доходные методы .....	346
12.4. Пересчет обновлений: онлайнный $\lambda$ -доходный алгоритм .....	348
12.5. Истинно онлайнный TD( $\lambda$ ).....	350
12.6. *Голландские следы в обучении методами Монте-Карло .....	352
12.7. Sarsa( $\lambda$ ).....	354
12.8. Переменные $\lambda$ и $\gamma$ .....	359
12.9. Следы с разделенной стратегией и переменным управлением.....	361
12.10. От Q( $\lambda$ ) Уоткинса к Tree-Backup( $\lambda$ ).....	364
12.11. Устойчивые методы с разделенной стратегией со следами приемлемости .....	367
12.12. Вопросы реализации.....	368
12.13. Выводы.....	369
Библиографические и исторические замечания.....	371

## **Глава 13. Методы градиента стратегии .....**

13.1. Аппроксимация стратегии и ее преимущества .....	374
13.2. Теорема о градиенте стратегии .....	376
13.3. REINFORCE: метод Монте-Карло на основе градиента стратегии .....	378
13.4. REINFORCE с базой.....	381
13.5. Методы исполнитель–критик.....	383
13.6. Метод градиента стратегии для непрерывных задач.....	385
13.7. Параметризация стратегии для непрерывных действий .....	388
13.8. Резюме .....	389
Библиографические и исторические замечания.....	390

## **Часть III. ЗАГЛЯНЕМ ПОГЛУБЖЕ.....**

### **Глава 14. Психология .....**

14.1. Предсказание и управление .....	394
14.2. Классическое обусловливание .....	395
14.2.1. Блокирующее обусловливание и обусловливание высшего порядка .....	397
14.2.2. Модель Рескорлы–Вагнера.....	399
14.2.3. TD-модель .....	401
14.2.4. Имитирование TD-модели.....	403
14.3. Инструментальное обусловливание .....	410
14.4. Отложенное подкрепление .....	415
14.5. Когнитивные карты .....	416
14.6. Привычное и целеустремленное поведение .....	418
14.7. Резюме.....	423
Библиографические и исторические замечания.....	425

<b>Глава 15. Нейронауки</b> .....	432
15.1. Основы нейронаук .....	433
15.2. Сигналы вознаграждения, сигналы подкрепления, ценности и ошибки предсказания .....	435
15.3. Гипотеза об ошибке предсказания вознаграждения .....	437
15.4. Дофамин .....	439
15.5. Экспериментальное подтверждение гипотезы об ошибке предсказания вознаграждения .....	443
15.6. Параллель между TD-ошибкой и дофамином .....	447
15.7. Нейронный исполнитель – критик .....	452
15.8. Правила обучения критика и исполнителя .....	456
15.9. Гедонистические нейроны .....	460
15.10. Коллективное обучение с подкреплением .....	462
15.11. Основанные на модели методы в мозге .....	466
15.12. Наркотическая зависимость .....	468
15.13. Резюме .....	469
Библиографические и исторические замечания .....	472
<b>Глава 16. Примеры и приложения</b> .....	481
16.1. TD-Gammon .....	481
16.2. Программы игры в шашки Сэмюэла .....	486
16.3. Стратегия выбора ставки в программе Watson .....	489
16.4. Оптимизация управления памятью .....	492
16.5. Игра в видеоигры на уровне человека .....	497
16.6. Мастерство игры в го .....	503
16.6.1. AlphaGo .....	506
16.6.2. AlphaGo Zero .....	509
16.7. Персонализированные веб-службы .....	513
16.8. Парение в восходящих потоках воздуха .....	516
<b>Глава 17. Передовые рубежи</b> .....	521
17.1. Общие функции ценности и вспомогательные задачи .....	521
17.2. Абстрагирование времени посредством опций .....	523
17.3. Наблюдения и состояния .....	526
17.4. Проектирование сигналов вознаграждения .....	532
17.5. Остающиеся вопросы .....	535
17.6. Экспериментальное подтверждение гипотезы об ошибке предсказания вознаграждения .....	539
Библиографические и исторические замечания .....	543
<b>Предметный указатель</b> .....	587

# Вступительное слово от ГК «Цифра»

Прошло уже несколько лет с тех пор, как наша команда ступила на путь применения искусственного интеллекта для совершенствования процессов в промышленности и логистике. В самом начале мы и представить не могли, насколько тернистой, но в то же время невероятно интересной окажется эта дорога. За это время мы успели поработать с различными производствами и решить множество задач – от оптимизации производства битумных материалов до улучшения системы распределения нефтепродуктов и внедрения систем машинного зрения на карьерные экскаваторы. Методы машинного обучения, которые используются для решения подобных задач, постоянно совершенствуются, и мы внимательно следим за развитием подходов в области искусственного интеллекта, в том числе за исследованиями в обучении с подкреплением.

Обучение с подкреплением – это один из разделов машинного обучения, исследующий вычислительный подход к обучению агента, который пытается максимизировать свою совокупную накопленную награду путем взаимодействия со сложной, зачастую стохастической средой. Последние несколько лет исследования этого подхода переживают настоящий ренессанс – ни одна научная конференция по искусственному интеллекту не обходится без секции на эту тему. Каждый год публикуются сотни научных статей, и все больше компаний в России и за рубежом начинают применять последние достижения этой области в своем бизнесе для улучшения различных внутренних процессов – от рекомендательных систем до оптимизации цепей поставок.

Мы видим огромный потенциал практического применения методов обучения с подкреплением для совершенствования процессов в промышленности и логистике, а также верим в решающее значение данных теоретических концепций и алгоритмов для прогресса искусственного интеллекта как области человеческого знания. Несмотря на огромный интерес к этой области в последнее время, по указанной теме издано не так много литературы. Именно поэтому мы решили поучаствовать в публикации этой замечательной книги на русском языке.

Данная книга представляет собой исчерпывающее введение в такую интересную и быстро развивающуюся область искусственного интеллекта, как обучение с подкреплением. Ее авторы, Ричард Саттон и Эндрю Барто, проделали невероятную работу, описав простым и понятным языком не только ключевые концепции и алгоритмы обучения с подкреплением, но и современные достижения этой области. В книге продемонстрирована связь дисциплины с психологией и нейронауками. Авторами подробно рассматриваются детали работы системы AlphaGo, обыгравшей чемпиона мира в японскую настольную игру го, а также алгоритма, играющего в игры Atari на уровне человека, и многие другие приложения.

Мы желаем читателю удачи на пути изучения такой сложной, но невероятно полезной и увлекательной дисциплины.

*Сергей Свиридов,*

директор по исследованиям и разработкам, группа компаний «Цифра»



Группа компаний «Цифра» разрабатывает технологии цифровизации промышленности, инвестирует в продукты и развивает среду промышленного интернета вещей и искусственного интеллекта. Компания создала самую крупную в России лабораторию промышленного AI. Сегодня решения «Цифры» повышают эффективность промышленных предприятий в 22 странах мира. Ключевые отрасли для группы – это горная добыча и металлургия, машиностроение, нефтегазовый сектор и химическая промышленность. «Цифра входит» в Industrial Internet Consortium и ряд других российских и международных отраслевых ассоциаций.

# Предисловие ко второму изданию

За двадцать лет, прошедших после выхода первого издания этой книги, мы стали свидетелями колоссального прогресса в области искусственного интеллекта, в немалой степени обусловленного достижениями машинного обучения, в т. ч. обучения с подкреплением. И этот прогресс был достигнут не только за счет впечатляющего роста вычислительных мощностей, но и благодаря развитию теории и алгоритмов. Поэтому необходимость во втором издании книги, вышедшей в 1998 году, давно назрела и созрела, и наконец-то в 2012 году мы решили приняться за нее. Во втором издании мы ставили себе ту же цель, что и в первом: дать простое и понятное изложение основных идей и алгоритмов обучения с подкреплением, которое было бы доступно специалистам из смежных дисциплин. Книга по-прежнему осталась введением, основное внимание уделяется базовым алгоритмам онлайн-обучения. Мы включили ряд новых вопросов, возникших и приобретших важность за прошедшие годы, а также расширили описание тем, которые теперь понимаем лучше. Но мы даже не пытались дать исчерпывающее изложение всего предмета, который стремительно развивался во многих направлениях. Приносим извинения за то, что были вынуждены оставить все эти достижения (за исключением небольшого числа) без внимания.

Как и в первом издании, мы решили отказаться от строго формального изложения теории обучения с подкреплением и от постановки задачи в самом общем виде. Но по мере углубления нашего понимания некоторых вопросов потребовалось включить больше математики; части, для которых необходимо более уверенное владение математическим аппаратом, оформлены в виде врезок; читатели, не склонные к математике, могут их пропустить. Мы также используем не совсем такую же нотацию, как в первом издании. В процессе преподавания мы поняли, что новая нотация помогает устранить ряд распространенных недоразумений. Она подчеркивает различие между случайными величинами, которые обозначаются заглавными буквами, и их экземплярами, обозначаемыми строчными буквами. Например, состояние, действие и вознаграждение на временном шаге  $t$  обозначаются  $S_t$ ,  $A_t$  и  $R_t$ , а их возможные значения –  $s$ ,  $a$  и  $r$ . Кроме того, строчными буквами записываются функции ценности (например,  $v_\pi$ ), а заглавными – их табличные представления (например,  $Q_t(s, a)$ ). Приближенные функции ценности являются детерминированными функциями случайных параметров, поэтому также записываются строчными буквами (например,  $\hat{v}(s, \mathbf{w}_t) \approx v_\pi(s)$ ). Векторы, например вектор весов  $\mathbf{w}_t$  (ранее обозначался  $\theta_t$ ) и вектор признаков  $\mathbf{x}_t$  (ранее  $\phi_t$ ), записываются строчными полужирными буквами, даже если являются случайными величинами. Заглавные полужирные буквы оставлены для матриц. В первом издании мы употребляли специальные обозначения  $\mathcal{P}_{SS'}^a$  и  $\mathcal{R}_{SS'}^a$  для вероятности перехода и ожидаемого вознаграждения. Один из недостатков этой нотации заключается в том, что она не полностью характеризует динамику вознаграждения,



а дает только математические ожидания – этого достаточно для динамического программирования, но не для обучения с подкреплением. Другой недостаток – чрезмерное количество верхних и нижних индексов. В этом издании мы ввели явное обозначение  $p(s', r | s, a)$  для совместной вероятности следующего состояния и вознаграждения при условии текущего состояния и действия. Все изменения нотации сведены в таблице на стр. 20.

Второе издание значительно дополнено, и организация материала претерпела изменения. После первой вводной главы появились три новые части. В первой части (главы 2–8) обучение с подкреплением рассматривается настолько полно, насколько возможно без выхода за пределы табличного случая, для которого можно найти точные решения. Мы включили методы обучения и планирования для табличного случая, а также их унификацию в  $n$ -шаговых методах и в архитектуре Дуна. Многих алгоритмов, представленных в этой части, в первом издании не было, например: UCB, Expected Sarsa, двойное обучение, обновление по дереву,  $Q(\sigma)$ , RTDP и MCTS. Подробное рассмотрение табличного случая в начале книги позволяет изложить основные идеи в простейшей постановке. Вторая часть книги (главы 9–13) посвящена обобщению этих идей на аппроксимации функций. В ней появились новые разделы об искусственных нейронных сетях, о базисе Фурье, LSTD, ядерных методах, методах Gradient-TD и Emphatic-TD, методах среднего вознаграждения, истинно онлайн-методе TD( $\lambda$ ) и методах градиента стратегии. Во втором издании намного подробнее рассмотрено обучение с разделенной стратегией, сначала в табличном случае (главы 5–7), а затем для аппроксимации функций в главах 11 и 12. Еще одно отличие второго издания заключается в отделении идеи прямого представления, связанной с  $n$ -шаговым бутстрэппингом (теперь она более полно рассмотрена в главе 7), от идеи обратного представления, связанной со следами приемлемости (она теперь независимо описана в главе 12). В третью часть книги включены новые большие главы о связях обучения с подкреплением с психологией (глава 14) и нейронауками (глава 15), а также переработанная глава с примерами, включающая игры Atari, стратегию ставок в программе Watson, а также две программы игры в го: AlphaGo и AlphaGo Zero (глава 16). Но по необходимости мы смогли включить лишь малую часть сделанного в этой области. Выбор отражает наш давний интерес к недорогим безмодельным методам, которые хорошо масштабируются на крупные приложения. Последняя глава посвящена обсуждению будущего влияния обучения с подкреплением на общество. Хорошо это или плохо, но второе издание получилось почти в два раза больше первого.

Эта книга задумывалась как основной учебник для одно- или двухсеместрового курса по обучению с подкреплением. В односеместровый курс следует включить первые десять глав и излагать их по порядку. Это составит хорошую основу, к которой можно добавить материал из других глав, а также из других книг, например Bertsekas and Tsitsiklis (1996), Wiering and van Otterlo (2012), Szepesvári (2010), или из литературы – сообразуясь со вкусами лектора. В зависимости от подготовки студентов может оказаться полезным дополнительный материал по онлайн-обучению с учителем. Естественным дополнением будут идеи опций и моделей опций (Sutton, Precup and Singh, 1999). В двухсеместровый курс можно включить все главы и дополнительные материалы. Эту книгу можно также включить как часть более широких курсов машинного обучения, искусственного интеллекта или ней-

ронных сетей. В таком случае имеет смысл рассматривать только некоторое подмножество глав. Мы рекомендуем главу 1 в качестве краткого обзора, главу 2 до раздела 2.4, главу 3, а затем избранные разделы остальных глав в зависимости от располагаемого времени и интересов лектора и аудитории. Глава 6 наиболее важна для предмета и всей книги. В курс, ориентированный на машинное обучение или нейронные сети, следует включить главы 9 и 10, а в курс, ориентированный на искусственный интеллект или планирование, – главу 8. Разделы и главы, которые мы считаем более трудными и не существенными для книги в целом, помечены звездочкой. Их можно опустить при первом чтении без ущерба для понимания последующего текста. Упражнения повышенной сложности также помечены звездочкой, они не существенны для усвоения основного материала главы.

Большинство глав заканчиваются разделом «Библиографические и исторические замечания», в которых мы перечисляем источники идей, изложенных в главе, приводим ссылки на литературу для дальнейшего чтения и на текущие исследовательские работы, а также даем историческую справку. Несмотря на все усилия сделать эти разделы полными и авторитетными, мы наверняка упустили какие-то важные работы предшественников. Приносим свои извинения и открыты для исправлений и дополнений, которые будут внесены в электронную версию книги.

Это издание, как и первое, посвящено памяти А. Гарри Клопфа. Именно Гарри познакомил нас друг с другом, и именно его идеи о мозге и искусственном интеллекте побудили нас отправиться в долгое путешествие по миру обучения с подкреплением. Гарри получил образование в области нейрофизиологии и очень интересовался машинным интеллектом, он работал старшим научным сотрудником в отделе авионики Управления научно-исследовательских работ ВВС США (AFOSR) при базе ВВС Райт-Паттерсон в штате Огайо. Он был недоволен тем, что процессам поиска равновесия, в т. ч. гомеостазу и методам классификации на основе исправления ошибок, придают чрезмерно большую важность при объяснении естественного интеллекта и закладывания фундамента машинного интеллекта. Он отмечал, что системы, пытающиеся что-то максимизировать (не важно, что именно), качественно отличаются от систем поиска равновесия, и доказывал, что именно в максимизирующих системах ключ к пониманию важных аспектов естественного интеллекта и построения искусственного. Гарри сыграл решающую роль в получении от AFOSR финансирования для проекта оценки научной ценности этих и родственных им идей. Этот проект был запущен в конце 1970-х годов в Массачусетском университете в Амхерсте (UMass Amherst), сначала под руководством Майкла Эрбиба (Michael Arbib), Уильяма Килмера (William Kilmer) и Нико Спинелли (Nico Spinelli), профессоров факультета компьютерных и информационных наук и членов-основателей университетского кибернетического центра нейронаучных систем, созданного с перспективой работы на стыке нейронаук и искусственного интеллекта. Барто, недавно получивший докторскую степень в Мичиганском университете, был принят в проект на должность младшего научного сотрудника. Тем временем Саттон, студент старшего курса, изучавший информатику и психологию в Стэнфорде, переписывался с Гарри на тему их общего интереса к роли временных характеристик возбудителя в классической теории обусловливания. Гарри убедил группу в UMass в том, что Саттон станет отличным приобретением для проекта. Так Саттон оказался аспирантом

в UMass и начал писать докторскую диссертацию под руководством Барто, который к тому времени занял должность доцента. Исследования обучения с подкреплением, описанные в этой книге, – закономерный итог проекта, начатого Гарри и питавшегося его идеями. Таким образом, Гарри свел нас, авторов книги, положив начало долгой и плодотворной совместной работе. Посвящая эту книгу Гарри, мы отдаем должное его существенному вкладу не только в дисциплину обучения с подкреплением, но и в наше сотрудничество. Мы также выражаем благодарность профессорам Эрбибу, Килмеру и Спинелли за предоставленную нам возможность начать разработку этих идей. Наконец, мы благодарны AFOSR за щедрую поддержку, которую управление оказывало на ранней стадии наших исследований, и Национальному научному фонду (NSF) за щедрое финансирование в течение ряда последующих лет.

Есть много людей, которым мы благодарны за их идеи и помощь в подготовке второго издания. Все, кого мы благодарили за помощь в первом издании, заслуживают нашей глубочайшей благодарности и за это издание тоже – оно бы просто не состоялось без их вклада в первое издание. К этому длинному перечню мы обязаны добавить многих, кто помогал готовить только второе издание. Студенты, которым мы много лет преподавали эту дисциплину, отметились самыми разными способами: находили ошибки, предлагали исправления и – не в последнюю очередь – испытывали затруднения, заставляя нас думать, как объяснить материал лучше. Мы выражаем особую благодарность Марте Стинструп (Martha Steenstrup), которая прочитала весь текст и поделилась подробными комментариями. Главы по психологии и нейронаукам не были бы написаны без помощи многочисленных специалистов в этих областях. Мы признательны Джону Муру (John Moore) за его многолетние терпеливые разъяснения теории и экспериментов по обучению животных и основ нейронауки, а также за внимательное прочтение нескольких черновых вариантов глав 14 и 15. Мы также благодарны Мэтту Ботвинику (Matt Botvinick), Натаниэлю Доу (Nathaniel Daw), Питеру Дайяну (Peter Dayan) и Йелю Ниву (Yael Niv) за пронизательные замечания к черновикам этих глав, помощь в освоении огромного массива литературы и указание на наши многочисленные ошибки в ранних вариантах рукописи. Разумеется, все оставшиеся ошибки в этих главах (а их не может не быть) – целиком наша вина. Мы выражаем благодарность Филу Томасу (Phil Thomas), который помог сделать эти главы доступными неспециалистам в области психологии и нейронаук, и Питеру Стерлингу (Peter Sterling), помогавшему сделать объяснения более понятными. Спасибо также Джиму Хоуку (Jim Houk) за знакомство с вопросами обработки информации в подкорковых ядрах головного мозга и за привлечение нашего внимания к смежным разделам нейронауки. Хоце Мартинес (José Martinez), Терри Сейновски (Terry Sejnowski), Дэвид Силвер (David Silver), Джерри Тезауро (Gerry Tesauro), Георгиос Теочарус (Georgios Theodoros) и Фил Томас (Phil Thomas) любезно помогли нам разобраться в деталях их приложений обучения с подкреплением, чтобы мы могли включить их в главу с примерами. Они же поделились ценными комментариями к черновым вариантам соответствующих разделов. Отдельное спасибо Дэвиду Силверу, который помог нам лучше понять дерево поиска Монте-Карло и программу DeepMind для игры в го. Мы также благодарны Джорджу Конидарису (George Konidaris) за помощь при написании раздела о базисе Фурье. Эмилио Картони (Emilio Cartoni), Томас Седерборг (Thomas Cederborg), Стефан Дерббах

(Stefan Dernbach), Клеменс Розенбаум (Clemens Rosenbaum), Патрик Тэйлор (Patrick Taylor), Томас Колин (Thomas Colin) и Пьер-Люк Бэкон (Pierre-Luc Bacon) помогли нам различными способами, за что мы им очень благодарны.

Саттон также выражает благодарность сотрудникам лаборатории обучения с подкреплением и искусственного интеллекта в университете Альберты за вклад во второе издание. Отдельное спасибо Рупаму Махмуду (Rupam Mahmood) за ценный вклад в обсуждение методов Монте-Карло обучения с разделенной стратегией в главе 5, Хамиду Мэю (Hamid Maei) за помощь в становлении взгляда на обучение с разделенной стратегией, представленного в главе 11, Эрику Грейвсу (Eric Graves) за постановку экспериментов в главе 13, Шан-тон Чжану (Shangton Zhang) за воспроизведение и, как следствие, проверку почти всех экспериментальных результатов, Крису де Асису (Kris De Asis) за улучшение нового технического наполнения глав 7–12 и Харму ван Сейну (Harm van Seijen) за идеи, которые привели к отделению  $n$ -шаговых методов от следов приемлемости и (совместно с Хадом ван Хасселтом [Hado van Hasselt]) – за идеи, касающиеся точной эквивалентности прямого и обратного представления следов приемлемости (глава 12). Саттон также выражает признательность за финансовую поддержку и свободу исследований, которые обеспечивали правительство провинции Альберты и Национальный совет научных и инженерных исследований Канады на протяжении всей работы над вторым изданием книги. В частности, он благодарен Рэнди Гебелю (Randy Goebel) за создание благоприятной среды для исследований в Альберте с прицелом на перспективу. Также он благодарен компании DeepMind за поддержку на протяжении последних шести месяцев работы над книгой.

Наконец, мы признательны многочисленным придирчивым читателям черновых вариантов второго издания, которые мы выкладывали в интернет. Они нашли немало пропущенных нами ошибок и указали места, где может возникнуть недопонимание.

# Предисловие к первому изданию

То, что теперь называется обучением с подкреплением, впервые привлекло наше внимание в конце 1979 года. Мы оба работали в Массачусетском университете над одним из ранних проектов воскрешения идеи сетей с нейроноподобными адаптивными элементами, которая могла оказаться многообещающим подходом к искусственному адаптивному интеллекту. Проект был посвящен исследованию «гетеростатической теории адаптивных систем» и разрабатывался под руководством А. Гарри Клопфа. Работа Гарри была богатейшим источником идей, а нам было позволено критически изучить их и сравнить с долгой историей предшествующих исследований в области адаптивных систем. Нашей задачей стало расчленение этих идей на составные части в попытке понять их взаимосвязи и сравнительную важность. Это продолжается и по сей день, но в 1979 году мы впервые осознали, что самая простая идея, которую долго считали чем-то само собой разумеющимся, удостоилась на удивление скромного внимания с вычислительной точки зрения. Это была идея обучающейся системы, которая хочет чего-то достичь и для этого адаптирует свое поведение так, чтобы максимизировать специальный сигнал со стороны окружающей среды. Иначе говоря, идея «гедонистической» обучающейся системы, или, как мы сказали бы теперь, идея обучения с подкреплением.

Как и многие другие, мы полагали, что обучение с подкреплением было всесторонне исследовано еще на заре развития кибернетики и искусственного интеллекта. Но при ближайшем рассмотрении оказалось, что его изучали очень поверхностно. Хотя обучение с подкреплением, безусловно, стало побудительным мотивом для некоторых ранних компьютерных исследований обучения, большая часть занимавшихся этим ученых затем обратилась к другим вещам: классификации образов, обучению с учителем или адаптивному управлению, а то и вовсе забросили исследования в области обучения. В результате специальным вопросам, связанным с тем, как обучиться получать что-нибудь от среды, было уделено сравнительно мало внимания. Оглядываясь назад, можно сказать, что интерес к этой идее стал важнейшим шагом, приведшим в движение всю эту ветвь исследований. Мало чего можно было бы достичь в плане вычислительного обучения с подкреплением, не осознав, что столь фундаментальная идея ранее не была досконально исследована.

С тех пор эта область науки прошла долгий путь, развивалась в нескольких направлениях и стала зрелой дисциплиной. Обучение с подкреплением постепенно стало одним из самых активных направлений исследований в машинном обучении, искусственном интеллекте и нейронных сетях. Было подведено солидное математическое основание и созданы впечатляющие приложения. Компьютерные исследования обучения с подкреплением превратились в обширную область, в которой трудятся сотни ученых по всему миру, занимающиеся такими разными дисциплинами, как психология, теория управления, искусственный интеллект

и нейронауки. Особенно важны результаты, устанавливающие и развивающие связи с теорией оптимального управления и динамическим программированием. В целом проблема обучения путем взаимодействия ради достижения поставленных целей еще далека от решения, но наше понимание стало значительно глубже. Мы теперь можем изучать отдельные направления, например обучение на основе временных различий, динамическое программирование и аппроксимацию функций, в контексте их вклада в решение общей проблемы.

Принимаясь за написание книги, мы ставили цель дать простое и ясное описание ключевых идей и алгоритмов обучения с подкреплением. Мы хотели, чтобы изложение было доступно читателям, работающим в смежных дисциплинах, но не могли охватить их все одинаково подробно. В основном мы ведем изложение с точки зрения искусственного интеллекта и технического конструирования. Рассмотрение связей с иными областями мы оставляем другим авторам или отложим до следующего раза. Мы также решили отказаться от строго формального изложения предмета. Мы не стремились к максимальному уровню математического абстрагирования и не пытались доказывать теоремы. Мы постарались выбрать такой уровень математических деталей, который указал бы читателям с математическим складом ума верное направление, но не отвлекал от простоты и потенциальной общности базовых идей.

В некотором смысле мы работали над этой книгой тридцать лет, так что нам есть кого благодарить. Прежде всего мы благодарим людей, лично помогавших нам в разработке идей, представленных в книге: Гарри Клопфа (Harri Klopff), который помог нам осознать, что обучение с подкреплением нуждается в новой жизни; Криса Уоткинса (Chris Watkins), Димитрия Бертсекаса (Dimitri Bertsekas), Джона Цициклиса (John Tsitsiklis) и Пода Вербоса (Paul Werbos), которые помогли нам понять ценность связей с динамическим программированием; Джона Мура (John Moore) и Джима Кехоу (Jim Kehoe) за идеи из теории обучения животных; Оливера Селфриджа (Oliver Selfridge) за подчеркивание широты и важности адаптации; и вообще наших коллег и студентов, помогавших самыми разными способами: Рона Уильямса (Ron Williams), Чарльза Андерсона (Charles Anderson), Сатиндера Сингха (Satinder Singh), Сридхара Махадевана (Sridhar Mahadevan), Стива Брадтке (Steve Bradtke), Боба Крайтца (Bob Crites), Питера Дайяна (Peter Dayan) и Лимона Бэрда (Leemon Baird). На наши воззрения на обучение с подкреплением оказали большое влияние беседы с Полом Коэном (Paul Cohen), Полом Утгоффом (Paul Utgoff), Мартой Стинструп (Martha Steenstrup), Джерри Тезауро (Gerry Tesauro), Майком Джорданом (Mike Jordan), Лесли Кэлблингом (Leslie Kaelbling), Эндрю Муром (Andrew Moore), Крисом Аткисоном (Chris Atkeson), Томом Митчеллом (Tom Mitchell), Нильсом Нильсоном (Nils Nilsson), Стюартом Расселом (Stuart Russell), Томом Диттерихом (Tom Dietterich), Томом Дином (Tom Dean) и Бобом Нарендра (Bob Narendra).

Мы благодарны Майклу Литтману, Джерри Тезауро, Майклу Крайтцу, Сатиндеру Сингху и Вэй Чжану (Wei Zhang) за наполнение конкретикой разделов 4.7, 15.1, 15.4, 15.5 и 15.6 соответственно. Мы благодарим Управление научно-исследовательских работ ВВС США, Национальный научный фонд и лаборатории GTE за длительную финансовую поддержку, нацеленную на перспективу. Мы также выражаем признательность многим людям, которые читали черновые варианты книги и делились ценными замечаниями: Тому Калту (Tom Kalt), Джону Цициклису,

Павлу Чихошу (Pawel Cichosz), Олле Гэллмо (Olle Gällmo), Чаку Андерсону (Chuck Anderson), Стюарту Расселу, Бену ван Рою (Ben Van Roy), Полу Стинструпу (Paul Steenstrup), Полу Коэну, Сридхару Махадевану, Джетте Рандлов (Jette Randlov), Брайану Шеппарду (Brian Sheppard), Томасу О'Коннелу (Thomas O'Connell), Ричарду Коггинсу (Richard Coggins), Кристине Версино (Cristina Versino), Джону Х. Хайетту (John H. Hiatt), Андреасу Баделту (Andreas Badelt), Джею Понте (Jay Ponte), Джо Беку (Joe Beck), Юстусу Пиатеру (Justus Piater), Марту Стинструп, Сатиндеру Сингху, Томми Яаколла (Jaakkola), Димитрию Бертсекасу (Dimitri Bertsekas), Торбьёрну Экману (Torbjörn Ekman), Кристине Бьёркман (Christina Björkman), Якобу Карлстрёму (Jakob Carlström) и Олле Палмгрену (Olle Palmgren). Наконец, мы благодарим Гвин Митчелл (Gwyn Mitchell) за разнообразную помощь, а также Гарри Стэнтона (Harry Stanton) и Боба Приора (Bob Prior), которые опекали нас в издательстве MIT Press.

# Обозначения

Заглавными буквами обозначаются случайные величины, строчными – значения случайных величин и скалярные функции. Вещественные векторы записываются строчными полужирными буквами (даже если они являются случайными величинами), матрицы – заглавными полужирными буквами.

$\doteq$	равенство, имеющее место по определению
$\approx$	приближенное равенство
$\propto$	пропорционально
$\Pr\{X = x\}$	вероятность, что случайная величина $X$ принимает значение $x$
$X \sim p$	случайная величина $X$ выбрана из распределения $p(x) \doteq \Pr\{X = x\}$
$\mathbb{E}[X]$	математическое ожидание случайной величины $X$ , т. е. $\mathbb{E}[X] \doteq \sum_x p(x)x$
$\operatorname{argmax}_a f(a)$	значение $a$ , в котором $f(a)$ достигает максимума
$\ln x$	натуральный логарифм $x$
$e^x$	основание натуральных логарифмов, число $e \approx 2.71828$ , возведенное в степень $x$ ; $e^{\ln x} = x$
$\mathbb{R}$	множество вещественных чисел
$f: X \rightarrow Y$	функция $f$ , отображающая элементы множества $X$ в элементы множества $Y$
$\leftarrow$	присваивание
$(a, b]$	интервал вещественной оси между $a$ и $b$ , включающий $b$ , но не включающий $a$
$\varepsilon$	вероятность предпринять случайное действие в $\varepsilon$ -жадной стратегии
$\alpha, \beta$	параметры, определяющие размер шага
$\gamma$	коэффициент обесценивания
$\lambda$	коэффициент затухания для следов приемлемости
$\mathbb{1}_{predicate}$	индикаторная функция ( $\mathbb{1}_{predicate} \doteq 1$ , если предикат $predicate$ равен true, в противном случае 0)

В задаче о многоруких бандитах:

$k$	количество действий (рук)
$t$	дискретный временной шаг или номер игры
$q_*(a)$	истинное значение (ожидаемое вознаграждение) действия $a$
$Q_t(a)$	оценка $q_*(a)$ в момент $t$
$N_t(a)$	сколько раз действие $a$ выбиралось до момента $t$
$H_t(a)$	обученное предпочтение действию $a$ в момент $t$
$\pi_t(a)$	вероятность выбора действия $a$ в момент $t$
$\bar{R}_t$	оценка ожидаемого вознаграждения в момент $t$ при условии $\pi_t$



В марковском процессе принятия решений:

$s, s'$	состояния
$a$	действие
$r$	вознаграждение
$\mathcal{S}$	множество всех незаключительных состояний
$\mathcal{S}^+$	множество всех состояний, включая заключительное
$\mathcal{A}(s)$	множество всех действий, допустимых в состоянии $s$
$\mathcal{R}$	множество всех возможных вознаграждений, конечное подмножество $\mathbb{R}$
$\subset$	подмножество (например, $\mathcal{R} \subset \mathbb{R}$ )
$\in$	элемент множества, например $s \in \mathcal{S}, r \in \mathcal{R}$
$ \mathcal{S} $	количество элементов во множестве $\mathcal{S}$ (мощность $\mathcal{S}$ )
$t$	дискретный временной шаг
$T, T(t)$	конечный временной шаг эпизода или эпизод, включающий временной шаг $t$
$A_t$	действие в момент $t$
$S_t$	состояние в момент $t$ , обычно стохастически зависящее от $S_{t-1}$ и $A_{t-1}$
$R_t$	вознаграждение в момент $t$ , обычно стохастически зависящее от $S_{t-1}$ и $A_{t-1}$
$\pi$	стратегия (правило принятия решения)
$\pi(s)$	действие, предпринятое в состоянии $s$ при <i>детерминированной</i> стратегии $\pi$
$\pi(a s)$	вероятность предпринять действие $a$ в состоянии $s$ при <i>стохастической</i> стратегии $\pi$
$G_t$	доход, начиная с момента $t$
$h$	горизонт, на какое время можно заглянуть вперед в прямом представлении
$G_{t:t+n}, G_{t,h}$	доход за $n$ шагов с $t + 1$ до $t + n$ или до $h$ (обесцененный и скорректированный)
$\bar{G}_{t,h}$	плоский доход (необесцененный и нескорректированный) за шаги от $t + 1$ до $h$ (раздел 5.8)
$G_t^\lambda$	$\lambda$ -доход (раздел 12.1)
$G_{t,h}^\lambda$	усеченный скорректированный $\lambda$ -доход (раздел 12.3)
$G_t^{\lambda s}, G_t^{\lambda a}$	$\lambda$ -доход, скорректированный на оценки ценности состояния или действия (раздел 12.8)
$p(s', r s, a)$	вероятность перехода в состояние $s'$ с вознаграждением $r$ из состояния $s$ после действия $a$
$p(s' s, a)$	вероятность перехода в состояние $s'$ из состояния $s$ после действия $a$
$r(s, a)$	ожидаемое немедленное вознаграждение в состоянии $s$ после действия $a$

$r(s, a, s')$	ожидаемое немедленное вознаграждение при переходе из $s$ в $s'$ после действия $a$
$v_\pi(s)$	ценность состояния $s$ при стратегии $\pi$ (ожидаемый доход)
$v_*(s)$	ценность состояния $s$ при оптимальной стратегии
$q_\pi(s, a)$	ценность выполнения действия $a$ в состоянии $s$ при стратегии $\pi$
$q_*(s, a)$	ценность выполнения действия $a$ в состоянии $s$ при оптимальной стратегии
$V, V_t$	массив оценок функции ценности состояний $v_\pi$ или $v_*$
$Q, Q_t$	массив оценок функции ценности действий $q_\pi$ или $q_*$
$\bar{V}_t(s)$	ожидаемая приближенная ценность действия, например $V_t(s) \doteq \sum_a \pi(a s) Q_t(s, a)$
$U_t$	цель для оценки в момент $t$
$\delta_t$	ошибка временного различия (TD-ошибка) в момент $t$ (случайная величина) (раздел 6.1)
$\delta_t^s, \delta_t^a$	формы TD-ошибки для состояния и действия (раздел 12.9)
$n$	в $n$ -шаговых методах $n$ – количество шагов бутстрэппинга
$\ v\ _\mu^2$	$\mu$ -взвешенная квадратичная норма функции ценности, $\ v\ _\mu^2 \doteq \sum_{s \in S} \mu(s) v(s)^2$
$d$	размерность – количество элементов $\mathbf{w}$
$d'$	альтернативная размерность – количество элементов $\theta$
$\mathbf{w}, \mathbf{w}_t$	$d$ -мерный вектор весов, определяющий приближенную функцию ценности
$w_i, w_{t,i}$	$i$ -й элемент обучаемого вектора весов
$\hat{v}(s, \mathbf{w})$	приближенная ценность состояния $s$ при условии вектора весов $\mathbf{w}$
$v_{\mathbf{w}}(s)$	альтернативное обозначение $\hat{v}(s, \mathbf{w})$
$\hat{q}(s, a, \mathbf{w})$	приближенная ценность пары состояние–действий $s, a$ при заданном векторе весов $\mathbf{w}$
$\nabla \hat{v}(s, \mathbf{w})$	вектор-столбец частных производных $\hat{v}(s, \mathbf{w})$ по $\mathbf{w}$
$\nabla \hat{q}(s, a, \mathbf{w})$	вектор-столбец частных производных $\hat{q}(s, a, \mathbf{w})$ по $\mathbf{w}$
$\mathbf{x}(s)$	вектор признаков, видимых в состоянии $s$
$\mathbf{x}(s, a)$	вектор признаков, видимых, когда в состоянии $s$ предпринимается действие $a$
$x_i(s), x_i(s, a)$	$i$ -й элемент вектора $\mathbf{x}(s)$ или $\mathbf{x}(s, a)$
$\mathbf{x}_t$	сокращенное обозначение $\mathbf{x}(S_t)$ или $\mathbf{x}(S_t, A_t)$
$\mathbf{w}^\top \mathbf{x}$	скалярное произведение векторов, $\mathbf{w}^\top \mathbf{x} \doteq \sum_i w_i x_i$ ; например, $\hat{v}(s, \mathbf{w}) \doteq \mathbf{w}^\top \mathbf{x}(s)$

$\mathbf{v}, \mathbf{v}_t$	вторичный $d$ -мерный вектор весов, используемый для обучения $\mathbf{w}$ (глава 11)
$\mathbf{z}_t$	$d$ -мерный вектор следов приемлемости в момент $t$ (глава 12)
$\theta, \theta_t$	вектор параметров целевой стратегии (глава 13)
$\pi(a s, \theta)$	вероятность выбора действия $a$ в состоянии $s$ при условии параметрического вектора $\theta$
$\pi_\theta$	стратегия, соответствующая параметрам $\theta$
$\nabla\pi(a s, \theta)$	вектор-столбец частных производных $\pi(a s, \theta)$ по $\theta$
$J(\theta)$	мера качества для стратегии $\pi_\theta$
$\nabla J(\theta)$	вектор-столбец частных производных $J(\theta)$ по $\theta$
$h(s, a, \theta)$	предпочтение выбору действия $a$ в состоянии $s$ , основанное на $\theta$
$b(a s)$	поведенческая стратегия, применяемая для выбора действий в процессе обучения целевой стратегии $\pi$
$b(s)$	базовая функция $b: \mathcal{S} \mapsto \mathbb{R}$ для методов градиента стратегии
$b$	коэффициент ветвления для МППР или дерева поиска
$\rho_{t:h}$	коэффициент выборки по значимости для временных шагов от $t$ до $h$ (раздел 5.5)
$\rho_t$	коэффициент выборки по значимости для одного только шага $t$ , $\rho_t = \rho_{t:t}$
$r(\pi)$	среднее вознаграждение (коэффициент вознаграждения) для стратегии $\pi$ (раздел 10.3)
$\bar{R}_t$	оценка $r(\pi)$ в момент $t$
$\mu(s)$	распределение состояний с единой стратегией (раздел 9.2)
$\mu$	$ \mathcal{S} $ -мерный вектор $\mu(s)$ для всех $s \in \mathcal{S}$
$\ v\ _\mu^2$	$\mu$ -взвешенная норма функции ценности $v$ , т. е. $\ v\ _\mu^2 \doteq \sum_s \mu(s)v(s)^2$ (раздел 11.4)
$\eta(s)$	ожидаемое количество посещений состояния $s$ в одном эпизоде (стр. 240)
$\Pi$	оператор проекции для функций ценности (стр. 316)
$B_\pi$	оператор Беллмана для функций ценности (раздел 11.4)
$\mathbf{A}$	матрица $\mathbf{A} \doteq \mathbb{E}[\mathbf{x}_t(\mathbf{x}_t - \gamma\mathbf{x}_{t+1})^\top]$ размерности $d \times d$
$\mathbf{b}$	$d$ -мерный вектор $\mathbf{b} \doteq \mathbb{E}[R_{t+1}\mathbf{x}_t]$
$\mathbf{w}_{\text{TD}}$	неподвижная точка TD $\mathbf{w}_{\text{TD}} \doteq \mathbf{A}^{-1}\mathbf{b}$ ( $d$ -мерный вектор, раздел 9.4)
$\mathbf{I}$	единичная матрица
$\mathbf{P}$	матрица $ \mathcal{S}  \times  \mathcal{S} $ вероятностей перехода состояний при стратегии $\pi$
$\mathbf{D}$	диагональная матрица $ \mathcal{S}  \times  \mathcal{S} $ со значениями $\mu$ на диагонали
$\mathbf{X}$	матрица $ \mathcal{S}  \times d$ с векторами-строками $\mathbf{x}(s)$
$\overline{\text{VE}}(\mathbf{w})$	среднеквадратическая ошибка $\overline{\text{VE}}(\mathbf{w}) \doteq \ \mathbf{v}_\mathbf{w} - \mathbf{v}_\pi\ _\mu^2$ (раздел 9.2)
$\bar{\delta}_\mathbf{w}(s)$	беллмановская ошибка (математическое ожидание TD-ошибки), когда состоянием $s$ является $\mathbf{v}_\mathbf{w}$ (раздел 11.4)

Конец ознакомительного фрагмента.

Приобрести книгу можно

в интернет-магазине

«Электронный универс»

[e-Univers.ru](http://e-Univers.ru)