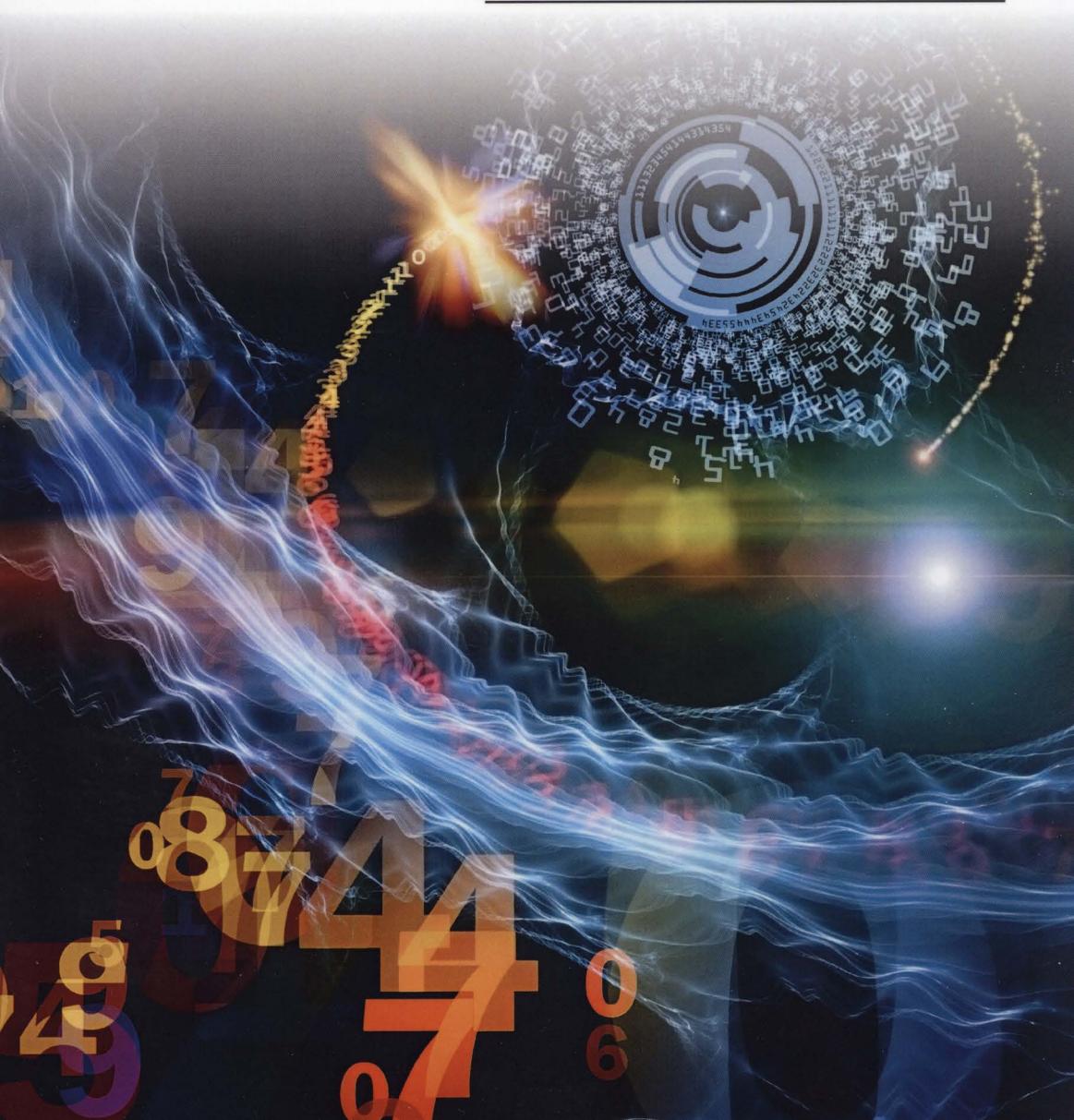


ВВЕДЕНИЕ В ГЛУБОКОЕ ОБУЧЕНИЕ

ЕВГЕНИЙ ЧЕРНЯК



ВВЕДЕНИЕ
в ГЛУБОКОЕ ОБУЧЕНИЕ

INTRODUCTION TO DEEP LEARNING

EUGENE CHARNIAK

The MIT Press
Cambridge, Massachusetts
London, England

ВВЕДЕНИЕ В ГЛУБОКОЕ ОБУЧЕНИЕ

ЕВГЕНИЙ ЧЕРНЯК



Москва · Санкт-Петербург
2020

ББК 32.973.26-018.2.75

Ч-49

УДК 681.3.07

ООО “Диалектика”

Зав. редакцией *С.Н. Тригуб*

Перевод с английского и редакция *В.А. Коваленко*

По общим вопросам обращайтесь в издательство “Диалектика” по адресу:

info@dialektika.com, <http://www.dialektika.com>

Черняк, Евгений.

Ч-49 Введение в глубокое обучение. : Пер. с англ. — СПб. : ООО “Диалектика”, 2020. — 192 с. : ил. — Парал. тит. англ.

ISBN 978-5-907203-10-5 (рус.)

ББК 32.973.26-018.2.75

Все названия программных продуктов являются зарегистрированными торговыми марками соответствующих фирм.

Никакая часть настоящего издания ни в каких целях не может быть воспроизведена в какой бы то ни было форме и какими бы то ни было средствами, будь то электронные или механические, включая фотокопирование и запись на магнитный носитель, если на это нет письменного разрешения издательства MIT Press.

Copyright © 2020 by Dialektika Computer Publishing.

Authorized translation from the English language edition *Introduction to Deep Learning* (ISBN 978-0-262-03951-2) published by MIT Press, Copyright © 2018 by Massachusetts Institute of Technology.

All rights reserved. No part of this book may be reproduced in any form or by any electronic or mechanical means (including photocopying, recording, or information storage and retrieval) without permission in writing from the publisher.

Научно-популярное издание

Евгений Черняк

Введение в глубокое обучение

ООО “Диалектика”, 195027, Санкт-Петербург, Магнитогорская ул., д. 30, лит. А, пом. 848

ISBN 978-5-907203-10-5 (рус.)

© ООО “Диалектика”, 2020,
перевод, оформление, макетирование

ISBN 978-0-262-03951-2 (англ.)

© 2018 The Massachusetts Institute of Technology

Оглавление

Введение	9
Глава 1. Нейронные сети с прямой связью	11
Глава 2. Язык Tensorflow	39
Глава 3. Сверточные нейронные сети	61
Глава 4. Векторное представление слов и рекуррентные NN	79
Глава 5. Обучение от последовательности к последовательности	103
Глава 6. Глубокое обучение с подкреплением	121
Глава 7. Модели нейронных сетей, обучаемых без учителя	145
Приложение. Выборочные ответы к упражнениям	169
Предметный указатель	175

Содержание

Введение	9
Посвящение	10
Ждем ваших отзывов!	10
Глава 1. Нейронные сети с прямой связью	11
1.1. Перцептроны	13
1.2. Функция кросс-энтропийных потерь для нейронных сетей	19
1.3. Производные и стохастический градиентный спуск	23
1.4. Написание программы	28
1.5. Матричное представление нейронных сетей	30
1.6. Независимость данных	33
1.7. Ссылки и что читать дальше	35
1.8. Упражнения	36
Глава 2. Язык Tensorflow	39
2.1. Вводная информация о Tensorflow	39
2.2. Программа TF	43
2.3. Многослойные NN	48
2.4. Другие части	51
2.4.1. Создание контрольных точек	51
2.4.2. <code>tensordot</code>	53
2.4.3. Инициализация переменных TF	55
2.4.4. Упрощение создания графов TF	57
2.5. Ссылки и что читать дальше	58
2.6. Упражнения	58
Глава 3. Сверточные нейронные сети	61
3.1. Фильтры, шаги и дополнения	62
3.2. Простой пример свертки TF	68
3.3. Многослойная свертка	70
3.4. Детали свертки	74
3.4.1. Смещения	74
3.4.2. Слои со сверткой	74
3.4.3. Объединение	75

3.5. Ссылки и что читать дальше	76
3.6. Упражнения	77
Глава 4. Векторное представление слов и рекуррентные NN	79
4.1. Векторное представление слова для языковых моделей	79
4.2. Создание языковых моделей с прямой связью	84
4.3. Улучшение языковых моделей с прямой связью	86
4.4. Переобучение	87
4.5. Рекуррентные сети	90
4.6. Долгая краткосрочная память	97
4.7. Ссылки и что читать дальше	100
4.8. Упражнения	101
Глава 5. Обучение от последовательности к последовательности	103
5.1. Парадигма Seq2Seq	104
5.2. Написание программы Seq2Seq MT	107
5.3. Внимание в Seq2seq	110
5.4. Seq2Seq с несколькими длинами окон	114
5.5. Упражнение по программированию	115
5.6. Упражнения	118
5.7. Ссылки и что читать дальше	118
Глава 6. Глубокое обучение с подкреплением	121
6.1. Итерация по значениям	122
6.2. Q-обучение	125
6.3. Основы глубокого Q-обучения	128
6.4. Методы градиента политики	131
6.5. Методы актор-критик	138
6.6. Воспроизведение опыта	140
6.7. Ссылки и что читать дальше	142
6.8. Упражнения	142
Глава 7. Модели нейронных сетей, обучаемых без учителя	145
7.1. Основы автокодировки	145
7.2. Сверточное автокодирование	149
7.3. Вариационное автокодирование	153
7.4. Генеративно-состязательные сети	161

7.5. Ссылки и что читать дальше	166
7.6. Упражнения	166
Приложение. Ответы к некоторым упражнениям	169
Глава 1	169
Глава 2	170
Глава 3	170
Глава 4	171
Глава 5	171
Глава 6	172
Глава 7	173
Предметный указатель	175

Введение

Ваш автор — давний исследователь искусственного интеллекта, специализирующийся на обработке естественного языка, революцию в котором сделало глубокое обучение. К сожалению, ему (мне) потребовалось много времени, чтобы это понять. Могу сказать в свое оправдание, что это уже третий раз, когда нейронные сети угрожают революцией, а отнюдь не первый. Тем не менее я внезапно оказался далеко позади и изо всех сил пытался наверстать упущенное. Именно поэтому я сделал то, что сделал бы на моем месте любой уважающий себя профессор: запланировал преподавание материала и начал ускоренный курс, просматривая веб-страницы, а также заставил своих студентов преподавать мне. (Последнее — не шутка. В частности, заслуживает особого упоминания старший преподаватель-ассистент по курсу, Сиддхарта (Сидд) Каррамхети (Siddharth (Sidd) Karramcheti).)

Этим объясняется несколько выдающихся особенностей этой книги. Во-первых, краткость. Я учусь медленно. Во-вторых, она сильно зависит от проекта. Многие публикации, особенно в области информатики, постоянно имеют противоречия между организацией темы и организацией материалов, связанных с конкретными проектами. Подобное разделение зачастую является хорошей идеей, но я считаю, что материал по информатике лучше изучать при написании программ, поэтому моя книга во многом отражает мои привычки в преподавании. Таков был самый удобный способ написания книги, и я надеюсь, что многие из читателей тоже найдут ее полезной.

Возникает вопрос об ожидаемой аудитории. Хотя я надеюсь, что многие практикующие CS сочтут эту книгу полезной по той же причине, по которой я ее написал, как преподаватель я, в первую очередь, верю своим ученикам, поэтому данная книга задумана в качестве учебника для курса по глубокому обучению. Курс, который я преподаю в Брауне (Brown), предназначен как для выпускников, так и для других студентов, и охватывает весь материал, а также некоторые лекции по “культуре” (для получения диплома студент должен выполнить серьезный итоговый проект). Требуются как линейная алгебра, так и многомерное исчисление. Хотя фактическое количество материала по линейной алгебре не так уж велико, студенты сказали мне, что без него им было бы довольно сложно разобраться в многослойных сетях и необходимых им тензорах. Тем не менее многовариантное исчисление было им гораздо понятней. Это явно проявляется только в главе 1, когда обратное распространение создается “с нуля”, и я не удивлюсь, если окажется полезной дополнительная лекция по

частным производным. И наконец, есть предпосылка для вероятности и статистики. Это упрощает диспозицию, и я, конечно же, хочу побудить студентов пройти такой курс. Я также предполагаю элементарные знания программирования на языке Python. Я не включаю это в текст, но у моего курса есть дополнительная “лаборатория” по основам языка Python.

То, что ваш автор играл в догонялки при написании этой книги, также объясняется тем фактом, что почти в каждой главе вы найдете раздел о том, что читать далее, помимо обычных ссылок на важные исследовательские работы и множества ссылок на вторичные источники — чужие учебные труды. Я никогда не узнал бы об этом материале без них.

Провиденс, Род-Айленд
Январь 2018 года

Посвящение

Моей семье, еще раз

Ждем ваших отзывов!

Вы, читатель этой книги, и есть главный ее критик. Мы ценим ваше мнение и хотим знать, что было сделано нами правильно, что можно было сделать лучше и что еще вы хотели бы увидеть изданным нами. Нам интересны любые ваши замечания в наш адрес.

Мы ждем ваших комментариев и надеемся на них. Вы можете прислать нам электронное письмо либо просто посетить наш веб-сайт и оставить свои замечания там. Одним словом, любым удобным для вас способом дайте нам знать, нравится ли вам эта книга, а также выскажите свое мнение о том, как сделать наши книги более интересными для вас.

Отправляя письмо или сообщение, не забудьте указать название книги и ее авторов, а также свой обратный адрес. Мы внимательно ознакомимся с вашим мнением и обязательно учтем его при отборе и подготовке к изданию новых книг.

Актуальность ссылок не гарантируется.

Наши электронные адреса:

E-mail: info@dialektika.com

WWW: <http://www.dialektika.com>

Глава 1

Нейронные сети с прямой связью

Изучение *глубокого обучения* (deep learning) (или *нейронных сетей* (neural net) — мы употребляем эти термины как синонимы) обычно принято начинать с их использования в компьютерном зрении. Эта область искусственного интеллекта претерпела техническую революцию, и ее базовая отправная точка — *интенсивность света* (light intensity) — вполне естественно представлена действительными числами, которыми манипулируют нейронные сети.

Рассмотрим конкретный пример — проблему идентификации рукописных цифр (от нуля до девяти). Если начать с нуля, то сначала нужно разработать камеру для фокусировки световых лучей, чтобы создать изображение того, что мы видим. Тогда нам потребуются датчики света, чтобы превратить световые лучи в электрические импульсы, которые компьютер может “ощущать”. И, наконец, поскольку мы имеем дело с цифровыми компьютерами, нужно *оцифровать* (discretize) изображение, т.е. представить цвета и интенсивности света в виде чисел в двумерном массиве. К счастью, у нас есть набор данных, в котором все это уже сделано, — набор данных Mnist (произносится как “эм нист”). (“Нист” (nist) здесь происходит от *Национального института стандартов* (National Institute of Standards или NIST) США, который отвечал за сбор данных.) В этом наборе данных каждое изображение представляет собой массив целых чисел размером 28×28 , как на рис. 1.1. (Я удалил левую и правую границы области, чтобы лучше разместить их на странице.)

На рис. 1.1 можно считать 0 белым цветом, 255 — черным, а числа между ними — оттенками серого. Мы называем эти числа *значениями пикселей* (pixel value), где *пиксель* (pixel) — это наименьшая часть изображения, которую может разобрать наш компьютер. Фактический “размер” области мира, представленной пикселям, зависит от камеры, от того, насколько далеко она находится от поверхности объекта, и т.д. Но для нашей простой задачи с цифрами не нужно об этом беспокоиться. Черно-белое изображение показано на рис. 1.2.

	7	8	9	10	11	12	13	14	15	16	17	18	19	20
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	185	159	151	60	36	0	0	0	0	0	0	0	0	0
8	254	254	254	254	241	198	198	198	198	198	198	198	198	170
9	114	72	114	163	227	254	225	254	254	254	250	229	254	254
10	0	0	0	0	17	66	14	67	67	67	59	21	236	254
11	0	0	0	0	0	0	0	0	0	0	0	83	253	209
12	0	0	0	0	0	0	0	0	0	0	22	233	255	83
13	0	0	0	0	0	0	0	0	0	0	129	254	238	44
14	0	0	0	0	0	0	0	0	0	59	249	254	62	0
15	0	0	0	0	0	0	0	0	0	133	254	187	5	0
16	0	0	0	0	0	0	0	0	9	205	248	58	0	0
17	0	0	0	0	0	0	0	0	126	254	182	0	0	0
18	0	0	0	0	0	0	0	75	251	240	57	0	0	0
19	0	0	0	0	0	0	19	221	254	166	0	0	0	0
20	0	0	0	0	0	3	203	254	219	35	0	0	0	0
21	0	0	0	0	0	38	254	254	77	0	0	0	0	0
22	0	0	0	0	31	224	254	115	1	0	0	0	0	0
23	0	0	0	0	133	254	254	52	0	0	0	0	0	0
24	0	0	0	61	242	254	254	52	0	0	0	0	0	0
25	0	0	0	121	254	254	219	40	0	0	0	0	0	0
26	0	0	0	121	254	207	18	0	0	0	0	0	0	0
27	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Рис. 1.1. Набор Mnist: цифровая версия изображения

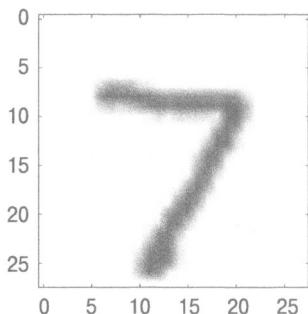


Рис. 1.2. Чёрно-белое изображение из пикселей на рис. 1.1

Внимательно посмотрев на это изображение, можно предложить некоторые простые способы выполнение нашей задачи. Например, обратите внимание, что пиксель в позиции [8, 8] темный. Учитывая, что это образ цифры “7”, это вполне резонно. Точно так семерки зачастую имеют светлый участок посередине, т.е. пиксель [13,13] имеет нуль в качестве значения интенсивности. Сравните это с цифрой “1”, которая зачастую имеет противоположные значения для этих двух позиций, поскольку стандартный рисунок числа не занимает верхний левый угол, но заполняет точно середину. Немного подумав, мы могли бы придумать много *эвристик* (heuristic) (правил, которые работают часто, но не всегда), таких как эти, а затем, используя их, написать программу классификации.

Но это не то, что мы собираемся делать, поскольку в этой книге мы сосредоточиваемся на *машином обучении* (machine learning), т.е. подходим к задачам, спрашивая себя, как мы можем позволить компьютеру учиться, приводя примеры с правильным ответом. В данном случае мы хотим, чтобы наша программа научилась идентифицировать изображения цифр размером 28×28 , приводя их примеры вместе с ответами (называемыми также *метками* (label)). В машинном обучении мы сказали бы, что это задача *обучения с учителем* (supervised learning) или, если быть совершенно точным, задача *полного обучения с учителем* (fully supervised learning), поскольку для каждого учебного примера мы также даем компьютеру правильный ответ. В последующих главах, например в главе 6, у нас нет такой роскоши. Там у нас будет задача *обучения с частичным привлечением учителя* (semi-supervised) или даже *обучение без учителя* (unsupervised learning) (в главе 7). В этих главах мы увидим, как это может работать.

Абстрагировав детали работы с миром световых лучей и поверхностей, мы остаемся с *задачей классификации* (classification problem), когда дан набор входных данных (*признаков* (feature)), идентифицирующих (или *классифицирующих* (classify)) сущность, которая породила эти входные данные (или имеет эти признаки) в качестве одной из конечного числа альтернатив. В нашем случае входные данные представляют собой пиксели, а классификация состоит из десяти возможных вариантов. Обозначим вектор из l входных данных (пикселей) как $\mathbf{x} = [x_1, x_2 \dots x_l]$ и ответ a . В общем случае входные данные являются действительными числами и могут быть как положительными, так и отрицательными, хотя в нашем случае все они являются положительными целыми числами.

1.1. Перцептроны

Начнем, однако, с более простой проблемы. Мы создаем программу, чтобы решить, чем является изображение: нулем или не нулем. Это *задача двоичной классификации* (binary classification problem). Одной из ранних схем машинного

обучения для двоичной классификации является *перцептрон* (perceptron), показанный на рис. 1.3.

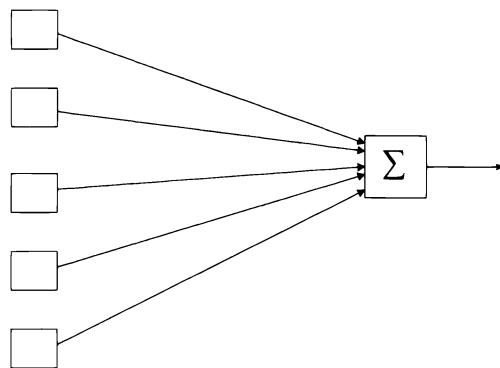


Рис. 1.3. Принципиальная схема перцептрана

Перцептроны были изобретены как простые компьютерные модели нейронов. Один нейрон (рис. 1.4) обычно имеет много входов (дendритов (dendrite)), тело клетки (cell body) и один выход (аксон (axon)). Повторяя это, перцептрон имеет много входов и только один выход. Простой перцептрон для определения, содержит ли наше изображение размером 28×28 нуль, будет иметь 784 входа, по одному на каждый пиксель, и один выход. Для простоты рисунка перцептрон на рис. 1.3 имеет пять входов.

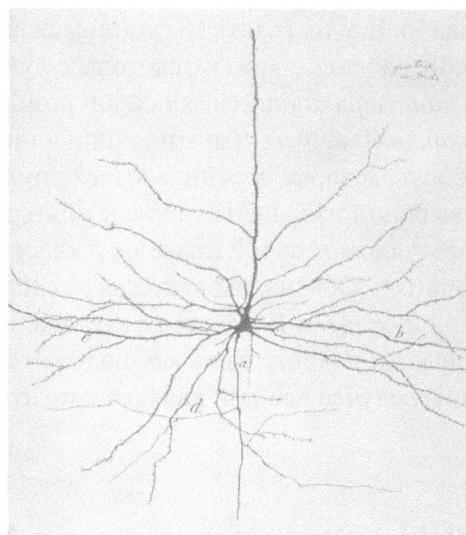


Рис. 1.4. Типичный нейрон

Перцептрон состоит из вектора *весов* (weight) $\mathbf{w} = [w_1 \dots w_m]$, по одному на каждый вход, и специального веса b , называемого *смещением* (bias). Мы называем \mathbf{w} и b *параметрами* (parameter) перцептрана. В более общем смысле мы используем для обозначения параметров Φ , где $\phi_i \in \Phi$ — это i -й параметр. Для перцептрана $\Phi = \{\mathbf{w} \cup b\}$.

При этих параметрах перцептрон вычисляет следующую функцию:

$$f_{\Phi}(\mathbf{x}) = \begin{cases} 1 & \text{если } b + \sum_{i=1}^l x_i w_i > 0 \\ 0 & \text{в противном случае} \end{cases} \quad (1.1)$$

Или, говоря словами, мы умножаем каждый вход перцептрана на вес входного значения и добавляем смещение. Если это значение больше нуля, возвращаем 1, в противном случае — 0. Помните, что перцептраны являются двоичными классификаторами, поэтому 1 указывает, что \mathbf{x} является членом класса, а 0 — не является.

Обычно *скалярное произведение* (dot product) двух векторов длиной l определяют как

$$\mathbf{x} \cdot \mathbf{y} = \sum_{i=1}^l x_i y_i \quad (1.2)$$

поэтому мы можем упростить формулу для вычисления перцептрана следующим образом:

$$f_{\Phi}(\mathbf{x}) = \begin{cases} 1 & \text{если } b + \mathbf{w} \cdot \mathbf{x} > 0 \\ 0 & \text{в противном случае} \end{cases} \quad (1.3)$$

Элементы, которые вычисляют $b + \mathbf{w} \cdot \mathbf{x}$, называют *линейными блоками* (linear unit), и, как на рис. 1.3, мы идентифицируем их с помощью S . Кроме того, обсуждая настройку параметров, полезно преобразовать смещение как другой вес в \mathbf{w} , значение которого всегда равно 1. (Таким образом, нужно поговорить только о настройке \mathbf{w} .)

Мы заботимся о перцептранах, потому что существует удивительно простой и надежный алгоритм — *алгоритм перцептрана* (perceptron algorithm) — для нахождения этих F приведенных *обучающих примеров* (training example). Мы указываем обсуждаемый пример с верхним индексом, поэтому ввод для k -го примера будет $\mathbf{x}^k = [x_1^k \dots x_l^k]$ и его ответ — a^k . Для двоичного классификатора, такого как перцептрон, ответом является 1 или 0, указывающим членство в классе или его отсутствие. При классификации по m классам ответом будет целое число от 0 до $m - 1$.

Иногда машинное обучение полезно характеризовать как задачу *аппроксимации функций* (function approximation). С этой точки зрения единичный блок перцептрана определяет *параметризованный класс* (parameterized class) функций. Изучение весов перцептрана — это выбор члена класса, который наилучшим образом аппроксимирует функцию решения — “истинную” функцию, которая при любом наборе значений пикселей правильно характеризует изображение как, скажем, нуль или не нуль.

Как и во всех исследованиях по машинному обучению, мы предполагаем, что есть как минимум два, а лучше три, набора примеров задач. Первый — это *обучающий набор* (training set), используемый для настройки параметров модели. Второй — это *набор разработки* (development set), используемый для проверки модели, когда мы пытаемся ее улучшить. (Он также называется *выбранным набором* (held-out set) или *проверочным набором* (validation set).) Третий — это *тестовый набор* (test set). Как только модель исправлена и (если нам повезет) дает хорошие результаты, мы оцениваем ее на примерах из тестового набора. Это предотвращает случайную разработку программы, которая работает на наборе разработки, но еще не сталкивалась с незнакомыми задачами. Эти наборы иногда называют *корпусами* (согрода), как “тестовый корпус”. Данные набора Mnist, которые мы используем, доступны в Интернете. Обучающие данные состоят из 60 тысяч изображений и их правильных меток, а набор для разработки и тестирования содержит 10 тысяч изображений и меток.

Отличное свойство алгоритма перцептрана — если существует набор значений параметров, который позволяет перцептрану правильно классифицировать весь обучающий набор, алгоритм гарантированно его найдет. К сожалению, для большинства реальных примеров такого набора нет. С другой стороны, даже тогда перцептраны зачастую работают удивительно хорошо в том смысле, что существуют настройки параметров, которые правильно маркируют очень высокий процент примеров.

Алгоритм работает, перебирая обучающий набор несколько раз и настраивая параметры так, чтобы увеличить количество правильно идентифицированных примеров. Получив обучающий набор без каких-либо параметров, которые необходимо изменить, мы знаем, что у нас есть правильный набор, и можно остановиться. Однако, если такого набора нет, они продолжат изменяться вечно. Чтобы предотвратить это, мы прекращаем обучение после N итераций, где N — это системный параметр, установленный программистом. Обычно N увеличивается с общим количеством изучаемых параметров. Отныне мы будем осторожны, чтобы различать системные параметры Φ и другие числа, связанные с нашей программой, которые мы могли бы иначе назвать “параметрами”, но не являющиеся частью Φ , такие как N — количество переборов обучающего набора. Мы называем последние *гиперпараметрами* (hyperparameter).

Псевдокод этого алгоритма представлен на рис. 1.5. Обратите внимание на стандартное использование Δx для обозначения “изменения в значении x ”.

1. Установить b и все w в 0.
2. Для N итераций или до тех пор, пока веса не перестанут изменяться
 - (а) для каждого учебного примера x^k с ответом a^k
 - i. если $a^k - f(x^k) = 0$ продолжить
 - ii. в противном случае для всех весов w_i , $\Delta w_i = (a^k - f(x^k))x_i$

Рис. 1.5. Алгоритм перцептрона

Критическими строками здесь являются 2(a)i и 2(a)ii. Здесь a^k равно 1 или 0, указывая, является ли изображение членом класса ($a^k = 1$). Таким образом, первая из двух строк говорит, что, если вывод перцептрана имеет правильную метку, ничего делать не нужно. Вторая указывает, как изменить вес w_i так, чтобы, если бы мы немедленно попытались повторить этот пример, перцептран либо понял бы его правильно, либо по крайней мере понял его менее неправильно, а именно — добавил $(a_k - f(x^k))x_i^k$ к каждому параметру w_i .

Наилучший способ убедиться, что строка 2(a)ii делает то, что нужно, — рассмотреть возможные варианты. Предположим, что учебный пример x^k является членом класса. Это означает, что его метка $a^k = 1$. Поскольку мы ошиблись, $f(x^k)$ (вывод перцептрана на примере k -го обучения) должен был быть 0, поэтому $(a^k - f(x^k)) = 1$ для всех i $\Delta w_i = x_i$. Поскольку все значения пикселей ≥ 0 , алгоритм увеличивает вес и в следующий раз $f(x^k)$ возвращает большее значение. Это и есть “менее неправильно”. (Мы оставляем читателю самостоятельно доказать, что формула делает то, что мы хотим в противном случае, когда пример не в классе, а перцептран указывает, что это так.)

Что касается смещения b , мы рассматриваем его как вес для мнимого признака x_0 , значение которого всегда равно 1, и приведенное выше обсуждение проходит без изменений.

Давайте рассмотрим небольшой пример. Здесь мы только смотрим (и корректируем) весовые коэффициенты для четырех пикселей [7, 7] (центр верхнего левого угла), [7, 14] (вверху по центру), [14, 7] и [4, 14]. Обычно имеет смысл делить значения пикселей так, чтобы они находились между нулем и единицей. Предположим, что наше изображение равно нулю, поэтому ($a = 1$), а значения пикселей для этих четырех местоположений равны 8, 9, 6 и 0 соответственно. Поскольку изначально все параметры равны нулю, когда мы оцениваем $f(x)$ на первом изображении $w \cdot x + b = 0$, поэтому $f(x) = 0$, и изображение классифицируется неправильно, а $a(1) - f(x_1) = 1$. При этом вес $w_{7,7}$ становится равным $(0 + 0,8 * 1) = 0,8$. Таким же образом следующие два w_j становятся равными 0,9 и 0,6. Вес центрального пикселя остается нулевым (поскольку значение

изображения здесь равно нулю). Смещение становится равным 1,0. В частности, обратите внимание, что если мы подадим это же изображение в перцептрон второй раз, с новыми весами оно будет классифицировано правильно.

Предположим, что следующее изображение — не нуль, а единица, и два центральных пикселя имеют значение один, а остальные — “ноль”. Вначале $b + \mathbf{w} \times \mathbf{x} = 1 + 0,8 * 0 + 0,9 * 1 + 0,6 * 0 + 0 * 1 = 1,9$, поэтому $f(x) > 0$ и перцептрон не правильно классифицирует пример как нуль. Таким образом, $f(x) - l_x = 0 - 1 = -1$, и мы корректируем каждый вес в соответствии со строкой 2 (a)ii. $w_{0,0}$ и $w_{14,7}$ неизменны, потому что значения пикселей равны нулю, тогда как $w_{7,14}$ теперь становится $0,9 - 0,9 * 1 = 0$ (предыдущее значение минус вес, умноженный на текущее значение пикселя). Новые значения b и $w_{14,14}$ мы оставляем читателю.

Мы перебираем учебные данные несколько раз. Каждый проход по данным называется эпохой (epoch). Обратите также внимание на то, что если учебные данные представлены программе в другом порядке, веса, которые мы изучаем, различны. Хорошей практикой является рандомизация порядка, в котором учебные данные представлены в каждой эпохе. Мы вернемся к этому вопросу в разделе 1.6. Тем не менее студентам, которые изучают этот материал впервые, мы даем здесь некоторую свободу и упускаем эту подробность.

Мы можем расширить перцептроны для решения многоклассовых задач (multiclass decision problem), создав не один перцептрон, а по одному для каждого класса, который мы хотим распознать. Для нашей первоначальной задачи с десятью цифрами у нас их было бы десять, по одному на каждую цифру, а возвращался бы в результате класс с самым высоким значением перцептрана. Графически это показано на рис. 1.6, где представлены три перцептрана для идентификации изображения в одном из трех классов объектов.

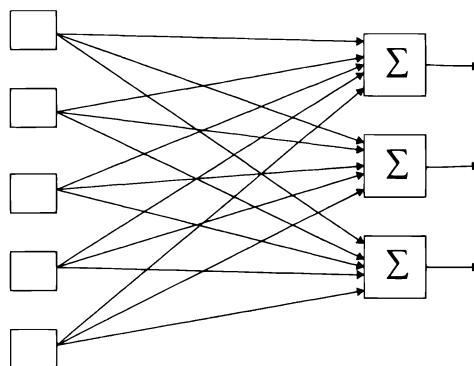


Рис. 1.6. Несколько перцептранов для идентификации нескольких классов

Хотя рис. 1.6 выглядит весьма запутанно, на самом деле он просто демонстрирует три отдельных перцептрана, которые получают одни и те же входные данные. За исключением того факта, что возвращаемый многоклассовым перцептраном ответ является номером линейного блока, который возвращает наибольшее значение, все перцептраны обучаются независимо от других, используя точно такой же алгоритм, который был показан ранее. Итак, учитывая изображение и метку, мы выполняем этап алгоритма перцептрана (a) десять раз для десяти перцептранов. Если меткой является, скажем, пять, но перцептран с наивысшим значением равен шести, то перцептраны от нуля до четырех не меняют своих параметров (поскольку они правильно указали, что это не нуль, не единица и т.д.). То же самое верно для значений от шести до девяти. С другой стороны, пятый и шестой перцептраны изменяют свои параметры, поскольку сообщают о неверных решениях.

1.2. Функция кросс-энтропийных потерь для нейронных сетей

Когда-то обсуждение *нейронных сетей* (далее мы будем сокращенно обозначать их как NN (neural net)) сопровождалось диаграммами, подобными представленной на рис. 1.6, с акцентом на отдельных вычислительных элементах (линейных блоках). Сегодня мы ожидаем, что количество таких элементов будет большим, поэтому говорим о вычислениях с точки зрения *слоев* (layer) — групп хранилищ или вычислительных блоков, которые можно считать работающими параллельно, а затем передающими значения на другой слой. Рис. 1.7 — это версия рис. 1.6, которая подчеркивает данное представление. Здесь показан входной слой, передающий на вычислительный слой.

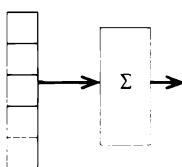


Рис. 1.7. Слои нейронной сети

Слово “слой” подразумевает, что может быть много слоев, каждый из которых добавляется к следующему. Это так и есть, и количество таких слоев является “глубиной” в “глубоком обучении”.

Однако множество слоев плохо работает с перцептранами, поэтому нам нужен еще один метод обучения изменению веса. В этом разделе рассмотрим, как это сделать с помощью простейшей конфигурации сети — *нейронной сети с прямой связью* (feed-forward neural network) — и относительно простой

техники обучения, градиентного спуска (gradient descent). (Некоторые исследователи, напротив, называют NN с прямой связью при обучении с градиентным спуском *многослойным перцептроном* (multilevel perceptron).)

Однако прежде, чем говорить о градиентном спуске, нужно обсудить функции потерь (loss function). Функция потерь оценивает, насколько “плох” для нас результат. Нашей целью в параметрах модели обучения является минимизация потерь. Функция потерь для перцепtronов имеет значение “нуль”, если мы получили правильный учебный пример, и “единица”, если он был неверным. Это *двоичная функция потерь* (zero-one loss). Преимущество двоичной функции потерь в том, что она довольно очевидна, настолько очевидна, что мы никогда не удосуживались оправдать ее использование. Но у нее есть и недостатки. В частности, она плохо работает с обучением градиентным спуском, где основная идея заключается в изменении параметра в соответствии с правилом

$$\Delta\phi_i = -\mathcal{L} \frac{\partial L}{\partial \phi_i} \quad (1.4)$$

Здесь \mathcal{L} — это *скорость обучения* (learning rate), действительное число, измеряющее, насколько мы меняем параметр в данный момент времени. Важной частью является частная производная потери L по параметру, который мы настраиваем. Или, другими словами, если мы можем выяснить, как на потери влияет рассматриваемый параметр, мы должны изменить его так, чтобы уменьшить потери (отсюда и знак “минус”, предваряющий \mathcal{L}). В нашем перцептроне, а в более общем случае — в NN, результат определяется Φ , параметрами модели, поэтому в таких моделях потери являются функцией $L(\Phi)$.

Для упрощения визуализации предположим, что у нашего перцептрана есть только два параметра. Тогда мы можем подумать об евклидовой плоскости с двумя осями, ϕ_1 и ϕ_2 , и для каждой точки на этой плоскости значение функции потерь будет нависать над (или находиться под) точкой. Скажем, наши текущие значения параметров — 1,0 и 2,2 соответственно. Посмотрите на плоскость в положении (1, 2,2) и обратите внимание на то, как ведет себя L в этой точке. На рис. 1.8 демонстрируется срез вдоль плоскости $\phi_2 = 2,2$ и показано, что мнимые потери ведут себя как функция от ϕ_1 . Посмотрим на потери, когда $\phi_1 = 1$. Мы видим, что касательная имеет наклон примерно $-1/4$. Если скорость обучения — $\mathcal{L} = 0,5$, то уравнение 1.4 указывает добавить $(-0,5) * (-1/4) = 0,125$, т.е. сдвинуться на 0,125 единиц вправо, что действительно уменьшает потери.

Чтобы уравнение 1.4 работало, потери должны быть дифференцируемой функцией параметров, а двоичная функция потерь таковой не является. Чтобы увидеть это, представьте график количества ошибок, которые мы совершаем, как функцию от некоторого параметра, ϕ . Скажем, мы только что оценили наш перцептран на примере и ошиблись. Хорошо, если, скажем, мы продолжаем

увеличивать ϕ (или, возможно, уменьшать его) и делаем это достаточно, $f(x)$ в конечном итоге меняет свое значение, и мы получаем правильный пример. Поэтому, на графике мы видим *пошаговую функцию* (step function). Но такие функции не дифференцируемы.

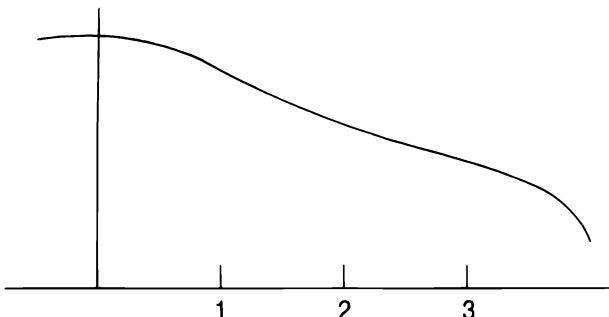


Рис. 1.8. Потери как функция от ϕ_1

Однако существуют и другие функции потерь. Наиболее популярной, наиболее близкой к “стандартной”, функцией потерь является функция *кросс-энтропийных потерь* (cross-entropy loss). В данном разделе мы объясним, что это такое и как наша сеть собирается ее вычислять. В следующем разделе она используется для обучения параметров.

В данный момент сеть на рис. 1.6 выводит вектор значений, по одному на каждый линейный блок, и мы выбираем класс с наибольшим выходным значением. Теперь изменим сеть так, чтобы выходные числа (оценка) были распределением вероятностей по классам, в нашем случае вероятности того, что правильна случайная величина класса — $C = c$ для $c \in [0, 1, 2, \dots, 9]$. *Распределение вероятностей* (probability distribution) — это набор неотрицательных чисел, сумма которых равна единице. В настоящее время сеть выводит числа, но они обычно бывают как положительными, так и отрицательными. К счастью, есть удобная функция для преобразования наборов чисел в распределения вероятностей, *softmax*:

$$\sigma(\mathbf{x})_j = \frac{e^{x_j}}{\sum_i e^{x_i}} \quad (1.5)$$

Softmax гарантированно возвращает распределение вероятности, поскольку даже если x отрицательное, e^x положительное, а значения суммируются в единицу, потому что знаменатель суммирует все возможные значения чисителя. Например, $\sigma([-1, 0, 1]) \approx [0,09; 0,244; 0,665]$. Особый случай, на который мы ссылаемся в нашем дальнейшем обсуждении, — это когда все выходы NN в softmax равны нулю. $e^0 = 1$, поэтому при наличии 10 вариантов все они получают вероятность 1/10, что естественно обобщается до 1/ n , если есть n вариантов.

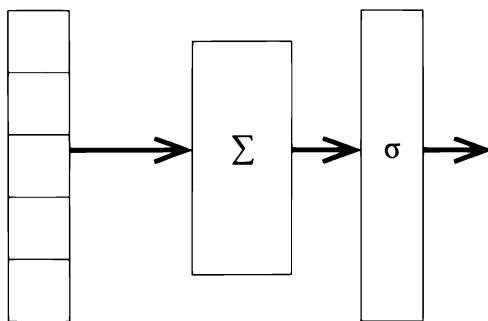


Рис. 1.9. Простая сеть со слоем softmax

На рис. 1.9 показана сеть с дополнительным слоем softmax. Как и ранее, входящие слева числа являются значениями пикселей изображения; однако теперь числа, выходящие справа, являются вероятностями класса. Также полезно иметь имена для чисел, выходящих из линейных блоков и входящих в функцию softmax. Обычно они называются *логитами* (logit) — термин для ненормализованных чисел, которые мы собираемся превратить в вероятности, используя функцию softmax. Для обозначения вектора логитов мы используем \mathbf{l} (по одному для каждого класса). Итак, имеем

$$p(l_i) = \frac{e^{l_i}}{\sum_j e^{l_j}} \quad (1.6)$$

$$\propto e^{l_i} \quad (1.7)$$

Здесь вторая строка отражает тот факт, что, поскольку знаменатель функции softmax является нормализующей константой, для гарантии того, что числа суммируются в единицу, вероятности пропорциональны числителю softmax.

Теперь мы можем определить функцию кросс-энтропийных потерь X :

$$X(\Phi, x) = -\ln p_\Phi(a_x) \quad (1.8)$$

Функция кросс-энтропийных потерь для примера x — это отрицательная логарифмическая вероятность, присвоенная метке x . Или, другими словами, мы вычисляем вероятности всех альтернатив, используя функцию softmax, а затем выбираем одну для правильного ответа. Потеря — это отрицательная логарифмическая вероятность данного числа.

Давайте посмотрим, почему это имеет смысл. В первую очередь, это следует в правильном направлении. Если X — это функция *потерь*, она должна увеличиваться по мере ухудшения модели. Что ж, модель, которая улучшается, должна присваивать более высокую вероятность правильного ответа. Поэтому мы ставим впереди знак “минус”, чтобы число становилось меньше по мере

увеличения вероятности. Кроме того, логарифм числа увеличивается/уменьшается по мере увеличения/уменьшения числа. Таким образом, $X(\Phi, x)$ действительно больше для плохих параметров, чем для хороших.

Но зачем применять логарифм? Мы привыкли думать о логарифмах как о сокращении расстояний между числами. Разница между $\log(10\,000)$ и $\log(1\,000)$ равна 1. Можно было бы подумать, что это будет плохим свойством для функции потерь: из-за этого плохие ситуации выглядят менее плохими. Но эта характеристика логарифмов вводит в заблуждение. Это верно, поскольку, когда x становится больше, $\ln x$ не увеличивается в той же степени. Однако рассмотрим график $-\ln(x)$ на рис. 1.10. По мере приближения x к нулю изменения логарифма намного больше, чем изменения x . А поскольку мы имеем дело с вероятностями, это тот участок, который нас интересует.

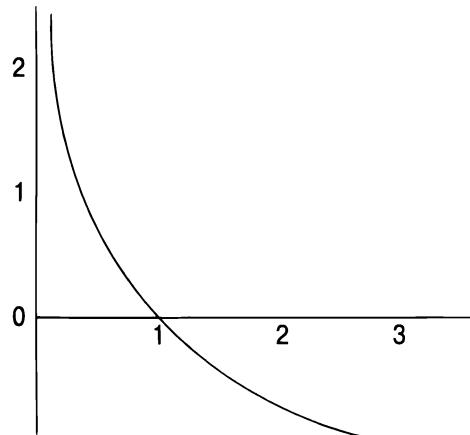


Рис. 1.10. График $-\ln x$

Что касается того, почему эта функция называется функцией *кросс-энтропийных потерь* (cross-entropy loss), в теории информации, когда распределение вероятностей предназначено для аппроксимации некоторого истинного распределения, *кросс-энтропия* (cross entropy) двух распределений является мерой того, насколько они различны. Кросс-энтропийные потери являются приближением отрицательного значения кросс-энтропии. Поскольку в этой книге нам не нужно углубляться в теорию информации, оставим это объяснение поверхностным.

1.3. Производные и стохастический градиентный спуск

Теперь у нас есть функция потерь, и ее можно вычислить, используя следующие уравнения:

$$X(\Phi, x) = -\ln p(a) \quad (1.9)$$

$$p(a) = \sigma_a(\mathbf{l}) = \frac{e^{l_a}}{\sum_i e^{l_i}} \quad (1.10)$$

$$l_j = b_j + \mathbf{x} \cdot \mathbf{w}_j \quad (1.11)$$

Сначала вычисляем логиты \mathbf{l} из уравнения 1.11. Затем они используются слоем softmax для вычисления вероятностей (уравнение 1.10) и мы вычисляем потери, отрицательный натуральный логарифм вероятности правильного ответа (уравнение 1.9). Обратите внимание, что ранее веса для линейного блока обозначались как \mathbf{w} . Теперь у нас есть много таких блоков, поэтому \mathbf{w}_j — это вес для j -го блока, а b_j — его смещение.

Этот процесс перехода от ввода к потере называется *прямым проходом* (forward pass) алгоритма обучения и вычисляет значения, которые будут использоваться в *обратном проходе* (backward pass) — проходе корректировки веса. Для этого имеется несколько методов. Здесь мы используем *стохастический градиентный спуск* (stochastic gradient descent). Термин “градиентный спуск” получил свое название от наклона функции потерь (ее *градиента* (gradient)), с последующим снижением потерь (спуск), следуя градиенту. Этот метод обучения в целом широко известен как *обратное распространение* (back propagation).

Мы начнем с рассмотрения простейшего случая оценки градиента — для одного из смещений, b_j . Из уравнений 1.9–1.11 видно, что b_j изменяет потери, сначала изменяя значение логита l_j , которое затем изменяет вероятность, а следовательно, и потери. Давайте рассмотрим все по этапам. (Здесь мы рассматриваем только ошибку, вызванную одним учебным примером, поэтому пишем $X(\Phi, x)$ как $X(\Phi)$.) Вначале

$$\frac{\partial X(\Phi)}{\partial b_j} = \frac{\partial l_j}{\partial b_j} \frac{\partial X(\Phi)}{\partial l_j} \quad (1.12)$$

Здесь используется цепное правило, чтобы сказать то, что мы сказали ранее словами: изменения в b_j вызывают изменения в X в силу изменений, которые они вызывают в логит l_j .

Теперь посмотрим на первую частную производную справа в уравнении 1.12. Ее значение, по сути, — всего 1:

$$\frac{\partial l_j}{\partial b_j} = \frac{\partial}{\partial b_j} (b_j + \sum_i x_i w_{i,j}) = 1 \quad (1.13)$$

Здесь $w_{i,j}$ — это i -й вес j -го линейного блока. Поскольку в $b_j + \sum_i x_i w_{i,j}$ изменяется только функция b_j , производная равна 1.

Далее рассмотрим, как X изменяется в зависимости от l_j :

$$\frac{\partial X(\Phi)}{\partial l_j} = \frac{\partial p_a}{\partial l_j} \frac{\partial X(\phi)}{\partial p_a} \quad (1.14)$$

Здесь p_i — это вероятность назначения сетью класса i . Таким образом, это указывает, что, поскольку X зависит только от вероятности правильного ответа, l_j влияет на X только за счет изменения этой вероятности. Далее,

$$\frac{\partial X(\phi)}{\partial p_a} = \frac{\partial}{\partial p_a} (-\ln p_a) = \frac{1}{p_a} \quad (1.15)$$

(из основ математики).

Это оставляет для оценки еще один член:

$$\frac{\partial p_a}{\partial l_j} = \frac{\partial \sigma_a(l)}{\partial l_j} = \begin{cases} (1 - p_j)p_a & a = j \\ -p_j p_a & a \neq j \end{cases} \quad (1.16)$$

Первое равенство уравнения 1.16 вытекает из того факта, что мы получаем вероятности, вычисляя функцию softmax для логитов. Второе равенство следует из Википедии. Вывод требует тщательного манипулирования членами, и мы его не выполняем. Тем не менее можно сделать это резонным. Мы спрашиваем, как изменения в логите l_j повлияют на вероятность, исходящую из softmax. Напомним себе, что

$$\sigma_a(l) = \frac{e^{l_a}}{\sum_i e^{l_i}}$$

имеет смысл, что есть два случая. Предположим, логит, который мы меняем (j), не равен a . То есть предположим, что это изображение числа 6, но мы интересуемся смещением, которое определяет логит 8. В этом случае l_j появляется только в знаменателе, а производная должна быть отрицательной (или нулевой), поскольку чем больше l_j , тем меньше p_a . Это второй случай в уравнении 1.16, и, конечно же, он дает число, меньшее или равное нулю, поскольку две вероятности, которые мы умножаем, не могут быть отрицательными.

С другой стороны, если $j = a$, то l_j появляется как в числителе, так и в знаменателе. Его появление в знаменателе стремится уменьшить выход, но в этом случае оно более чем компенсируется увеличением числителя. Таким образом, для этого случая мы ожидаем положительную (или нулевую) производную, и это то, что дает первый случай уравнения 1.16.

Получив этот результат, мы можем теперь вывести уравнение для изменения параметров смещения b_j . Подстановка уравнений 1.15 и 1.16 в уравнение 1.14 дает

$$\frac{\partial X(\Phi)}{\partial l_j} = -\frac{1}{p_a} \begin{cases} (1-p_j)p_a & a=j \\ -p_j p_a & a \neq j \end{cases} \quad (1.17)$$

$$= \begin{cases} -(1-p_j) & a=j \\ p_j & a \neq j \end{cases} \quad (1.18)$$

Остальное довольно просто. Мы заметили в уравнении 1.12, что

$$\frac{\partial X(\Phi)}{\partial b_j} = \frac{\partial l_j}{\partial b_j} \frac{\partial X(\Phi)}{\partial l_j}$$

и далее, что первая из производных справа имеет значение 1. Таким образом, производная потеря по b_j определяется уравнением 1.14. Наконец, используя правило для изменения весов (уравнение 1.12), получаем правило для обновления параметров смещения NN:

$$\Delta b_j = \mathcal{L} \begin{cases} (1-p_j) & a=j \\ -p_j & a \neq j \end{cases} \quad (1.19)$$

Уравнение для изменения весовых параметров (в отличие от смещения) является незначительным изменением уравнения 1.19. Вот уравнение, соответствующее уравнению 1.12 для весов:

$$\frac{\partial X(\Phi)}{\partial w_{i,j}} = \frac{\partial l_j}{\partial w_{i,j}} \frac{\partial X(\Phi)}{\partial l_j} \quad (1.20)$$

В первую очередь, обратите внимание, что крайняя справа производная та-
кая же, как в уравнении 1.12. Это означает, что на этапе корректировки веса мы
должны сохранить этот результат, когда смещаем изменения, чтобы использо-
вать его здесь. Первая из двух производных справа преобразуется в

$$\frac{\partial X(\Phi)}{\partial w_{i,j}} = \frac{\partial}{\partial w_{i,j}} (b_j + (w_{1,j} x_1 + \dots + w_{i,j} x_i + \dots)) = x_i \quad (1.21)$$

(Если принять близко к сердцу идею о том, что смещение — это просто вес, соответствующее значение признака которого всегда равно 1, можно было бы просто вывести это уравнение, и тогда уравнение 1.13 следовало бы непосредственно из 1.21 применительно к этому новому псевдовесу.)

Используя этот результат, получаем уравнение для обновлений веса:

$$\Delta w_{i,j} = -\mathcal{L} x_i \frac{\partial X(\Phi)}{\partial l_j} \quad (1.22)$$

Теперь мы выяснили, как параметры нашей модели должны быть скорректированы в свете одного примера обучения. Затем алгоритм *градиентного спуска* (gradient descent) заставил бы нас пройти все обучающие примеры, в которых записывалось, как каждый из них рекомендовал бы перемещать значения параметров, но не менять их на самом деле, пока мы не сделаем полный проход через них всех. На этом этапе мы модифицируем каждый параметр суммой изменений от отдельных примеров.

Проблема здесь в том, что этот алгоритм может быть очень медленным, особенно если обучающий набор велик. Как правило, нам необходимо часто настраивать параметры, поскольку они будут взаимодействовать по-разному, ведь каждый из них увеличивается или уменьшается в результате конкретных тестовых примеров. Таким образом, на практике мы почти никогда не используем градиентный спуск, а используем только *стохастический градиентный спуск* (stochastic gradient descent), который обновляет параметры каждые m примеров для m , намного меньшего, чем размер всего обучающего набора. Обычно для m выбирают значение 20. Это называется *размером партии* (batch size).

В общем, чем меньше размер партии, тем меньшей должна быть установлена скорость обучения \mathcal{L} . Идея состоит в том, что любой пример подтолкнет весы к правильной классификации этого примера за счет других. Если скорость обучения низкая, это не имеет большого значения, так как изменения, внесенные в параметры, соответственно малы. И наоборот, при большем размере партии мы неявно усредняем по m различным примерам, поэтому опасность перекоса параметров для идиосинкразий одного примера уменьшается, а изменения, вносимые в параметры, могут быть больше.

1. Для j от 0 до 9 установить b_j случайным (но близким к нулю).
2. Для j от 0 до 9 и для i от 0 до 783 установить $w_{i,j}$ аналогично.
3. Пока точность разработки не прекратит расти
 - (a) для каждого учебного примера k в партиях из m примеров
 - i. выполнить прямой проход, используя уравнения 1.9–1.11
 - ii. выполнить обратный проход, используя уравнения 1.22, 1.19 и 1.14
 - iii. каждые m примеров модифицировать все Φ суммарными дополнениями
 - (b) вычислить точность модели, выполнив прямой проход для всех примеров в корпусе разработки
4. Вывести Φ из итерации до снижения точности разработки.

Рис. 1.11. Псевдокод для простого распознавания цифр с прямой связью

1.4. Написание программы

Теперь у нас есть широкий охват нашей первой программы NN. Псевдокод представлен на рис. 1.11. Начиная сверху, первое, что мы делаем, — инициализируем параметры модели. Иногда хорошо инициализировать все нулями, как мы это делали в алгоритме перцептрона. Хотя это относится и к нашей текущей задаче, это не всегда так. Таким образом, хорошей общепринятой практикой является установка случайных значений весов, близких к нулю. Возможно, вы также захотите задать для генератора случайных чисел Python начальное значение, поэтому при отладке вы всегда устанавливаете параметры равными одним и тем же начальным значениям, а следовательно, должны получать точно такие же выходные данные. (Если вы этого не сделаете, Python использует в качестве начального числа значения из окружения, такие как последние несколько цифр из системных часов.)

Обратите внимание, что на каждой итерации обучения мы сначала изменяем параметры, а затем используем модель на наборе разработки, чтобы увидеть, насколько хорошо он работает с текущим набором параметров. Запуская примеры разработки, мы *не* запускаем обратный обучающий проход. Если мы на самом деле собирались использовать нашу программу для каких-то реальных целей (например, для чтения почтовых индексов на почте), то примеры, которые мы видим, являются не теми, на которых мы обучались, а потому мы хотим знать, насколько хорошо наша программа работает “в дикой природе”. Наши данные для разработки являются приближением к этой ситуации.

Здесь пригодятся некоторые эмпирические знания. В первую очередь, это общепринятая практика, когда значения пикселей не отклоняются слишком далеко от -1 до 1 . В нашем случае, поскольку исходные значения пикселей — от 0 до 255 , мы просто делим их на 255 , прежде чем использовать в нашей сети. Это пример процесса *нормализации данных* (data normalization). Жестких и быстрых правил не существует, но зачастую имеет смысл хранить значения в диапазонах от -1 до 1 или от 0 до 1 . Можно понять, почему это действительно так, в уравнении 1.22, в котором мы видели, что различие в уравнением для корректировки члена смещения и веса, поступающего от одного из входов NN, заключалось в том, что последний имел мультипликативный член x_i , значение входного члена. Хотя мы сказали, что если бы мы приняли к сведению наш комментарий о том, что член смещения был просто весовым членом, входное значение которого всегда было равно 1 , уравнение для обновления параметров смещения следовало бы из уравнения 1.22. Таким образом, если мы оставляем входные значения неизмененными, а один из пикселей имеет значение 255 , мы изменили его весовое значение в 255 раз больше, чем изменили смещение. Учитывая, что у нас нет априорной причины полагать, что одному нужно большее исправление, чем другому, это покажется странным.

Далее стоит вопрос об установке \mathcal{L} , скорости обучения. Это может быть сложно. В нашей реализации мы использовали значение 0,0001. Первое, что нужно отметить, — установка слишком большого значения намного хуже, чем слишком малого. Сделав это, вы получите от softmax математическую ошибку переполнения. Снова обратимся к уравнению 1.5. Одна из первых вещей, которые должны вас поразить, — это экспоненты как в числителе, так и в знаменателе. Повышение до значения, большого e ($\approx 2,7$), — это надежный способ получить переполнение, что мы и будем делать, если какой-либо из логитов станет большим, что, в свою очередь, может произойти, если у нас слишком высокая скорость обучения. Даже если сообщение об ошибке не предупреждает вас о том, что что-то не так, слишком высокая скорость обучения может привести к тому, что ваша программа окажется в невыгодной области кривой обучения.

По этой причине стандартной практикой является наблюдение за тем, что происходит с потерями на отдельных примерах в процессе вычислений. Давайте начнем с того, чего ожидать от первого обучающего изображения. Числа проходят через NN и попадают на слой логитов. Все веса и смещения будут нулевыми плюс или минус совсем немного (скажем 0,1). Это означает, что все значения логита очень близки к нулю, поэтому все вероятности очень близки к 1/10. (См. обсуждение на с. 13.) Потеря составляет минус натуральный логарифм вероятности, назначенной для правильного ответа, $-\ln(1/10) \approx 2,3$. Согласно общей тенденции мы ожидаем, что индивидуальные потери сократятся, поскольку мы обучаемся на новых примерах. Но, естественно, одни изображения находятся дальше от нормы, чем другие, и, таким образом, классифицируются NN с меньшей уверенностью. Таким образом, мы видим индивидуальные потери, которые расположены выше или ниже, и эту тенденцию может быть трудно распознать. Таким образом, вместо того чтобы выводить по одной потере за раз, мы суммируем их все по мере продвижения и выводим среднее значение, скажем, каждые 100 партий. Это среднее значение должно уменьшаться вполне заметным образом, хотя даже здесь вы можете увидеть дребезг (девиации).

Возвращаясь к нашему обсуждению скорости обучения и опасностей, связанных с ее слишком высоким значением, слишком низкая скорость обучения может реально замедлить скорость, с которой ваша программа сходится к хорошему набору параметров. Поэтому смотреть на маленькие значения и экспериментировать с большими — это, как правило, наилучший образ действий.

Поскольку одновременно изменяется так много параметров, алгоритмы NN могут быть сложными для отладки. Как и при любой отладке, хитрость заключается в том, чтобы сделать как можно меньше изменений, прежде чем ошибка проявится. Во-первых, помните о том, что, когда мы изменяем веса, если вы сразу же запускаете тот же обучающий пример, потери будут меньше. Если это не так, то либо есть ошибка, либо вы установили слишком высокую скорость

обучения. Во-вторых, помните, что нет необходимости менять все веса, чтобы увидеть уменьшение потерь. Вы можете изменить только один из них или одну группу из них. Например, при первом запуске алгоритма меняются только смещения. (Однако, если вы подумаете об этом, смещение в однослойной сети в основном будет отражать тот факт, что разные классы встречаются с разными частотами. В данных Mnist этого не так много, поэтому мы не добьемся значительного улучшения, просто участь на смещениях в данном случае.)

Если ваша программа работает правильно, вы должны получить точность данных для разработки порядка 91% или 92%. Это не очень хорошо для данной задачи. В последующих главах мы увидим, как достичь примерно 99%. Но это только начало.

В действительно простых NN есть одна приятная вещь: иногда мы можем напрямую интерпретировать значения отдельных параметров и решить, являются ли они резонными. Возможно, вы помните, что в нашем обсуждении рис. 1.1 мы отметили, что пиксель (8,8) был темным — у него было значение пикселя 254. Мы отметили, что это было диагностическим признаком для изображений цифры 7 в отличие от, например, цифры 1, которая обычно не имеет маркировки в верхнем левом углу. Мы можем превратить это наблюдение в прогноз значений в нашей весовой матрице $w_{i,j}$, где i — это номер пикселя, а j — значение ответа. Если значения пикселей изменяются от 0 до 784, то позиция (8,8) будет равна пикселию $8 \times 28 + 8 = 232$, а вес, соединяющий его с ответом 7 (правильный ответ), будет $w_{232,7}$, тогда как его применение к 1 дало бы $w_{232,1}$. Вы должны убедиться, что видите, что теперь предполагается, что $w_{232,7}$ должно быть больше, чем $w_{232,1}$. Мы запускали нашу программу несколько раз с низкой дисперсией случайной инициализации наших весов. В каждом случае первое число было положительным (например, 0,25), а второе — отрицательным (например, -17).

1.5. Матричное представление нейронных сетей

Линейная алгебра дает другой способ представления того, что происходит в NN: использование матриц. *Матрица* (matrix) — это двумерный массив элементов. В нашем случае эти элементы являются действительными числами. Размеры матрицы — это количество строк и столбцов соответственно. Итак, матрица l на m выглядит так:

$$\mathbf{X} = \begin{pmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,m} \\ x_{2,1} & x_{2,2} & \dots & x_{2,m} \\ \dots \\ x_{l,1} & x_{l,2} & \dots & x_{l,m} \end{pmatrix} \quad (1.23)$$

Основными операциями над матрицами являются сложение и умножение. Сложение двух матриц (которые должны иметь одинаковые размеры) происходит поэлементно, т.е., если мы сложим две матрицы $\mathbf{X} = \mathbf{Y} + \mathbf{Z}$, то $x_{i,j} = y_{i,j} + z_{i,j}$.

Умножение двух матриц $\mathbf{X} = \mathbf{Y}\mathbf{Z}$ определяется, когда \mathbf{Y} имеет размеры l и m , а таковые у \mathbf{Z} равны m и n . В результате получается матрица размером l на n , где

$$x_{i,j} = \sum_{k=1}^{k=m} y_{i,k} z_{k,j} \quad (1.24)$$

Вот простой пример:

$$(1 \ 2) \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} + (7 \ 8 \ 9) = (9 \ 12 \ 15) + (7 \ 8 \ 9) = \\ = (16 \ 20 \ 24)$$

Мы можем использовать эту комбинацию умножения и сложения матриц, чтобы определить работу линейных блоков. В частности, входными признаками являются матрица \mathbf{X} размером $1 * l$. В задаче о цифрах $l = 784$. Весами в блоках являются \mathbf{W} , где $w_{i,j}$ — это i -й вес для блока j . Таким образом, размеры \mathbf{W} — это количество пикселей, умноженное на количество цифр, $784 * 10$. \mathbf{B} — это вектор смещений длиной 10, и

$$\mathbf{L} = \mathbf{X}\mathbf{W} + \mathbf{B} \quad (1.25)$$

Здесь \mathbf{L} — это вектор логитов длиной 10. Когда вы впервые видите такое уравнение, имеет смысл убедиться, что размеры работают.

Теперь мы можем выразить потери (L) для нашей модели Mnist с прямой связью следующим образом:

$$\text{Pr}(A(x)) = \sigma(\mathbf{x}\mathbf{W} + \mathbf{b}) \quad (1.26)$$

$$L(x) = -\log(\text{Pr}(A(x) = a)) \quad (1.27)$$

Здесь первое уравнение дает распределение вероятностей по возможным классам ($A(x)$), а второе задает кросс-энтропийные потери.

Мы также можем выразить обратный проход более компактно. Сначала введем *оператор градиента* (gradient operator)

$$\nabla_l X(\Phi) = \left(\frac{\partial X(\Phi)}{\partial l_1} \dots \frac{\partial X(\Phi)}{\partial l_m} \right) \quad (1.28)$$

Перевернутый треугольник, $\nabla_x f(\mathbf{x})$, обозначает вектор, созданный в результате взятия частной производной от f по всем значениям в \mathbf{x} . Ранее мы говорили о частной производной по индивидуальному l_j . Здесь мы определяем

производную по всем \mathbf{l} как вектор отдельных производных. Мы также напоминаем читателю о транспонировании матрицы — превращении строк матрицы в столбцы и наоборот:

$$\begin{pmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,m} \\ x_{2,1} & x_{2,2} & \dots & x_{2,m} \\ \dots & & & \\ x_{l,1} & x_{l,2} & \dots & x_{l,m} \end{pmatrix}^T = \begin{pmatrix} x_{1,1} & x_{2,1} & \dots & x_{l,1} \\ x_{1,2} & x_{2,2} & \dots & x_{l,2} \\ \dots & & & \\ x_{1,m} & x_{2,m} & \dots & x_{l,m} \end{pmatrix} \quad (1.29)$$

С их помощью мы можем переписать уравнение 1.22 как

$$\Delta \mathbf{W} = -\mathcal{L} \mathbf{X}^T \nabla_{\mathbf{l}} X(\Phi) \quad (1.30)$$

Справа мы умножаем $784 * 1$ раз матрицу $1 * 10$, чтобы получить матрицу $784 * 10$ изменений для матрицы $784 * 10$ весов \mathbf{W} .

Это элегантное резюме происходящего, когда входной слой передает результат на слой линейных блоков для создания логитов, за которыми следуют производные потерь, распространяющиеся обратно к изменениям параметров. Но есть и практическая причина для предпочтения этой новой нотации. При работе с большим количеством линейных блоков линейная алгебра вообще и глубокое обучение в частности могут отнять очень много времени. Однако в матричной нотации можно выразить очень много задач, и многие языки программирования имеют специальные пакеты, которые позволяют программировать с использованием конструкций линейной алгебры. Кроме того, эти пакеты оптимизированы, чтобы сделать их более эффективными, чем если бы вы программируали их вручную. В частности, если вы программируете на языке Python, стоит использовать пакет *Numpy* и его матричные операции. Обычно вы получаете ускорение на порядок.

Кроме того, одним конкретным применением линейной алгебры являются компьютерная графика и ее использование в игровых программах. Это привело к созданию специализированного оборудования, *графических процессоров* (Graphics Processing Unit — GPU). Графические процессоры содержат медленные процессоры по сравнению с *центральными процессорами* (CPU), но их много, а также есть программное обеспечение для их эффективного параллельного применения в линейных алгебраических вычислениях. Некоторые специализированные языки для NN (например, *Tensorflow*) имеют встроенное программное обеспечение, которое определяет наличие графических процессоров и использует их без каких-либо изменений в коде. Это обычно дает еще один порядок увеличения скорости.

В этом случае есть еще одна, третья, причина для принятия матричной нотации. Как пакеты специального программного обеспечения (например, *Numpy*),

так и аппаратные средства (графические процессоры) более эффективны, если параллельно обрабатывать несколько обучающих примеров. Кроме того, это согласуется с идеей, что мы хотим обработать некоторое количество m обучающих примеров (размер партии), прежде чем обновить параметры модели. Поэтому обычной практикой является ввод всех m из них в нашу матричную обработку совместно. В уравнении 1.25 мы представили изображение x в виде матрицы размером 1×784 . Это был один учебный пример с 784 пикселями. Теперь мы изменим это так, чтобы матрица имела размеры $m \times 784$. Интересно, что это почти работает без каких-либо изменений в нашей обработке (и необходимые изменения уже встроены, например, в NumPy и Tensorflow). Посмотрим, почему.

Вначале рассмотрим умножение матриц XW , где теперь X имеет m строк, а не 1. Конечно, с одной строкой мы получаем вывод размером 1×784 . С m строками получается $m \times 784$. Кроме того, как вы помните из линейной алгебры, хотя это можно подтвердить, посмотрев определение матричного умножения, выходные строки выглядят так, как будто в каждом случае мы делаем умножение одной строки, а затем складываем их вместе, чтобы получить матрицу $m \times 784$.

Добавление члена смещения в уравнение не очень хорошо работает. Мы сказали, что для сложения матриц необходимо, чтобы обе матрицы имели одинаковые размеры. Это больше не верно для уравнения 1.25, так как XW имеет теперь размер $m \times 10$, тогда как B , член смещения, имеет размер 1×10 . Именно здесь происходят скромные изменения.

NumPy и Tensorflow имеют *широковещание* (broadcasting). Когда для какой-либо арифметической операции требуется, чтобы массивы имели размеры, отличные от имеющихся, размеры массивов иногда можно отрегулировать. В частности, когда один из массивов имеет размерность $1 \times n$, а нам требуется $m \times n$, первый получает $m - 1$ (виртуальных) копий, сделанных из одной строки или столбца так, чтобы она имела правильный размер. Это именно то, что мы хотим сделать здесь. Это фактически придает B размер $m \times 10$. Таким образом, мы добавляем смещение ко всем членам в результате умножения $m \times 10$. Помните, что мы делали, когда это был размер 1×10 ? Каждое из десяти решений было одним из возможных решений о том, каким может быть правильный ответ, и мы добавили смещение к числу для этого решения. Сейчас мы делаем то же самое, но для каждого возможного решения и для всех m примеров мы работаем параллельно.

1.6. Независимость данных

Все теоремы о том, что если выполняются следующие предположения, то наши модели NN фактически сходятся к правильному решению в зависимости от *предположения iid* (*iid assumption*), что наши данные *независимы и*

одинаково распределены (independent and identically distributed — iid). Каноническим примером являются измерения космических лучей, ниспадающие потоки которых случайны и неизменны.

Наши данные редко (почти никогда) выглядят так: представьте, что Национальный институт стандартов предоставляет постоянный поток новых примеров. Для данных первой эпохи предположение iid выглядит довольно хорошо, но как только мы начинаем со второй, наши данные становятся идентичными первой. В некоторых случаях наши данные могут не нарушать iid начиная с обучающего примера 2. Это часто имеет место при *глубоком обучении с подкреплением* (deep reinforcement learning — deep RL) (см. главу 6), и по этой причине сети в данной ветви глубокого RL зачастую страдают от *нестабильности* (instability) — отказа сети, когда она приходит к правильному решению не всегда. Здесь мы рассмотрим сравнительно небольшой пример, когда простой ввод данных в неслучайном порядке может привести к катастрофическим результатам.

Предположим, что для каждого изображения Mnist мы добавили второе, которое идентично первому, но с обращением черных и белых цветов, т.е. если оригинал имеет значение пикселя v , то обращенное изображение имеет значение $-v$. Теперь мы обучим наш перцептрон Mnist на этом новом корпусе, но с использованием другого порядка учебных примеров. (И мы предполагаем, что размер партии равен некоторому четному целому числу.) При первом упорядочении каждого исходного изображения цифре Mnist сразу же следует его обращенная версия. Утверждаю (и мы подтвердили это эмпирически), что наша простая NN Mnist работает не лучше, чем случайность. Мысль о нескольких моментах должна сделать это резонным. Мы видим изображение единицы, и обратный проход изменяет веса. Теперь мы обрабатываем второе, обращенное изображение. Поскольку вход — это минус предыдущий вход, а все остальное одинаково, изменения всех весов в точности отменяют предыдущие, и в конце обучающего набора не будет изменений ни одного из весов. Так что обучения нет, есть лишь случайный выбор, с которого мы начали.

С другой стороны, на самом деле не должно быть ничего слишком сложного в том, чтобы научиться обрабатывать каждый набор данных, обычный и обращенный, отдельно, и для одного набора весов должно быть лишь немного сложнее справиться с обоими. Действительно, простой рандомизации порядка ввода достаточно, чтобы вернуть производительность почти на уровень исходной задачи. Если мы увидим обращенное изображение, скажем, через 10 тысяч выборок, веса изменятся достаточно, поэтому обращенное изображение не полностью компенсирует исходное обучение. Если бы у нас был бесконечный источник изображений и мы подбросили монетку, чтобы принять решение передать NN оригинал или инверсию, то даже эта небольшая отмена исчезла бы.

1.7. Ссылки и что читать дальше

В этом и последующих разделах “Ссылки и что читать дальше” я постараюсь сделать несколько вещей более или менее одновременно: а) указать учащемуся на последующий материал по теме главы, б) определить наиболее важные вклады в эту область и в) привести ссылки, которые я сам использовал для изучения этого материала. Во всех случаях, особенно в случае б, я не претендую на полноту или объективность. Я понял это, когда, готовясь к написанию данного раздела, начал читать историю нейронных сетей. В частности, я прочитал запись в блоге Андрея Куренкова (Andrey Kurenkov) [1], чтобы освежить свои воспоминания (и, возможно, внести дополнения).

Одной из ключевых ранних работ по NN была работа Мак-Каллока и Питтса [2], которые предложили то, что мы называем здесь линейным блоком, в качестве формальной модели нейрона. Это было еще в 1943 году. Но у них не было алгоритма обучения, который мог бы обучить одного или нескольких из них выполнению задачи. Это был большой вклад Розенблатта в его статью о перцептроне 1958 года [3]. Однако, как мы отмечали выше, его алгоритм работал только для однослойного NN.

Следующим большим шагом было изобретение обратного распространения, которое *работает* для многослойных NN. Это была одна из тех ситуаций, когда многие исследователи пришли к идеи независимо друг от друга в течение нескольких лет. (Это, конечно, происходит только тогда, когда первоначальные статьи не привлекают достаточного внимания, чтобы все остальные узнали, что проблема уже решена.) Документ, который завершил этот период, был подготовлен Румельхартом, Хинтоном и Вильямсом и явно отмечает их повторное открытие [4]. Эта статья была одной из многих в группе из Университета Сан-Диего, которая была ответственна за второй расцвет нейронных сетей под рубрикой *параллельной распределенной обработки* (Parallel Distributed Processing — PDP). Двухтомная коллекция этих работ была весьма влиятельной [5].

Что касается того, как я узнал все, что я знаю об NN, я приведу больше подробностей в последующих главах. Для этой главы, я помню, сначала читал блог Стивена Миллера [6], в котором очень медленно и с множеством примеров излагались прямой и обратный проходы обратного распространения. В целом позвольте мне отметить два общих учебника по NN, с которыми я ознакомился. Одним из них является *Deep Learning* Яна Гудфеллоу, Йошуа Бенджио и Аарона Курвиля [7]; второй — это *Hands-On Machine Learning with Scikit-Learn and Tensorflow* Орельена Жерона [8].

1.8. Упражнения

Упражнение 1.1. Рассмотрим нашу программу прямой связи Mnist с размером партии 1. Предположим, мы смотрим на переменные смещения до и после обучения на первом примере. Если они установлены правильно (т.е. если в нашей программе нет ошибок), опишите изменения, которые вы должны увидеть в их значениях.

Упражнение 1.2. Мы упрощаем наши вычисления Mnist, предполагая, что наше “изображение” имеет два пикселя с двоичными значениями, 0 и 1, нет параметров смещения и мы выполняем задачу двоичной классификации. а) Вычислите логиты и вероятности прямого прохода, когда значения пикселей равны [0,1], а веса

$$\begin{matrix} 0,2 & -0,3 \\ -0,1 & 0,4 \end{matrix}$$

Здесь $w[i, j]$ — это вес связи между i -м пиксели и j -м блоком. Например, $w[0, 1]$ здесь равно $-0,3$. б) Предположим, что правильный ответ — 1 (а не 0) и использовалась скорость обучения 0,1. Каковы потери? Вычислите также $\Delta w_{0,0}$ на обратном проходе.

Упражнение 1.3. Те же вопросы, что и в упражнении 1.2, за исключением изображения [0,0].

Упражнение 1.4. Один из студентов спрашивает вас: “В элементарном исчислении мы нашли минимумы функции, дифференцируя ее, устанавливая полученное выражение равным нулю, а затем решая уравнение. Поскольку наша функция потерь является дифференцируемой, почему бы нам не сделать это, а не заниматься градиентным спуском?” Объясните, почему это на самом деле невозможно.

Упражнение 1.5. Вычислите следующее:

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 2 & 3 \end{pmatrix} + (4 \quad 5) \tag{1.31}$$

Вы должны принять широковещание, чтобы вычисления были четко определены.

Упражнение 1.6. В этой главе мы ограничились проблемами классификации, для которых кросс-энтропия, как правило, является функцией потерь по выбору. Есть также проблемы, когда мы хотим, чтобы наша NN прогнозировала конкретные значения. Например, несомненно, многим людям нужна программа, которая, учитывая цену конкретной акции сегодня, а также всевозможные другие факты о мире, выводит цену акции завтра. Если бы мы обучали для этого однослойную NN, мы обычно использовали бы *квадратичную функцию потерь* (squared-error loss):

$$L(\mathbf{X}, \Phi) = (t - l(\mathbf{X}, \Phi))^2 \quad (1.32)$$

Здесь t — это фактическая цена, которая была достигнута в этот день, а $l(\mathbf{X}, \Phi)$ — это выход одного слоя NN при $\Phi = \{\mathbf{b}, \mathbf{W}\}$. (Это также известно как *квадратичная потеря* (quadratic loss).) Выведите уравнение для производной потери по b_i .

Глава 2

Язык Tensorflow

2.1. Вводная информация о Tensorflow

Tensorflow — это язык программирования с открытым исходным кодом, специально разработанный Google для упрощения программирования в области глубокого обучения. Начнем с традиционной первой программы:

```
import tensorflow as tf
x = tf.constant("Hello World")
sess = tf.Session()
print(sess.run(x)) # Выводит "Hello World"
```

Если это похоже на код языка Python, то это потому, что так и есть. Фактически язык Tensorflow (далее TF) — это набор функций, которые можно вызывать из разных языков программирования. Самый полный интерфейс — у языка Python, и поэтому мы используем здесь именно его.

Следующее, что необходимо отметить, — это то, что функции TF не столько выполняют программу, сколько определяют вычисления, которые возможны только при вызове команды `run`, как в последней строке вышеупомянутой программы. Точнее, функция TF `Session` в строке 3 создает сеанс, и с этим сеансом связан граф, определяющий вычисления. Такие команды, как `constant`, добавляют элементы в это вычисление. В данном случае элемент является просто постоянным элементом данных, значением которого является строка Python “Hello World”. Третья строка указывает TF оценить переменную TF, на которую указывает `x` внутри графа, связанного с сеансом `sess`. Как и следовало ожидать, это приводит к выводу слов “Hello World”.

Поучительно сопоставлять это поведение с тем, что происходит при замене последней строки строкой `print(x)`. Она выведет

```
Tensor("Const:0", shape=(), dtype=string)
```

Дело в том, что переменная Python `x` связана не со строкой, а с частью графа вычислений Tensorflow. Только когда мы интерпретируем эту часть графа, выполняя функцию `sess.run(x)`, мы получаем доступ к значению константы TF.

Таким образом для пояснения, возможно очевидного, в приведенном выше коде `x` и `sess` являются переменными Python, а поэтому могут быть присвоены так, как мы хотели. `import` и `print` являются функциями Python и должны быть написаны таким образом, чтобы Python понимал, какую функцию мы хотим выполнить. Наконец, `constant`, `Session` и `run` являются командами TF, и снова написание должно быть точным (включая прописную “S” в `Session`). Кроме того, нам всегда сначала нужна строка `import tensorflow`. Поскольку это упомянуто, мы опускаем ее вперед.

```
x = tf.constant(2.0)
z = tf.placeholder(tf.float32)
sess= tf.Session()
comp=tf.add(x,z)
print(sess.run(comp,feed_dict={z:3.0})) # Выводит 5.0
print(sess.run(comp,feed_dict={z:16.0})) # Выводит 18.0
print(sess.run(x)) # Выводит 2.0
print(sess.run(comp)) # Выводит очень длинное сообщение об ошибке
```

Рис. 2.1. Заполнители в TF

В коде на рис. 2.1 `x` снова является переменной Python, значение которой является константой TF, в данном случае это число с плавающей запятой 2,0. Далее, `z` — это переменная Python, значение которой является *заполнителем* (*placeholder*) TF. Заполнитель в TF подобен формальной переменной в функции языка программирования. Предположим, у нас был следующий код Python:

```
x = 2.0
def sillyAdd(z):
    return z+x
print(sillyAdd(3)) # Выводит 5.0
print(sillyAdd(16)) # Выводит 18.0
```

Здесь `z` — это имя аргумента функции `sillyAdd`, и когда мы вызываем функцию, как в случае `sillyAdd(3)`, оно заменяется значением 3. Версия TF работает аналогично, за исключением того, что способ присвоения заполнителям TF значения несколько иной, как видно из строки 5 на рис. 2.1:

```
print(sess.run(comp,feed_dict={z:3.0}))
```

Здесь `feed_dict` является именованным аргументом функции `run` (поэтому его имя должно быть написано правильно). Он принимает в качестве возможных значений словарь Python. В словаре каждому заполнителю, требуемому для вычисления, должно быть присвоено значение. Таким образом, первый вызов `sess.run` выводит сумму 2,0 и 3,0, а второй — 18,0. Третье замечание: если вычисление не требует значения заполнителя, его не нужно указывать. С другой

стороны, как указано в комментарии к четвертому оператору `print`, если для вычисления требуется значение, а оно не указано, вы получите ошибку.

Tensorflow называется так потому, что его фундаментальными структурами данных являются *тензоры* (*tensor*) — многомерные типизированные массивы. Существует порядка пятнадцати или около того типов тензоров. Определив заполнитель `z` выше, мы указали его тип как `float32`. Наряду со своим типом тензор имеет *форму*. Итак, рассмотрим матрицу 2×3 . Она имеет форму [2, 3]. Вектор длиной 4 имеет форму [4]. Это отличается от матрицы 1×4 , которая имеет форму [1, 4], а матрица 4×1 имеет форму [4, 1]. Массив $3 \times 17 \times 6$ имеет форму [3, 17, 6]. Все они являются тензорами. Скаляры (т.е. числа) имеют нулевую форму и также являются тензорами. Имейте, кроме того, в виду, что тензоры не делают линейно алгебраических различий между векторами строк и векторами столбцов. Есть тензоры, форма которых имеет один компонент, например [5], и это все. То, как мы их рисуем на странице, не имеет значения для математики. Иллюстрируя массив тензоров, мы всегда соблюдаем правило, согласно которому нулевое измерение изображается вертикально, а первое — горизонтально. Но это предел нашей последовательности. (Обратите также внимание на то, что при подсчете тензорные компоненты упоминаются с нуля.)

Возвращаясь к нашему обсуждению заполнителей: большинство заполнителей являются не простыми скалярами наших предыдущих примеров, а многомерными тензорами. Например, следующий раздел начинается с простой программы Tensorflow для распознавания цифр набора Mnist. Основной код TF получает изображение и выполняет прямой проход NN, чтобы получить предположение о том, на какую цифру мы смотрим. Кроме того, на этапе обучения он запускает обратный проход и изменяет параметры программы. Чтобы передать программе изображение, мы определяем заполнитель. Он будет иметь тип `float32` и форму [28,28] или, возможно, [784] в зависимости от того, передали ли мы ему двумерный или одномерный список Python, например

```
img=tf.placeholder(tf.float32, shape=[28, 28])
```

Обратите внимание, что `shape` является именованным аргументом функции `placeholder`.

Прежде чем мы погрузимся в настоящую программу, вот еще одна структура данных TF. Как отмечалось ранее, модели NN определяются своими параметрами и архитектурой программы — как параметры объединяются с входными значениями для получения ответа. Параметры (например, веса `w`, которые соединяют входное изображение с логитами ответов), как правило, инициализируются случайным образом, и NN модифицирует их, чтобы минимизировать потери на обучающих данных. Существует три этапа создания параметров TF. Сначала создайте тензор с начальными значениями. Затем превратите тензор в

`Variable` (то, что TF называет параметрами) и инициализируйте переменные/параметры. Давайте, например, создадим параметры, которые нам нужны для псевдокода Mnist с прямой связью на рис. 1.11. Сначала — слагаемые смещения `b`, а затем — веса `W`:

```
bt = tf.random_normal([10], stddev=.1)
b = tf.Variable(bt)
W = tf.Variable(tf.random_normal([784,10],stddev=.1))
sess=tf.Session()
sess.run(tf.global_variables_initializer())
print(sess.run(b))
```

Первая строка добавляет инструкцию для создания тензора формы [10], десятью значениями которого являются случайные числа, созданные из *нормального распределения* (normal distribution) со стандартным отклонением 0,1. (Нормальное распределение, или *гауссово распределение* (Gaussian distribution), представляет собой известную кривую в форме колокола. Числа, выбранные из нормального распределения, будут сосредоточены вокруг среднего значения (m), и как далеко они удаляются от среднего, определяется *стандартным отклонением* (standard deviation) (s). Конкретнее, около 68% выбранных значений будут в пределах одного стандартного отклонения от среднего, а вероятность того, что они пойдут дальше, быстро уменьшается.)

Вторая строка приведенного выше кода получает `bt` и добавляет фрагмент графа TF, который создает переменную с теми же формой и значениями. Поскольку нам редко нужен исходный тензор после создания переменной, обычно мы объединяем два события, не сохраняя указатель на тензор, как в строке 3, которая создает параметры `W`. Прежде чем использовать `b` или `W`, нужно их инициализировать в сеансе, который мы создали. Это сделано в строке 5. Стока 6 осуществляет вывод (при каждом запуске она будет разной):

```
[ -0.05206999  0.08943175 -0.09178174 -0.13757218  0.15039739
  0.05112269 -0.02723283 -0.02022207  0.12535755 -0.12932496 ]
```

Если бы мы изменили порядок последних двух строк, то получили бы сообщение об ошибке, когда попытались бы вычислить переменную, на которую указывает `b` в команде `print`.

Поэтому в программах TF мы создаем переменные, в которых храним параметры модели. Первоначально их значения неинформативны, обычно случайны с небольшим стандартным отклонением. В соответствии с предыдущим обсуждением обратный проход градиентного спуска их изменяет. После изменения сеанс, на который указывает `sess`, сохраняет новые значения и использует их при следующем запуске сеанса.

2.2. Программа ТF

На рис. 2.2 мы даем (почти) полную программу ТF для программы NN с прямым проходом Mnist. Это должно работать так, как написано. Ключевым элементом, которого вы здесь не видите, является код `mnist.train.next_batch`, который осуществляет чтение данных Mnist. Просто чтобы сориентироваться, обратите внимание, что все, что находится перед пунктирной линией, связано с настройкой графа вычислений ТF; все, что после, граф использует для обучения параметров, а затем запускает программу, чтобы увидеть, насколько она точна на тестовых данных. Теперь пройдем по ней строкой за строкой.

```

0 import tensorflow as tf
1 from tensorflow.examples.tutorials.mnist import input_data
2 mnist = input_data.read_data_sets("MNIST_data/", one_hot=True)
3
4 batchSz=100
5 W = tf.Variable(tf.random_normal([784, 10], stddev=.1))
6 b = tf.Variable(tf.random_normal([10], stddev=.1))
7
8 img=tf.placeholder(tf.float32, [batchSz,784])
9 ans = tf.placeholder(tf.float32, [batchSz, 10])
10
11 prbs = tf.nn.softmax(tf.matmul(img, W) + b)
12 xEnt = tf.reduce_mean(-tf.reduce_sum(ans * tf.log(prbs),
13                                     reduction_indices=[1]))
14 train = tf.train.GradientDescentOptimizer(0.5).minimize(xEnt)
15 numCorrect= tf.equal(tf.argmax(prbs,1), tf.argmax(ans,1))
16 accuracy = tf.reduce_mean(tf.cast(numCorrect, tf.float32))
17
18 sess = tf.Session()
19 sess.run(tf.global_variables_initializer())
20 #-----
21 for i in range(1000):
22     imgs, anss = mnist.train.next_batch(batchSz)
23     sess.run(train, feed_dict={img: imgs, ans: anss})
24
25 sumAcc=0
26 for i in range(1000):
27     imgs, anss= mnist.test.next_batch(batchSz)
28     sumAcc+=sess.run(accuracy, feed_dict={img: imgs, ans: anss})
29 print "Test Accuracy: %r" % (sumAcc/1000)
```

Рис. 2.2. Код Tensorflow для программы NN с прямым проходом Mnist

После импорта Tensorflow и кода для чтения данных Mnist мы определяем наши два набора параметров в строках 5 и 6. Это незначительный вариант того,

что мы только что видели в нашем обсуждении переменных TF. Затем мы создаем заполнители для данных, которые мы передаем в NN. Вначале, в строке 8, у нас есть заполнитель для данных изображения. Это тензор формы [batchSz, 784]. В обсуждении того, почему линейная алгебра была хорошим средством представления вычислений NN (стр. 32), мы отметили, что наши вычисления ускоряются, когда мы обрабатываем несколько примеров одновременно, а кроме того, это хорошо согласуется с понятием размера партии стохастического градиентного спуска. Здесь мы видим, как это работает в TF, а именно — наш заполнитель для изображения занимает не одну строку из 784 пикселей, а 100 из них (поскольку это значение batchSz). Аналогично в строке 9 мы видим, что мы даем программе 100 ответов на изображение одновременно.

Еще один момент о строке 9. Мы представляем ответ вектором длиной 10 со всеми значениями 0, кроме a -го, где a — правильная цифра для этого изображения. Например, мы начали главу 1 с изображения 7 (см. рис. 1.1). Соответствующее представление правильного ответа — $(0, 0, 0, 0, 0, 0, 0, 1, 0, 0)$. Векторы этой формы называют *унитарными векторами* (one-hot vector), поскольку позволяют выбирать только одно значение в качестве активного.

Строка 9 завершается параметрами и исходными значениями нашей программы, и код переходит к размещению в графе фактических вычислений. В частности, строка 11 начинает показывать мощь TF для вычислений NN. Она определяет большую часть прямого прохода NN модели. В частности, она указывает, что мы хотим передать (размер партии) изображения в наши линейные блоки (как определено W и b), а затем применить softmax ко всем результатам, чтобы получить вектор вероятностей.

При просмотре такого кода мы рекомендуем сначала проверять формы задействованных тензоров, чтобы убедиться, что они разумны. Здесь самым глубоко вложенным вычислением является умножение матриц `matmul` входных изображений [100, 784] на W [784, 10], чтобы получить матрицу формы [100, 10], к которой мы добавляем смещения, завершающиеся матрицей формы [100, 10]. Это 10 логитов для 100 изображений в нашей партии. Затем мы пропускаем это через функцию softmax и в итоге получаем матрицу [100, 10] вероятностей назначений меток для наших изображений.

В строке 12 вычисляется средняя кросс-энтропийная потеря за 100 примеров, которые мы обрабатываем параллельно. Работая изнутри, функция `tf.log(x)` возвращает тензор с каждым элементом x , замененным его натуральным логарифмом. На рис. 2.3 мы демонстрируем работу функции `tf.log` с размером партии из трех векторов, каждый из которых содержит распределение вероятностей из пяти элементов.

$$\begin{array}{ccccccccc} 0,20 & 0,10 & 0,20 & 0,10 & 0,40 & & -1,6 & -2,3 & -1,6 & -2,3 & -0,9 \\ 0,20 & 0,10 & 0,20 & 0,10 & 0,40 & \rightarrow & -1,6 & -2,3 & -1,6 & -2,3 & -0,9 \\ 0,20 & 0,10 & 0,20 & 0,10 & 0,40 & & -1,6 & -2,3 & -1,6 & -2,3 & -0,9 \end{array}$$

Рис. 2.3. Работа функции tf.log

Далее, стандартный символ умножения “*” в `ans * tf.log(prbs)` выполняет поэлементное умножение двух тензоров. На рис. 2.4 показано, как поэлементное умножение унитарных векторов для каждой метки в пакетном режиме отрицательной матрицы натуральных логарифмов создает строки, в которых содержатся только нули, за исключением отрицательного логарифма вероятности правильного ответа.

$$\begin{array}{ccccccccc} 0 & 0 & 1 & 0 & 0 & & 1,6 & 2,3 & 1,6 & 2,3 & 0,9 & & 0 & 0 & 1,6 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & * & 1,6 & 2,3 & 1,6 & 2,3 & 0,9 & = & 0 & 0 & 1,6 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & & 1,6 & 2,3 & 1,6 & 2,3 & 0,9 & & 0 & 0 & 0 & 0 & 0,9 \end{array}$$

Рис. 2.4. Расчет отрицательного логарифма вероятности правильного ответа

На данный момент, чтобы получить кросс-энтропию для каждого изображения, просто нужно сложить все значения в массиве. Первая операция, к которой мы обращаемся, — это функция

```
tf.reduce_sum( A, reduction_indices = [1]),
```

Она суммирует строки `A`, как на рис. 2.5. Критически важная часть здесь такова:

```
reduction_indices = [1].
```

При знакомстве с тензорами ранее, мы мимоходом упоминали, что в измерениях тензоров используется нумерация с нуля. Теперь `reduce_sum` может суммировать по столбцам (стандартно) с помощью `reduction_index=[0]`¹ или, как в этом случае, суммировать по строкам, `reduction_index=[1]`. В результате получается массив `[100,1]` с логарифмом правильной вероятности в качестве единственной записи в каждой строке. (На рис. 2.5 используется размер партии только 3 и предполагает 5 альтернативных вариантов, а не 10.) В качестве последней части вычисления кросс-энтропии функция `reduce_mean` в строке 13 на рис. 2.2 суммирует все столбцы (снова стандартно) и возвращает среднее (1,1 или около того).

¹ Вероятно, имелось в виду `reduction_indices`. — Примеч. ред.

$$\begin{array}{ccccccccc}
 0 & 0 & 1,6 & 0 & 0 & & 1,6 \\
 0 & 0 & 1,6 & 0 & 0 & \rightarrow & 1,6 \\
 0 & 0 & 0 & 0 & 0,9 & & 0,9
 \end{array}$$

Рис. 2.5. Вычисление `tf.reduce_sum` при `reduction_index=[1]`

Наконец, мы можем перейти к строке 14 на рис. 2.2, и именно там TF действительно показывает свои достоинства: эта единственная строка — все, что нам нужно, чтобы включить весь обратный проход:

```
tf.train.GradientDescentOptimizer(0.5).minimize(xEnt)
```

Это говорит о необходимости расчета изменений веса с использованием градиентного спуска и минимизации функции кросс-энтропийных потерь, определенной в строках 12 и 13. Она также определяет скорость обучения 0,5. Нам не нужно беспокоиться о вычислении производных или о чем-либо еще. Если выразить и прямое вычисление, и потерю в TF, то компилятор TF будет знать, как вычислить необходимые производные и связать их вместе в правильном порядке, чтобы внести изменения. Мы можем изменить этот вызов функции, выбрав другую скорость обучения или, если у нас была другая функция потерь, заменить `xEnt` чем-то, что указывало бы на другое вычисление TF.

Естественно, существуют ограничения на способность TF получать обратный проход на основе прямого прохода. Повторим, что это можно сделать только в том случае, если все вычисления прямого прохода выполняются функциями TF. Для начинающих, таких как мы, это не слишком большое ограничение, так как TF имеет множество встроенных операций, которые он умеет различать и объединять.

Строки 15 и 16 вычисляют точность модели (`accuracy`), т.е. подсчитывают количество правильных ответов и делят на количество обработанных изображений. Сначала сосредоточимся на стандартной математической функции `argmax`, как в $\arg \max_x f(x)$, возвращающая значение x , которое максимизирует $f(x)$. В нашем случае `tf.argmax(prbs, 1)` получает два аргумента. Первый — это тензор, по которому мы берем `argmax`. Вторая — это ось (`axis`) тензора для использования в функции `argmax`. Это работает как именованный аргумент, который мы использовали для `reduce_sum`, — он позволяет суммировать по разным осям тензора. Например, если тензор равен $((0, 2, 4), (4, 0, 3))$ и мы используем ось 0 (стандартно), мы возвращаем $(1, 0, 0)$. Сначала мы сравнили 0 с 4 и вернули 1, так как 4 было больше. Затем мы сравнили 2 с 0 и вернули 0, так как 2 было больше. Если бы мы использовали ось 1, мы бы вернули $(2, 0)$. В строке 15 у нас есть массив размера партии логитов. Функция `argmax` возвращает массив размера партии логитов. Затем мы применяем функцию `tf.equal`, чтобы сравнить максимальные логиты с правильным ответом. Он возвращает

вектор размера партии логических значений (True, если они равны), который превращает `tf.cast(tensor, tf.float32)` в числа с плавающей запятой, так что `tf.reduce_mean` может сложить их и получить правильный процент. Не преобразуйте логические значения в целые числа, так как при получении среднего значения оно возвращает также целое число, которое в этом случае всегда будет равно нулю.

Далее, как только мы определили наш сеанс (строка 18) и инициализировали значения параметров (строка 19), мы можем обучить модель (строки с 21 по 23). Здесь мы используем код, полученный из библиотеки TF Mnist, для одновременного извлечения 100 изображений и их ответов, а затем запускаем их, вызывая функцию `sess.run` для фрагмента графа вычислений, на который указывает `train`. По завершении этого цикла мы прошли 1000 итераций обучения по 100 изображений на итерацию или всего 100 000 тестовых изображений. На моем четырехядерном Mac Pro это занимает около 5 секунд (или меньше, если сначала поместить нужные вещи в кеш). Я упоминаю “четыре ядра”, поскольку TF учитывает доступную вычислительную мощность и в целом хорошо ее использует.

Обратите внимание на одну немного странную вещь в строках с 21 по 23: мы никогда не упоминаем явно о выполнении прямого прохода! TF также выясняет это, основываясь на графике вычислений. Из `GradientDescentOptimizer` он знает, что ему необходимо выполнить вычисления, на которые указывает `xEnt` (строка 12), что требует вычисления `probs`, а это, в свою очередь, определяет вычисления прямого прохода в строке 11.

Наконец, строки с 25 по 29 показывают, насколько хорошо мы справляемся с тестовыми данными с точки зрения процентного соотношения (91% или 92%). Вначале, просто взглянув на организацию графа, обратите внимание, что вычисление `accuracy` в конечном итоге требует прямого прохода при вычислении `probs`, но не обратного прохода `train`. Таким образом, как и следовало ожидать, весовые коэффициенты не модифицируются для лучшей работы с тестовыми данными.

Как упоминалось в главе 1, вывод частоты ошибок при обучении модели является хорошей практикой при отладке. Как правило, она уменьшается. Для этого изменим строку 23 на

```
acc,ignore= sess.run([accuracy,train],  
                    feed_dict={img: imgs, ans: anss})
```

Это обычный синтаксис Python для объединения вычислений. Значение первого вычисления (для `accuracy`) присваивается переменной `acc`, а для второго — `ignore`. (Распространенной идиомой Python является замена `ignore` символом подчеркивания (`_`), универсальным символом Python, используемым,

когда синтаксис требует, чтобы переменная получала значение, но нам не нужно его запоминать.) Естественно, нам также необходимо добавить команду для вывода значения acc.

Мы упомянули об этом, чтобы предупредить читателя о распространенной ошибке (по крайней мере, ваш автор и некоторые из его начинающих учеников сделали ее). Ошибка в том, чтобы оставить строку 23 одинокой и добавить новую строку 23.5:

```
acc= sess.run(accuracy, feed_dict={img: imgs, ans: anss}).
```

Это, однако, менее эффективно, поскольку TF теперь делает прямой проход дважды: один раз — когда мы указываем ему обучаться, и один раз — когда мы запрашиваем его о точности. Есть также более важная причина, чтобы избежать этой ситуации. Обратите внимание, что первый вызов изменяет веса и, таким образом, повышает вероятность правильной метки для этого изображения. После запроса о точности программист получает преувеличенное представление о том, насколько хорошо работает программа. Когда у нас только один вызов sess.run, но мы запрашиваем оба значения, этого не происходит.

2.3. Многослойные NN

Программа, которую мы разработали в псевдокоде в главе 1 и только сейчас в TF, является однослойной. Существует только один слой линейных блоков. Естественный вопрос — можем ли мы добиться большего успеха с несколькими слоями таких блоков? Ранее исследователи NN поняли, что ответ — “Нет”. Это следует почти сразу после того, как мы видим, что линейные блоки могут быть преобразованы в матрицы линейной алгебры, т.е. как только мы увидим, что однослойная NN с прямой связью просто вычисляет $y = \mathbf{X}\mathbf{W}$. В нашей модели Mnist \mathbf{W} имеет форму [784, 10], чтобы преобразовать значения 784 пикселей в 10 значений логитов и добавить дополнительный вес для замены члена смещения. Предположим, мы добавили дополнительный слой линейных блоков \mathbf{U} с формой [784, 784], который, в свою очередь, передает на слой \mathbf{V} той же формы, что и \mathbf{W} , [784, 10]:

$$\mathbf{y} = (\mathbf{x}\mathbf{U})\mathbf{V} \quad (2.1)$$

$$= \mathbf{x}(\mathbf{U}\mathbf{V}) \quad (2.2)$$

Вторая строка следует из ассоциативного свойства умножения матриц. Дело в том, что любые возможности, полученные в двухслойной ситуации комбинацией \mathbf{U} с последующим умножением на \mathbf{V} , могут быть получены однослойным NN при $\mathbf{W} = \mathbf{UV}$.

Оказывается, есть простое решение — добавить некоторые нелинейные вычисления между слоями. Одним из наиболее часто используемых является `tf.nn.relu` (или `r`), который обозначает *блок линейной ректификации* (Rectified Linear Unit — ReLU), или *линейный выпрямитель*. Он определяется как

$$\rho(x) = \max(x, 0) \quad (2.3)$$

и показан на рис. 2.6.

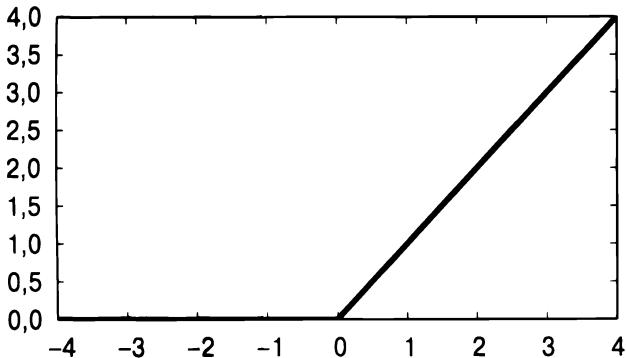


Рис. 2.6. Поведение `tf.nn.relu`

Нелинейные функции, помещаемые между слоями в глубоком обучении, называются *функциями активации* (activation function). В то время как блок ReLU в настоящее время довольно популярен, другие также используются, например *сигмовидная функция* (sigmoid function) определена как

$$S(x) = \frac{e^{-x}}{1 + e^{-x}} \quad (2.4)$$

и показана на рис. 2.7. Во всех случаях активация применяется к отдельным действительным числам в аргументе тензора, например $\rho([1, 17, -3]) = [1, 17, 0]$.

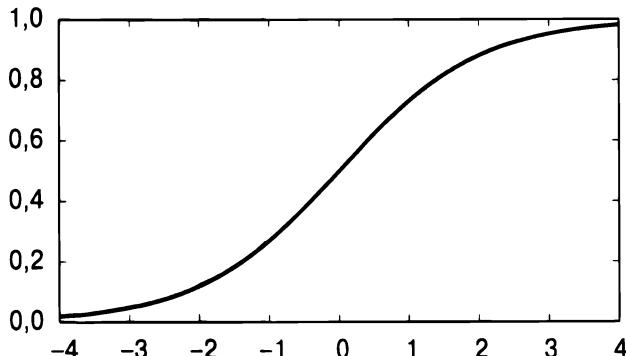


Рис. 2.7. Сигмовидная функция

До того как было обнаружено, что даже такая простая нелинейность, как ReLU, будет работать, сигмовидная функция стала очень популярной. Однако диапазон значений, которые может выдавать сигмовидная функция, весьма ограничен, от нуля до единицы, тогда как ReLU имеет диапазон от нуля до бесконечности. Это очень важно, когда мы делаем обратный проход, чтобы найти градиент влияния параметров на потери. Обратное распространение через функции с ограниченным диапазоном значений может сделать градиент практически нулевым — процесс, известный как проблема *исчезающего градиента* (vanishing gradient). Простые функции активации очень помогли. По этой причине `tf.nn.lrelu`, *ReLU с утечкой* (leaky relu), также используется очень часто, поскольку имеет еще более широкий диапазон значений, чем ReLU, как видно на рис. 2.8.

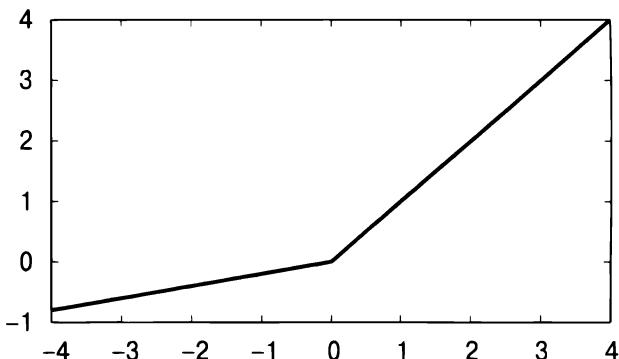


Рис. 2.8. Функция *lrelu*

Собрав фрагменты нашего многослойного NN вместе, получаем новую модель:

$$\text{Pr}(A(x)) = \sigma(\rho(\mathbf{x}\mathbf{U} + \mathbf{b}_u)\mathbf{V} + \mathbf{b}_v) \quad (2.5)$$

Здесь σ — это функция softmax, \mathbf{U} и \mathbf{V} — веса первого и второго слоев линейных блоков, а \mathbf{b}_u и \mathbf{b}_v — их смещения.

Давайте теперь сделаем это в TF. На рис. 2.9 мы заменили определения \mathbf{W} и \mathbf{b} в строках 5 и 6 на рис. 2.2 двумя слоями, \mathbf{U} и \mathbf{V} , строки с 1 по 4 на рис. 2.9. Мы также заменим вычисление `prbs` в строке 11 на рис. 2.2 строками 5 и 7 на рис. 2.9. Это превращает наш код в многослойный NN. (Кроме того, чтобы отразить большее количество параметров, нужно снизить коэффициент обучения в 10 раз.) Хотя старая программа показала точность около 92% после обучения на 100 000 изображений новая достигает точности около 94% на 100 000 изображениях. Кроме того, если мы увеличим количество обучающих изображений, производительность тестового набора будет увеличиваться примерно до 97%. Обратите внимание, что единственная разница между этим кодом и кодом

без нелинейной функции — это строка 6. Если мы удалим его, производительность действительно снизится примерно до 92%. Этого достаточно, чтобы заставить вас поверить в математику!

```

1 U = tf.Variable(tf.random_normal([784, 784], stddev=.1))
2 bU = tf.Variable(tf.random_normal([784], stddev=.1))
3 V = tf.Variable(tf.random_normal([784, 10], stddev=.1))
4 bV = tf.Variable(tf.random_normal([10], stddev=.1))
5 L1Output = tf.matmul(img, U) + bU
6 L1Output = tf.nn.relu(L1Output)
7 prbs = tf.nn.softmax(tf.matmul(L1Output, V) + bV)

```

Рис. 2.9. Код построения графа TF для многоуровневого распознавания цифр

Еще один момент. В однослойной сети с параметрами массива W форма W фиксируется количеством входов с одной стороны (784) и количеством возможных выходов — с другой (10). С двумя слоями мы можем сделать еще один выбор — *скрытый размер* (hidden size). Таким образом, U — это размер ввода по скрытому размеру, а V — скрытый размер по выходному размеру. На рис. 2.9 мы просто установили скрытый размер 784, тот же самый, что и размер ввода, но ничего не требовало такого выбора. Как правило, его увеличение улучшает производительность, но до определенного предела.

2.4. Другие части

В этом разделе мы рассмотрим аспекты TF, которые очень полезны при выполнении заданий по программированию, предложенных в остальной части книги (например, *контрольные точки* (checkpointing)), и которые мы будем использовать в следующих главах.

2.4.1. Создание контрольных точек

Зачастую в вычисления TF полезно создавать контрольные точки — сохранять тензоры в вычислении, чтобы вычисление могло быть возобновлено позже или для повторного использования в другой программе. В TF мы делаем это, создавая и используя *хранимые объекты* (saver object):

```
saveOb = tf.train.Saver()
```

Как и ранее, `saveOb` — это переменная Python, и выбор имени за вами. Объект может быть создан в любое время до его использования, но по причинам, объясненным ниже, сделать это непосредственно перед инициализацией переменных (вызвав `global_variable_initialize`) вполне логично. Затем после каждой n эпох обучения сохраняйте текущие значения всех своих переменных:

```
saveOb.save(sess, "mylatest.ckpt")
```

Метод `save` получает два аргумента: сеанс, который нужно сохранить, а также имя и местоположение файла. В приведенном выше случае информация находится в том же каталоге, что и программа Python. Если бы аргументом было `tmp/model.checkpt`, он был бы помещен в подкаталог `tmp`.

Вызов метода `save` создает четыре файла. Наименьший файл, по имени `checkpoint`, представляет собой файл Ascii, в котором указаны некоторые подробные сведения о контрольной точке, сделанной для этого каталога. Имя `checkpoint` является фиксированным. Если вы назовете один из своих файлов “`checkpoint`”, он будет переписан. Три других файла используют имена, заданные в методе `save`. В данном случае они названы

```
mylatest.ckpt.data-00000-of-00001  
mylatest.ckpt.index  
mylatest.chpt.meta
```

Первый из них содержит значения сохраняемых параметров. Два других содержат метаинформацию, которую TF использует, когда вы захотите импортировать эти значения (это будет описано ниже). Если ваша программа вызывает метод `save` несколько раз, эти файлы каждый раз переписываются.

Далее мы хотим, скажем, провести дальнейшее обучение по той же модели NN, которую мы уже начали обучать. Самое простое, что нужно сделать, — это изменить оригинальную программу обучения. Вы сохранили хранимые объекты, но теперь мы хотим инициализировать все переменные TF сохраненными ранее значениями. Таким образом, обычно удаляют вызов `global_variable_initialize` и заменяют его вызовом метода `restore` нашего хранимого объекта:

```
saveOb.restore(sess, "mylatest.ckpt")
```

В следующий раз, когда вы вызываете программу обучения, она возобновляет работу с переменными TF, значения которых они имели во время последнего сохранения при предыдущем обучении. Однако больше ничего не меняется. Итак, если ваш обучающий код вывел, скажем, номер эпохи, следующий за потерей, на этот раз он выводит номера эпох, начиная с единицы, если вы не переписали свой код так, чтобы поступить иначе. (Естественно, вы можете, если хотите, это исправить или вообще сделать все более элегантно, но написание лучшего кода на языке Python не является здесь нашей главной задачей.)

2.4.2. tensordot

`tensordot` — это обобщение умножения матриц на тензоры. Мы знакомы со стандартным матричным умножением `matmul` из предыдущей главы. Мы можем вызвать `tf.matmul(A, B)`, когда A и B имеют одинаковое количество измерений, скажем, n , последнее измерение A имеет тот же размер, что и второе последнее измерение B, и первые $n - 2$ измерений будут идентичны. Таким образом, если размеры A равны [2, 3, 4], а размеры B — [2, 4, 6], то размеры произведения равны [2, 3, 6]. Матричное умножение можно рассматривать как получение повторяющихся скалярных произведений. Например, умножение матриц

$$\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} -1 & -2 \\ -3 & -4 \\ -5 & -6 \end{pmatrix} \quad (2.6)$$

можно осуществить, взяв скалярное произведение векторов $\langle 1, 2, 3 \rangle$ и $\langle -1, -3, -5 \rangle$ и поместив ответ в верхнюю левую позицию в матрице результатов. Продолжая таким образом, мы берем скалярное произведение i -й строки с j -м столбцом, и это есть i, j -е значение в ответе. Таким образом, если A является первой из вышеуказанных матриц, а B — второй, это вычисление также может быть выражено как

```
tf.tensordot(A, B, [[1], [0]])
```

Первые два аргумента, конечно, являются тензорами, с которыми мы работаем. Третий аргумент представляет собой двухэлементный список: первый элемент — это список размерностей из первого аргумента, а второй — соответствующий список из второго аргумента. Это указывает `tensordot` взять скалярные произведения этих двух измерений. Естественно, указанные размерности должны иметь одинаковый размер, если мы хотим получить их скалярное произведение. Поскольку 0-е измерение — это то, что мы выводим в виде строк, а 1-е — как столбцы, это говорит о том, что нужно взять скалярное произведение каждой из строк A с каждым из столбцов B. `tensordot` размещает выходные размерности в порядке слева направо, начиная с размерности A, а затем переходя к размерности B. То есть в этом случае у нас есть входные размерности [2, 3], за которыми следуют [3, 2]. Два измерения, включенные в скалярное произведение, “исчезают” (измерения 1 и 0), давая ответ с размерностями [2, 2].

На рис. 2.10 приведен более сложный пример, который `matmul` не может обработать в одной инструкции. Мы заимствовали его из главы 5, в которой имена переменных будут иметь смысл. Здесь мы смотрим на это, чтобы увидеть, что делает `tensordot`. Не глядя на числа, просто посмотрите на третий аргумент в

вызове функции `tensordot`, $\begin{bmatrix} 1 \end{bmatrix} \begin{bmatrix} 0 \end{bmatrix}$. Это означает, что мы берем скалярное произведение из 1 размерности `encOut` и 0 размерности `AT`. Это вполне законно, так как они оба имеют размер 4. Т.е. мы берем скалярное произведение двух тензоров с размерностями $[2, 4, 4]$ и $[4, 3]$ соответственно. (Цифры, выделенные курсивом, являются размерностями, подвергающимися скалярному произведению.) Поскольку эти размерности уходят после скалярного произведения, результирующий тензор имеет размерности $[2, 4, 3]$, которые при их выводе в нижней части примера верны. Согласно элементарной арифметике мы получаем скалярное произведение столбцов двух тензоров. То есть первое скалярное произведение находится между $[1, 1, 1, -1]$ и $[.6, .2, .1, .1]$. Результат, .8, появляется как первое числовое значение в результирующем тензоре.

```

eo= (( 1, 2, 3, 4),
      ( 1, 1, 1, 1),
      ( 1, 1, 1, 1),
      (-1, 0,-1, 0)),
      (( 1, 2, 3, 4),
      ( 1, 1, 1, 1),
      ( 1, 1, 1, 1),
      (-1, 0,-1, 0))
encOut=tf.constant(eo, tf.float32)

AT = ( (.6, .25, .25 ),
       (.2, .25, .25 ),
       (.1, .25, .25 ),
       (.1, .25, .25 ))
wAT = tf.constant(AT, tf.float32)

encAT = tf.tensordot(encOut,wAT, [1], [0])
sess= tf.Session()

print sess.run(encAT)
[[ 0.80000001   0.5    0.5     ]
 [ 1.50000012   1.      1.      ]
 [ 2.          1.      1.      ]
 [ 2.70000005   1.5    1.5     ]]
 ...

```

Рис. 2.10. Пример `tensordot`

Наконец, `tensordot` не ограничивается получением скалярного произведения одного измерения от каждого тензора. Если размерности A равны $[2, 4, 4]$, а таковые у B — $[4, 4]$, то операция `tensordot(A, B, [[1,2],[0,1]])` приводит к тензору размерности $[2]$.

2.4.3. Инициализация переменных TF

В разделе 1.4 мы сказали, что, как правило, хорошей практикой является инициализация параметров NN (т.е. переменных TF) случайными, но близкими к нулю значениями. В нашей первой программе TF (рис. 2.9) мы отменили этот запрет с помощью такой команды, как:

```
b = tf.Variable(tf.random_normal([10], stddev=.1))
```

И мы предположили, что стандартное отклонение 0,1 было достаточно “близко к нулю”.

Тем не менее существуют теория и практика выбора стандартных отклонений в этих случаях. Здесь мы представляем правило *инициализации Ксавье* (Xavier initialization). Обычно оно используется для установки стандартного отклонения при случайной инициализации переменных. Пусть n_i будет количеством соединений, входящих в слой, а n_o — количеством исходящих. Для переменной W на рис. 2.9 $n_i = 784$ — количество пикселей, а $n_o = 10$ — количество альтернативных классификаций. Для инициализации Ксавье мы устанавливаем стандартное отклонение с следующим образом:

$$\sigma = \sqrt{\frac{2}{n_i + n_o}} \quad (2.7)$$

Например, для W , поскольку рассматриваемые значения равны 784 и 10, мы получаем $\sigma \approx 0,0502$, что округляется до 0,1. Обычно рекомендуемые стандартные отклонения могут варьироваться от 0,3 для слоя $10 * 10$ до 0,03 для одного из слоев $1000 * 1000$. Чем больше входных и выходных значений, тем ниже стандартное отклонение.

Инициализация Ксавье была первоначально предназначена для использования с сигмовидной функцией активации (см. рис. 2.7). Как отмечалось ранее, $\sigma(x)$ становится относительно невосприимчивым к x , когда x значительно ниже -2 или выше +2. То есть если значения, подаваемые в сигмовидную функцию, слишком высоки или слишком низки, изменение в них может практически не влиять на значения потерь. Если двигаться в обратном направлении, то при обратном проходе изменения в потерях не будут влиять на параметры, которые поступают в сигмовидную функцию, если изменение в потерях удалено сигмовидной функцией. Вместо этого мы хотим, чтобы дисперсия отношения между входом и выходом слоя была около единицы. Здесь мы используем дисперсию в ее техническом смысле: ожидаемое значение квадрата разности между значением числовой случайной величины и ее средним значением. Кроме того, *ожидаемое значение (expected value)* случайной величины X (обозначается $E[X]$) является вероятностным средним ее возможных значений:

$$E[X] = \sum_x p(X=x) * x. \quad (2.8)$$

Стандартным примером является ожидаемое значение броска честного шестигранного кубика:

$$E[R] = \frac{1}{6} * 1 + \frac{1}{6} * 2 + \frac{1}{6} * 3 + \frac{1}{6} * 4 + \frac{1}{6} * 5 + \frac{1}{6} * 6 = 3,5 \quad (2.9)$$

Поэтому мы хотим сохранить отношение входной дисперсии к выходной дисперсии примерно около 1, чтобы слой не способствовал чрезмерному ослаблению сигнала сигмовидной функцией. Это накладывает ограничения на инициализацию. Мы имеем тот факт (вы можете посмотреть вывод), что для линейного блока с весовой матрицей \mathbf{W} дисперсия в прямом проходе (V_f) и обратном проходе (V_b) соответственно составляют:

$$V_f(\mathbf{W}) = \sigma^2 \cdot n_i \quad (2.10)$$

$$V_b(\mathbf{W}) = \sigma^2 \cdot n_o \quad (2.11)$$

Здесь σ — это стандартное отклонение весов \mathbf{W} . (Это имеет смысл, учитывая, что дисперсия одного гауссиана равна (σ^2) .) Если установить оба значения, V_f и V_b , равными нулю и решить для σ , получится

$$\sigma = \sqrt{\frac{1}{n_i}} \quad (2.12)$$

$$\sigma = \sqrt{\frac{1}{n_o}}. \quad (2.13)$$

Естественно, это не имеет решения, если количество входов не совпадает с количеством выходов. Поскольку чаще всего это *не так*, мы берем “среднее” между двумя значениями согласно правилу Ксавье

$$\sigma = \sqrt{\frac{2}{n_i + n_o}} \quad (2.14)$$

Эквивалентные уравнения есть и для других функций активации. С появлением блоков ReLU и других функций активации, которые не насыщаются так же легко, как сигмовидная, проблема не так важна, как это было раньше. Тем не менее правило Ксавье дает хорошее представление о том, каким должно быть стандартное отклонение, и его TF версии часто используются наряду с его родственниками.

2.4.4. Упрощение создания графов TF

Возвращаясь к рис. 2.9, мы видим, что потребовалось семь строк кода, чтобы создать нашу двухслойную сеть прямой связи. В общей схеме это не так много — подумайте, что потребовалось бы, если бы мы программировали на языке Python без TF. Однако, если бы мы создавали, скажем, восьмислойную сеть (и к концу этой книги вы будете делать именно это), для этого потребовалось бы 24 строки кода или около того.

TF имеет удобную группу функций, модуль `layers`, для более компактного кодирования очень распространенных многослойных ситуаций. Здесь мы представляем

```
tf.contrib.layers.fully_connected.
```

Слой считается *полносвязанным* (*fully connected*), если все его элементы связаны со всеми модулями в следующем слое. Все слои, которые мы используем в первых двух главах, полносвязаны между собой, поэтому нет необходимости проводить различия между ними и сетями, в которых это свойство отсутствует. Чтобы определить такой слой, мы обычно делаем следующее: а) создаем веса **W**, б) создаем смещения **b**, в) осуществляя умножение матриц и добавляем смещения и наконец г) применяем функцию активации. Предполагая, что мы импортировали `tensorflow.contrib.layers` как `layers`, все это можно сделать одной строкой:

```
layerOut=layers.fully_connected(layerIn,outSz,activeFn)
```

Вышеуказанный вызов создает матрицу, инициализированную инициализацией Ксавье, и вектор смещений, инициализированных нулем. Он возвращает `layerIn`, умноженную на матрицу, плюс смещения, к которым была применена функция активации, заданная `activeFn`. Если вы не укажете функцию активации, она использует `relu`. Если в качестве функции активации вы указываете `None`, то активация не используется.

С помощью `fully connected` мы можем записать все семь строк на рис. 2.9 как

```
L1Output=layers.fully_connected(img,756)
prbs=layers.fully_connected(L1Output,10,tf.nn.softmax)
```

Обратите внимание, что мы указали использование `tf.nn.softmax` для применения к выходу второго уровня, используя его в качестве функции активации для второго уровня.

Конечно, если у нас есть NN, состоящий из ста слоев, и это происходит, даже если писать 100 вызовов `fully_connected` утомительно. К счастью, мы можем использовать Python или любой другой API TF для спецификации нашей сети.

Приведем несколько необычный пример. Предположим, что мы хотим создать 100 скрытых слоев, каждый на 1 меньше предыдущего, и размер первого является системным параметром. Мы могли бы написать

```
outpt = input
for i in range(100):
    outpt = layers.fully_connected(outpt, sysParam - i)
```

Этот пример глуп, но суть серьезна: фрагменты графов TF могут передаваться и обрабатываться в Python, как списки или словари.

2.5. Ссылки и что читать дальше

Tensorflow был запущен в рамках исследовательского проекта *Google Brain* компании Google, созданным двумя ее исследователями, Джейфом Дином и Грегом Коррадо, а также профессором Стэнфорда Эндрю Ын. В то время он назывался “DistBelief”. Когда его использование вышло за рамки этого одного проекта, компания Google собственно взялась за дальнейшее развитие и наняла Джейфри Хинтона из Университета Торонто, которого мы упомянули в главе 1 за его новаторский вклад в глубокое обучение.

Инициализация Ксавье берет свое имя от имени Ксавье Глоро (Xavier Glorot), первого автора [9], который представил эту технику.

В наши дни Tensorflow является лишь одним из многих языков программирования, нацеленных на глубокое обучение (например, [10]). С точки зрения количества пользователей язык Tensorflow является, безусловно, самым популярным. Далее — Keras, язык более высокого уровня, построенный поверх Tensorflow, занимает второе место, за ним следует Caffe, первоначально разработанный в Калифорнийском университете в Беркли. Facebook теперь поддерживает версию Caffe с открытым исходным кодом Caffe2. Pytorch — это интерфейс Python для Torch, языка, который завоевал популярность в сообществе специалистов по обработке естественных языков.

2.6. Упражнения

Упражнение 2.1. Каков был бы результат, если бы на рис. 2.5 мы вместо этого вычислили `tf.reduce_sum(A)`, где A — массив слева от фигуры?

Упражнение 2.2. Что плохого в том, чтобы взять строку 14 из рис. 2.2 и вставить ее между строками 22 и 23, чтобы цикл теперь выглядел следующим образом:

```
for i in range(1000):
    imgs, anss = mnist.train.next_batch(batchSz)
    train = tf.train.GradientDescentOptimizer(0.5).minimize(xEnt)
    sess.run(train, feed_dict={img: imgs, ans: anss})
```

Упражнение 2.3. Вот еще один вариант в тех же строках кода. Это нормально? Если нет, то почему?

```
for i in range(1000):
    img, anss= mnist.test.next_batch(batchSz)
    sumAcc+=sess.run(accuracy, feed_dict={img:img, ans:anss})
```

Упражнение 2.4. Какова будет форма тензорного вывода операции на рис. 2.10

`tensordot(wAT, encOut, [[0],[1]]) ?`

Объясните.

Упражнение 2.5. Покажите вычисление, которое подтверждает, что первое число в тензоре, распечатанном в нижней части примера на рис. 2.10 (0,8), является правильным (до трех мест).

Упражнение 2.6. Предположим, что `input` имеет форму [50,10]. Сколько переменных TF создается следующим кодом:

`O1 = layers.fully_connected(input, 20, tf.sigmoid) ?`

Каково будет стандартное отклонение для переменных в созданной матрице?

Глава 3

Сверточные нейронные сети

Рассмотренные до сих пор NN были полносвязными. То есть они обладают тем свойством, что все линейные блоки в слое связаны со всеми линейными блоками в следующем слое. Однако вовсе не требуется, чтобы NN имели именно эту форму. Мы, конечно, можем представить себе прямой проход, при котором линейный блок передает свой выходной сигнал только некоторым блокам следующего слоя. Немного сложнее, но не чрезмерно, увидеть, что, скажем, Tensorflow, если он знает, какие блоки связаны с какими, может правильно вычислить производные веса на обратном проходе.

Один частный случай частично связанных NN — это *сверточные нейронные сети* (convolutional neural network). Сверточные NN особенно часто применяются в компьютерном зрении, поэтому мы продолжаем обсуждение набора данных Mnist.

Одноуровневая полносвязная NN Mnist учится связывать определенные интенсивности света в определенных позициях на изображении с определенными цифрами, например высокие значения в позиции (8, 14) с числом 1. Но это явно не то, как работают люди. Фотографирование цифр в более яркой комнате может добавить 10 к каждому значению пикселя, но это мало повлияет на нашу классификацию, поскольку при распознавании сцены важны различия в значениях пикселей, а не их абсолютные значения. Кроме того, различия имеют смысл только между близкими значениями. Предположим, вы находитесь в маленькой комнате, освещенной одной лампочкой в одном углу комнаты. То, что мы воспринимаем как светло пятно, скажем, на некоторых обоях в противоположном конце комнаты, может фактически отражать меньше фотонов, чем “темное” пятно рядом с лампой. Для определения того, что происходит на сцене, важны локальные различия интенсивности света с акцентом на “местные” и “различия”. Естественно, исследователи компьютерного зрения вполне осведомлены об этом, и стандартным ответом на эти факты было почти повсеместное применение сверточных методов.¹

¹ В этом обсуждении термин “свертка” употребляется так, как в глубоком обучении. Это близко, но не совсем то же самое, что и его использование в математике, где глубокая свертка называется *взаимной корреляцией* (cross-correlation).

3.1. Фильтры, шаги и дополнения

Для наших целей *сверточный фильтр* (convolutional filter) (называемый также *сверточным ядром* (convolutional kernel)) представляет собой (как правило, небольшой) массив чисел. Если мы имеем дело с черно-белым изображением, то это двумерный массив. Набор данных Mnist — черно-белый, так что это все, что нам здесь нужно. Если бы у нас был цвет, потребовался бы трехмерный массив — или, что эквивалентно, три двумерных массива — по одному для длин волн красного, синего и зеленого (RGB) света, из которых можно создать все цвета. На данный момент мы игнорируем сложности цвета. Мы вернемся к ним позже.

Рассмотрим сверточный фильтр, показанный на рис. 3.1. Для *свертки* (convolve) фильтра с фрагментом изображения мы берем скалярное произведение фильтра и фрагмент изображения равного размера. Следует помнить, что скалярное произведение двух векторов подразумевает попарное умножение соответствующих элементов векторов и суммирование произведения, чтобы получить одно число. Здесь мы обобщаем это понятие для массивов двух или более измерений, поэтому умножаем соответствующие элементы массивов и затем суммируем все произведения.

$$\begin{matrix} 1,0 & 1,0 & 1,0 & 1,0 \\ 1,0 & 1,0 & 1,0 & 1,0 \\ -1,0 & -1,0 & -1,0 & -1,0 \\ -1,0 & -1,0 & -1,0 & -1,0 \end{matrix}$$

Рис. 3.1. Простой фильтр для определения горизонтальной линии

Более формально мы считаем ядро свертки функцией, *функцией ядра* (kernel function). Мы получаем значение V этой функции в позиции x, y на изображении I следующим образом:

$$V(x, y) = (I \cdot K)(x, y) = \sum_m \sum_n I(x + m, y + n)K(m, n) \quad (3.1)$$

То есть формально свертка — это операция (здесь она представлена центральной точкой), которая получает две функции, I и K , и возвращает третью функцию, которая выполняет операцию справа. Для наших обычных целей мы можем пропустить формальное определение и просто перейти к операциям с правой стороны. Кроме того, мы обычно думаем о точке x, y как о середине фрагмента (или рядом с серединой), над которым мы работаем, поэтому для ядра $4 * 4$, показанного выше, и m , и n могут варьироваться от -2 до $+1$ включительно.

0,0	0,0	0,0	0,0	0,0	0,0
0,0	2,0	2,0	2,0	0,0	0,0
0,0	2,0	2,0	2,0	0,0	0,0
0,0	2,0	2,0	2,0	0,0	0,0
0,0	0,0	0,0	0,0	0,0	0,0
0,0	0,0	0,0	0,0	0,0	0,0

Рис. 3.2. Изображение небольшого квадрата

Давайте свернем фильтр на рис. 3.1 с правой нижней частью простого изображения квадрата, показанного на рис. 3.2. Два нижних ряда фильтра перекрываются нулями на изображении. Но четыре верхних левых элемента фильтра перекрываются квадратом 2,0, поэтому значение фильтра в этом фрагменте равно 8. Естественно, если бы все значения пикселей были равны нулю, значение применения фильтра тоже было бы равно нулю. Но если бы весь фрагмент изображения состоял из всех десяток, он все равно был бы нулевым. На самом деле нетрудно увидеть, что этот фильтр имеет самые высокие значения для фрагментов с горизонтальной линией, проходящей через середину фрагмента, идущей от высоких значений сверху и значений ниже снизу. Дело, конечно, в том, что фильтры можно сделать такими, чтобы они были чувствительны к изменениям интенсивности света, а не к их абсолютным значениям, и поскольку фильтры обычно намного меньше, чем полные изображения, они концентрируются на локальных изменениях. Мы можем, конечно, спроектировать ядро фильтра, которое имеет высокие значения для фрагментов изображения с прямыми линиями, идущими от верхнего левого угла к нижнему правому или как угодно.

В приведенном выше обсуждении мы представили фильтр так, как если бы он был разработан программистом, чтобы выделить особый вид изображения, и действительно, это было сделано до появления глубокой сверточной фильтрации. Однако что делает подходы глубокого обучения особенноими, так это то, что *значения фильтра являются параметрами NN — они изучаются во время обратного прохода*. В нашей текущей дискуссии о том, как работает свертка, это лучше игнорировать, и мы продолжаем представлять наши фильтры “уже готовыми” до следующего раздела.

Помимо свертки фильтра с *фрагментом* (patch) изображения, мы говорим о свертке фильтра с *изображением* (image). Это подразумевает применение фильтра ко многим фрагментам изображения. Обычно у нас есть много разных фильтров, и цель каждого фильтра — подобрать определенный признак в изображении. Сделав это, мы можем затем передать все значения признаков в один или несколько полно связных слоев, в softmax и, следовательно, в функцию потерь. Эта архитектура показана на рис. 3.3. Здесь мы представляем слой

сверточного фильтра в виде объемного прямоугольника (поскольку банк фильтров является (как минимум) трехмерным тензором, высота и ширина, а также количество различных фильтров).

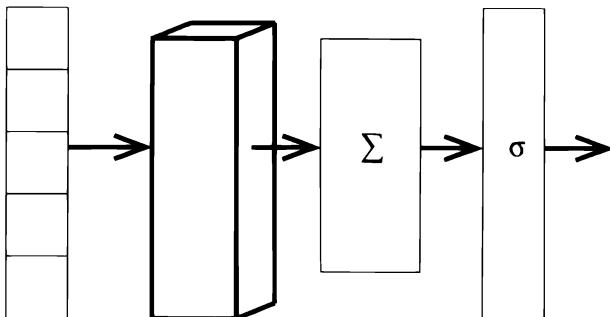


Рис. 3.3. Архитектура распознавания изображений с фильтрами свертки

Обратите внимание на преднамеренную неопределенность, описанную выше, когда мы сказали, что сворачиваем фильтр со “многими” фрагментами изображения. Точнее, мы сначала определим *шаг* (stride) — расстояние между двумя применениеями фильтра. Скажем, два шага означают, что мы применяем фильтр к каждому следующему пикселю. Чтобы стать еще более конкретными, мы говорим как о горизонтальном шаге, s_h , так и о вертикальном шаге, s_v . Идея в том, что при прохождении изображения мы применяем фильтр после каждого s_h пикселей. Достигнув конца строки, мы спускаемся вертикально на s_v строк и повторяем процесс. Применяя фильтр на шаге 2, мы по-прежнему применяем фильтр ко всем пикселям в области. Единственное, на что влияет шаг, — это на место, где именно фильтр применяется в следующий раз.

Далее мы определяем, что подразумеваем под “концом строки” при применении фильтра. Это осуществляется за счет указания *дополнения* (padding) для свертки. TF допускает два возможных дополнения: *Valid* (допустимое) и *Same* (одинаковое). После свертки фильтров с определенным участком изображения мы перемещаем s_h вправо. Существует три варианта: а) мы не находимся вблизи границы изображения, поэтому продолжаем работать над этой строкой, б) крайний слева пиксель для следующей свертки выходит за край изображения и в) крайний слева пиксель, который видит фильтр, находится в изображении, но крайний справа располагается за концом изображения. Дополнение *Same* останавливает в случае б, *Valid* прекращает в случае в. Например, на рис. 3.4 показана ситуация, когда изображение имеет ширину 28 пикселей, фильтр имеет ширину 4 пикселя и высоту 2 пикселя, а шаг равен 1. При дополнении *Valid* мы останавливаемся после пикселя 24 при отсчете с нуля. Это потому, что наш шаг приведет нас к пикселя 25, а для размещения фильтра шириной 4 потребуется 29-й пиксель, которого не существует. Дополнения

Same будут продолжать свертку до тех пор, пока не встретится пиксель 27. Естественно, мы делаем тот же выбор в вертикальном направлении, когда достигаем нижней части изображения.

0	1	23	24	25	26	27
.	.	3,2	3,1	2,5	2,0	0
.	.	3,2	3,1	2,5	2,0	0
.	.	3,2	3,1	2,5	2,0	0
.	.	3,2	3,1	2,5	2,0	0
.	.	3,2	3,1	2,5	2,0	0

Рис. 3.4. Конец строки при дополнениях *Valid* и *Same*

Решение о том, где остановиться, называется *дополнением* (padding) потому, что при горизонтальном движении с дополнением *Same* к моменту остановки мы должны использовать “мнимые” пиксели. Левая сторона фильтра находится внутри границы изображения, а правая — нет. В TF мнимые пиксели имеют значение “нуль”. Таким образом, при дополнении *Same* нужно дополнить границу изображения мнимыми пикселями. При дополнении *Valid* фактическое дополнение почти никогда не требуется, поскольку мы прекращаем свертку до того, как какая-либо часть фильтра выйдет за границы изображения. Когда требуется дополнение (при дополнении *Same*), оно применяется ко всем краям как можно более одинаково.

Поскольку это понадобится нам позже, мы дадим количество сверток фрагментов, которые мы применяем горизонтально для дополнения *Same*:

$$\lceil i_h / s_h \rceil \quad (3.2)$$

Здесь $\lceil x \rceil$ — это *функция потолка* (ceiling function). Она возвращает наименьшее целое число $\geq x$. Чтобы убедиться в необходимости функции потолка, рассмотрим случай, когда изображение имеет нечетное количество пикселей в ширину, скажем, пять, а шаг — два. Сначала фильтр применяется к фрагменту 0–2 в горизонтальном направлении. Затем он перемещается на две позиции вправо и применяется к фрагменту 2–4. Когда мы доберемся до позиции 4, фильтр следует применить к фрагменту 4–6. Поскольку ширина равна 5, позиция 6 отсутствует. Однако для дополнения *Same* мы добавляем нуль в конец строки, чтобы фильтр работал в позициях 4–6, а общее количество применений равно 4. Если дополнение *Same* не добавляет дополнительные нули, выше приведенное уравнение будет иметь *функцию пола* (floor function), а не потолка.

Естественно, те же самые рассуждения применимы и в вертикальном направлении, давая нам $\lceil i_v / s_v \rceil$.

Для дополнения *Valid* количество по горизонтали составляет

$$\lfloor (i_h - f_h + 1) / s_h \rfloor \quad (3.3)$$

Если вы не видите это последнее уравнение сразу, сначала убедитесь, что вы видите, что $i_h - f_h$ — это то, как часто вы можете сдвигаться (до того, как закончится свободное пространство), если шаг равен единице. Но количество применений составляет один плюс количество смещений.

Несмотря на использование мнимых пикселей, дополнение *Same* довольно популярно, поскольку в сочетании с одним шагом оно обладает тем свойством, что размер вывода такой же, как и у исходного изображения. Зачастую мы объединяем множество слоев свертки, каждый из которых становится входом для следующего слоя. Независимо от того, каков размер шага, дополнение *Valid* всегда имеет вывод, который меньше, чем ввод. При повторяющихся слоях свертки результат постепенно уменьшается.

Прежде чем перейти к реальному коду, обсудим, как свертка влияет на то, как мы представляем изображения. Сердцем сверточного NN в TF является двумерная функция свертки (convolution function)

```
tf.nn.conv2d(input, filters, strides, padding)
```

плюс необязательные именованные аргументы, которые мы здесь игнорируем. 2d в имени указывает, что мы сворачиваем изображение. (Существуют также версии 1d и 3d, которые сворачивают одномерные объекты, такие как аудиосигнал, или трехмерные, такие как видеоклип.) Как и следовало ожидать, первый аргумент — это размер партии отдельных изображений. До сих пор мы рассматривали отдельное изображение как двумерный массив чисел, в котором каждое число — это интенсивность света. Если вы укажете тот факт, что у нас есть размер партии, входом будет трехмерный тензор.

Но `tf.nn.conv2d` требует, чтобы отдельные изображения были трехмерными объектами и последним измерением был вектор *каналов* (*channel*). Как упоминалось ранее, нормальные цветные изображения имеют три канала — по одному для красного, синего и зеленого (RGB). Отныне, обсуждая изображения, мы все еще говорим о двумерном массиве пикселей, но каждый пикセル представляет собой список интенсивностей. Этот список содержит одно значение для черно-белых изображений и три значения для цветных.

То же самое верно для фильтров свертки. Фильтр $m * n$ совпадает с размером $m \times n$ пикселей, но теперь и пиксели, и фильтр могут иметь несколько каналов. В данном случае мы создадим фильтр, чтобы найти горизонтальные края бутылок кетчупа на рис. 3.5. Первый сверху ряд фильтра активируется наиболее

высоко, когда входной свет интенсивен только для красного и менее интенсивен для синего и зеленого. Следующие два ряда хотят меньше красного (поэтому есть некоторый контраст) и больше синего и зеленого.

$$\begin{array}{llll} (1, -1, -1) & (1, -1, -1) & (1, -1, -1) & (1, -1, -1) \\ (-1, 1, 1) & (-1, 1, 1) & (-1, 1, 1) & (-1, 1, 1) \\ (-1, 1, 1) & (-1, 1, 1) & (-1, 1, 1) & (-1, 1, 1) \end{array}$$

Рис. 3.5. Простой фильтр для определения горизонтальной линии кетчупа

На рис. 3.6 показан простой пример TF применения функции небольшой свертки к маленькому изображению. Как отмечалось выше, первый вход в conv2D — это четырехмерный тензор, где I — константа. В комментарии перед объявлением I мы демонстрируем, что она будет выглядеть как простой двумерный массив, без дополнительных измерений, добавляемых размером партии (здесь — единица) и размер каналов (снова единица). Второй аргумент — это четырехмерный тензор фильтров, здесь W также с комментарием, показывающим двухмерную версию, на этот раз без дополнительных измерений количества каналов и количества фильтров (по одному на каждый). Затем следует вызов conv2D с горизонтальными и вертикальными шагами, оба как единица и с дополнением *Valid*. Глядя на результат, мы видим, что это четырехмерность (размер партии (1), высота (2), ширина (2), каналы (1)). Высота и ширина значительно уменьшены по сравнению с размером изображения, как и следует ожидать, когда мы используем дополнение *Valid*, а также фильтр достаточно активен (со значением 6), опять же, как и следовало ожидать, поскольку он предназначен для обнаружения вертикальных² линий, которые являются именно тем, что появляется на изображении.

```
ii = [[ [0], [0], [2], [2] ],
       [ [0], [0], [2], [2] ],
       [ [0], [0], [2], [2] ],
       [ [0], [0], [2], [2] ] ]
''' ((0 0 2 2)
      (0 0 2 2)
      (0 0 2 2)
      (0 0 2 2))'''
I = tf.constant(ii, tf.float32)

ww = [ [[[ -1]], [[ -1]], [[ 1]]],
       [[[ -1]], [[ -1]], [[ 1]]],
       [[[ -1]], [[ -1]], [[ 1]]]]
```

² Вероятно, имелось в виду “горизонтальных”. — Примеч. ред.

```
''''((-1 -1 1)
 (-1 -1 1)
 (-1 -1 1))'''
W = tf.constant(ww, tf.float32)

C = tf.nn.conv2d(I, W, strides=[1, 1, 1, 1], padding='VALID')
sess = tf.Session()
print sess.run(C)
'''[[[ 6.] [ 0.]]
 [[ 6.] [ 0.]]]]'''
```

Рис. 3.6. Простое упражнение с использованием conv2D

3.2. Простой пример свертки TF

Теперь мы выполним упражнение по превращению прямой программы TF Mnist из главы 2 в сверточную модель NN. Код, который мы создаем, приведен на рис 3.7.

```
1 image = tf.reshape(img, [100, 28, 28, 1])
2     #Превращает img в четырехмерный тензор
3 flts=tf.Variable(tf.truncated_normal([4,4,1,4],stddev=0.1))
4     #Создает параметры для фильтров
5 convOut = tf.nn.conv2d(image, flts, [1, 2, 2, 1], "Same")
6     #Создает граф для свертки
7 convOut= tf.nn.relu(convOut)
8     #Не забудьте добавить нелинейность
9 convOut=tf.reshape(convOut,[100, 784])
10    #Назад к 100 одномерным векторам изображений
11 prbs = tf.nn.softmax(tf.matmul(convOut, W) + b)
```

Рис. 3.7. Основной код, необходимый для преобразования рис. 2.2 в сверточную NN

Как уже отмечалось, ключевым вызовом функции TF является `tf.nn.conv2d`. На рис. 3.7 мы видим в строке 5

```
convOut = tf.nn.conv2d(image, flts, [1, 2, 2, 1], "Same")
```

Рассмотрим каждый аргумент по очереди. Как обсуждалось только что, `image` является четырехмерным тензором — в данном случае вектором трехмерных изображений. Мы выбираем размер партии равным 100, поэтому `tf.nn.conv2d` ожидает 100 трехмерных изображений. Функции, которые читают данные в главе 2, читают векторы одномерных изображений (длиной 784), поэтому строка 1 на рис. 3.7

```
image = tf.reshape(img,[100, 28, 28, 1])
```

преобразует вход в форму $[100, 28, 28, 1]$, где последняя единица указывает, что у нас есть только один входной канал. `tf.reshape` работает почти так же, как `reshape` NumPy.

Следующий аргумент `tf.nn.conv2d` в строке 5 — это указатель на фильтры, которые будут использоваться. Это тоже четырехмерный тензор, в данном случае — формы

[высота, ширина, каналы, количество]

Параметры фильтра создаются в строке 3. Мы выбрали фильтры 4×4 ($[4, 4]$), каждый пиксель имеет один канал ($[4, 4, 1]$), и мы выбрали четыре фильтра ($[4, 4, 1, 4]$). Обратите внимание, что высота и ширина фильтра, а также количество создаваемых нами фильтров являются гиперпараметрами. Количество каналов (в данном случае — 1) определяется количеством каналов в изображении, а потому является фиксированным. Очень важно, что мы наконец-то сделали то, что обещали в начале: строка 3 создает значения фильтра в качестве параметров модели NN (с начальными значениями, равными нулю и стандартным отклонением 0,1), поэтому они запоминаются моделью.

Аргумент *шага* (stride) в `tf.nn.conv2d` представляет собой список из четырех целых чисел, указывающих размер шага в каждом из четырех измерений ввода. В строке 5 мы видим, что выбраны шаги 1, 2, 2 и 1. На практике первыми и последними почти всегда являются 1. Во всяком случае, трудно представить себе случай, когда они не будут 1. В конце концов, первое измерение — это отдельные трехмерные изображения в пакете. Если бы шаг по этому измерению равнялся двум, мы бы пропустили все остальные изображения! Столь же странно, что если бы последний шаг был больше единицы, скажем 2, и у нас было три цветовых канала, то мы смотрели бы только на красный и синий свет, пропуская зеленый. Таким образом, типичные значения для шага будут $(1, 1, 1, 1)$ или, если мы хотим свертывать только каждый второй фрагмент изображения в горизонтальном и вертикальном направлениях, $(1, 2, 2, 1)$. Вот почему вы часто видите в обсуждениях инструкций `tf.nn.conv2d`, что первый и последний шаги должны быть единицей.

Последний аргумент, *дополнение* (padding), является строкой, равной одному из типов дополнения, которые распознает TF, например *Same*.

Вывод `conv2d` очень похож на ввод. Это тоже четырехмерный тензор, и, как и ввод, первое измерение вывода — это размер партии. Или, другими словами, выходные данные представляют собой вектор выходов свертки, по одному для каждого входного изображения. Следующие два измерения — это количество применений фильтров, за которыми следуют горизонталь и вертикаль; они могут быть определены, как в уравнениях 3.2 и 3.3. Последнее измерение выходного тензора — это количество фильтров, свернутых с изображением. Выше мы указали, что будем использовать четыре, т.е. выходная форма такова:

[размер партии, размер по горизонтали,
размер по вертикали, количество фильтров]

В нашем случае это будет $(100, 14, 14, 4)$. Если мы думаем о выходе как об “изображении”, то на входе будет 28×28 с одним каналом, но на выходе будет (14×14) и 4 канала. Это означает, что в обоих случаях входное изображение представлено 784 числами. Мы выбрали это намеренно, чтобы сохранить сходство с главой 2, но мы не обязаны были делать это. Скажем, мы могли бы выбрать 16, а не четыре разных фильтра, и в этом случае у нас было бы изображение, представленное $(14 * 14 * 16 = 3136)$ числами.

В строке 11 мы подаем эти 784 значения в полносвязный слой, который создает логиты для каждого изображения, которые, в свою очередь, передаются в softmax, а затем мы вычисляем кросс-энтропийную потерю (на рис. 3.7 не показана), и у нас есть очень простая сверточная NN для Mnist. Код имеет общую форму, показанную на рис. 2.2. Кроме того, строка 7 выше помещает нелинейность между выходом свертки и входом полносвязного слоя. Это важно. Как было показано ранее, без нелинейных функций активации между линейными блоками улучшения не происходит.

Производительность этой программы значительно выше, чем в главе 2, — 96% или чуть больше, в зависимости от случайной инициализации (по сравнению с 92% для версии с прямой связью). Количество параметров модели практически одинаково для двух версий. В обоих слоях с прямой связью используется 7840 весовых коэффициентов W и 100 смещений b ($784 + 10$ весовых коэффициентов в каждом блоке в полносвязном слое, умноженное на 10 единиц). Свертка добавляет четыре сверточных фильтра, каждый с $4 * 4$ весами или еще 64 параметрами. Вот почему мы устанавливаем выходной размер свертки равным 784. В нулевой аппроксимации качество NN повышается, поскольку мы даем ей больше параметров для использования. Здесь, однако, количество параметров практически не изменилось.

3.3. Многослойная свертка

Как уже упоминалось, мы можем еще больше повысить точность, перейдя от одного уровня свертки к нескольким. В этом разделе мы строим модель с двумя слоями.

Ключевым моментом многослойной свертки является то, что мы пропустили при обсуждении вывода из `tf.conv2d`: он имеет тот же формат, что и ввод. Оба являются трехмерными векторами размера партии изображений, изображения являются двумерными плюс одно дополнительное измерение для количества каналов. Таким образом, выход из одного слоя свертки может быть входом для второго слоя, и это именно то, что он делает. Когда мы говорим о заполнителе,

полученном из данных, последнее измерение — это количество цветовых каналов. Говоря о выводе conv2d, мы имеем в виду, что последнее измерение — это количество различных фильтров в слое свертки. Слово “фильтр” здесь вполне уместно. В конце концов, чтобы пропустить через объектив только синий свет, мы буквально поместили бы перед ним цветной фильтр. Таким образом, три фильтра дают нам изображения в спектрах RGB. Теперь мы получаем “изображения” в псевдоспектрах, как “горизонтальные линейно-границевые спектры”. Это будет воображаемое изображение, создаваемое фильтром, например, на рис. 3.1. Кроме того, фильтры для изображений RGB также имеют веса, связанные со всеми тремя спектрами, второй слой свертки имеет веса для каждого канала, выводимого из первого.

Мы предоставляем код для превращения NN прямой связи Minst в двухслойную модель свертки на рис. 3.8. Строки 1–4 являются повторениями первых строк на рис. 3.7, за исключением того, что в строке 2 мы увеличиваем количество фильтров в первом слое свертки до 16 (с 4 в более ранней версии). Страна 2 отвечает за создание фильтров второго сверточного слоя flts2. Обратите внимание, что мы создали 32 из них. Это отражено в строке 5, в которой значения 32 фильтров становятся 32 значениями входного канала для второго слоя свертки.

```

1 image = tf.reshape(img, [100, 28, 28, 1])
2 flts=tf.Variable(tf.normal([4, 4, 1, 16], stddev=0.1))
3 convOut = tf.nn.conv2d(image, flts, [1, 2, 2, 1], "Same")
4 convOut= tf.nn.relu(convOut)
5 flts2=tf.Variable(tf.normal([2, 2, 16, 32], stddev=0.1))
6 convOut2 = tf.nn.conv2d(convOut, flts2, [1, 2, 2, 1], "Same")
7 convOut2 = tf.reshape(convOut2, [100, 1568])
8 W = tf.Variable(tf.normal([1568,10],stddev=0.1))
9 prbs = tf.nn.softmax(tf.matmul(convOut2, W) + b)
```

Рис. 3.8. Первичный код, необходимый для преобразования рис. 2.2 в двухслойную сверточную NN

Когда мы линеаризуем эти выходные значения в строке 7, их будет $(784 * 4)$. Помните, что мы начали с 784 пикселей и каждый слой свертки использовал шаг 2 как по горизонтали, так и по вертикали. Таким образом, полученные размеры трехмерного изображения после первой свертки были $(14, 14, 16)$. Вторая свертка также имела шаг 2 на изображении 14×14 и имела 32 канала, поэтому выходной сигнал равен $[100, 7, 7, 32]$, а линеаризованная версия одного изображения в строке 7 имеет $7 * 7 * 32 = 1568$ скалярных значений, которые также являются высотой W, которая превращает эти значения изображения в 10 логитов.

Отстраняясь от деталей, давайте рассмотрим общий поток модели. Мы начинаем с изображения $28 * 28$. В конце мы имеем “картинку” $7 * 7$, но у нас

также есть 32 различных значения фильтра в каждой точке двумерного массива. Или, другими словами, в конце мы разделили изображение на 49 фрагментов; каждый фрагмент изначально был 4×4 пикселя и теперь характеризуется 32 значениями фильтра. Поскольку это повышает производительность, мы вправе предположить, что эти значения говорят важные вещи о том, что происходит в соответствующем фрагменте 4×4 .

Действительно, похоже, что так и есть. Хотя, на первый взгляд, фактические значения в фильтрах могут сбивать с толку, по крайней мере на начальных уровнях изучения можно выявить некоторую логику в их “построении”. На рис. 3.9 показаны веса 4×4 для четырех из восьми фильтров свертки первого слоя, которые были изучены за один прогон кода на рис. 3.7. Вы можете потратить на них несколько секунд, чтобы разобрать, что они ищут. Для одних это получится. Другие не имеют особого смысла для меня. Однако взаимная корреляция с рис. 3.10 должна помочь. Рис. 3.10 был создан в результате вывода стандартного для нас теперь изображения семерки, фильтра с самым высоким значением для всех 14×14 точек в изображении после первого слоя свертки. Из тумана нулей довольно быстро проступает образ 7, поэтому фильтр 0 связан с областью всех нулей. Затем мы можем заметить, что правый край диагонали 7 в значительной степени соответствует всем семеркам, тогда как нижняя горизонтальная часть фрагментов на изображении соответствует единицам. Посмотрите еще раз на значения фильтра. Мне кажется, 1, 2 и 7 соответствуют результатам на рис. 3.9. С другой стороны, в фильтре 0 нет ничего, что могло бы предложить пустое значение. Однако это тоже имеет смысл. Мы использовали функцию `arg-max` от `Numpy`, которая возвращает позицию наибольшего числа в списке чисел. Все значения пикселей для пустых областей равны нулю, поэтому все фильтры возвращают 0. Если функция `arg-max` возвращает первое значение в случае, когда все значения равны, то это то, что мы ожидаем.

-0.152168	-0.366335	-0.464648	-0.531652
0.0182653	-0.00621072	-0.306908	-0.377731
0.482902	0.581139	0.284986	0.0330535
0.193956	0.407183	0.325831	0.284819
0.0407645	0.279199	0.515349	0.494845
0.140978	0.65135	0.877393	0.762161
0.131708	0.638992	0.413673	0.375259
0.142061	0.293672	0.166572	-0.113099
0.0243751	0.206352	0.0310258	-0.339092
0.633558	0.756878	0.681229	0.243193
0.894955	0.91901	0.745439	0.452919
0.543136	0.519047	0.203468	0.0879601

0.334673	0.252503	-0.339239	-0.646544
0.360862	0.405571	-0.117221	-0.498999
0.520955	0.532992	0.220457	0.000427301
0.464468	0.486983	0.233783	0.101901

Рис. 3.9. Фильтры 0, 1, 2 и 7 из восьми фильтров, созданных за один проход двухслойной свертки NN

```

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 5 2 2 2 2 2 2 2 2 2 2 0 0
0 0 1 1 4 4 4 4 4 2 2 2 2 0 0
0 0 1 1 1 1 1 1 1 1 2 7 0 0
0 0 0 0 0 0 0 5 1 4 2 7 0 0
0 0 0 0 0 0 0 5 1 2 7 0 0 0
0 0 0 0 0 0 5 1 4 2 7 0 0 0
0 0 0 0 0 5 2 1 2 7 0 0 0 0
0 0 0 0 0 5 1 4 2 0 0 0 0 0
0 0 0 0 5 1 4 2 7 0 0 0 0 0
0 0 0 0 2 1 2 2 0 0 0 0 0 0
0 0 0 0 1 1 1 7 0 0 0 0 0 0

```

Рис. 3.10. Наиболее активная функция для всех 14×14 точек в слое 1 после обработки рис. 1.1

Рис. 3.11 аналогичен рис. 3.10, за исключением того, что он показывает наиболее активные фильтры на уровне 2 модели. Он менее интерпретируемый, чем уровень 1. Существуют различные аргументы в пользу того, почему это может иметь место. Мы включили его в основном потому, что первый сверточный слой иллюстраций гораздо более интерпретируемый, чем большинство, но мы не должны предполагать, что то, что мы видели для слоя 1, является типичным.

0	0	0	0	0	0	0
17	11	31	17	17	16	16
6	16	12	6	6	5	5
17	17	17	5	24	5	10
0	0	11	26	3	5	0
0	17	11	24	5	10	0
0	6	24	8	5	0	0

Рис. 3.11. Наиболее активные функции для всех 7×7 точек в слое 2 после обработки рис. 1.1

3.4. Детали свертки

3.4.1. Смещения

Мы также можем иметь смещения с нашими сверточными ядрами. Мы не упоминали об этом до сих пор потому, что только в последнем примере несколько фильтров были применены к каждому фрагменту, т.е. мы указали 16 различных фильтров в строке 2 на рис. 3.8. Смещение может заставить программу придавать одному каналу фильтра больший или меньший вес, чем другому, добавляя другое значение к выходу свертки канала. Таким образом, число переменных смещения на конкретном сверточном слое равно количеству выходных каналов. Например, на рис. 3.8 мы можем добавить смещения к первому слою свертки, добавив следующее между строками 3 и 4:

```
bias = tf.Variable(tf.zeros [16])
convOut += bias
```

Широковещание неявно. Несмотря на то что convOut имеет форму [100, 14, 14, 16], bias имеет форму [16], поэтому сложение неявно создает [100, 14, 14] его копий.

3.4.2. Слои со сверткой

В разделе 2.4.4 показано, как один стандартный компонент архитектуры NN, полносвязные слои, может быть эффективно написан с использованием layers. Для сверточных слоев есть эквивалентные функции

```
tf.contrib.layers.conv2d(inpt, numFlts, fltDim, strides, pad)
```

плюс необязательные именованные аргументы. Например, строки со 2 по 4 на рис. 3.8 можно заменить строкой

```
convOut = layers.conv2d(image, 16, [4, 4], 2, "Same") ?
```

На вывод свертки указывает convOut. Как и прежде, мы создаем 16 разных ядер, каждое размером 4×4 . Шагов в обоих направлениях два, и мы используем одинаковые дополнения. Это не совсем то же самое, что и версия без слоев, так как layers.conv2d предполагает, что вы хотите смещения, если не укажете обратное. Настаивая на отсутствии смещения, мы просто присваиваем значение соответствующему именованному аргументу use_bias=False.

3.4.3. Объединение

Как можно ожидать для больших изображений (например, $1000 * 1000$ пикселей), уменьшение размера изображения от исходного до значения, подаваемого в полно связанный слой, за которым следует softmax в конце, является гораздо более значительным. Есть функции TF, которые могут помочь справиться с этим сокращением. Обратите внимание, что в нашей программе сокращение было потому, что шаги при свертке рассматривали лишь каждый следующий фрагмент. Вместо этого мы могли бы сделать следующее:

```
convOut = tf.nn.conv2d(image, flts, [1,1,1,1], "Same")
convOut = tf.nn.max_pool(convOut, [1,2,2,1], [1,2,2,1], "Same").
```

Эти две строки предназначены для замены строки 3 на рис. 3.8. Вместо свертки со вторым шагом мы сначала применили свертку с первым шагом. Таким образом, convOut имеет форму [batchSz, 28, 28, 1] — без уменьшения размера изображения. Следующая строка дает нам уменьшение размера изображения, точно равное тому, которое было получено с помощью шага два, который мы первоначально использовали.

Ключевая функция здесь, max_pool, находит максимальное значение для фильтра по области в изображении. Требуется четыре аргумента, три из четырех совпадают с conv2d. Первый — наш стандартный четырехмерный тензор изображений, третий — шаги, а последний — дополнения. В приведенном выше случае max_pool рассматривает convOut, четырехмерный вывод первой свертки. Это делается с шагами [1, 2, 2, 1]. Первый элемент списка указывает просмотреть каждое изображение размера партии, а последний — каждый канал. Две двойки указывают передвинуться на два блока перед повторением операции и делать так как по горизонтали, так и по вертикали. Второй аргумент задает размер области, в которой он должен найти максимум. Как обычно, первая и последняя единицы в значительной степени принудительны, в то время как средние две двойки указывают, что мы должны взять максимум из фрагмента $2 * 2$ convOut.

На рис. 3.12 противопоставляются два разных способа, которыми мы можем добиться снижения четырехмерности в нашей программе Mnist, хотя здесь мы делаем это на изображении $4 * 4$. (Числа придуманы.) В верхнем ряду мы применили фильтр со второго шага (дополнение *Same*) и сразу получили массив значений фильтра $2 * 2$. В строке 2 мы применили фильтр с шагом 1, что создает массив значений $4 * 4$. Затем для каждого отдельного фрагмента $2 * 2$ мы выводим наибольшее значение, которое дает нам окончательный массив в правом нижнем углу рисунка.

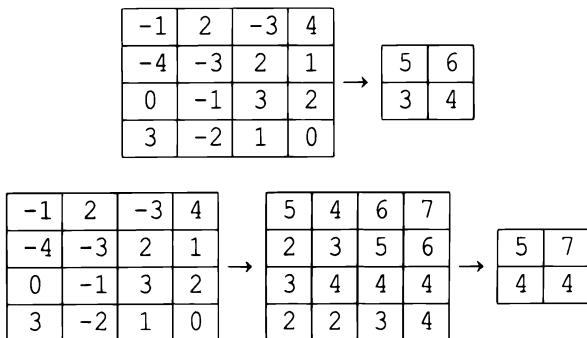


Рис. 3.12. Снижение четырехмерности с использованием max_pool и без

Прежде чем двигаться дальше, отметим, что есть также функция avg_pool, которая работает аналогично max_pool, за исключением того, что значение для пула является средним значением отдельных значений, а не максимальным.

3.5. Ссылки и что читать дальше

Вводная статья в изучение сверточных ядер с использованием NN и обратного распространения была подготовлена Яном Лекуном и другими [11], хотя гораздо более полное исследование темы принадлежит Яну Лекуну и другим [12] и является окончательной ссылкой. Часть моего образования в сверточных NN была получена из учебника Google по распознаванию цифр Mnist [13].

Если вы находитите идею получить NN, способные распознавать изображения, действительно актуальной и хотите, чтобы следующий проект работал, я бы порекомендовал набор данных CFAIR 10 (Канадский институт перспективных исследований) [14]. Это также задача классификации изображений в десяти направлениях, но объекты, которые нужно распознать, более сложные (самолет, кошка, лягушка), изображения — цветные, фоны могут быть сложными, а объект, который нужно классифицировать, не отцентрирован. Размеры изображений также больше [32, 32, 3]. Набор данных можно загрузить с [15]. Общее количество изображений примерно такое же, как у Mnist, 60 000, поэтому общий импринт данных управляем. Существует также сетевое руководство от Google по созданию NN для этой задачи [16].

Если вы действительно амбициозны, попробуйте поработать с набором данных Imagenet Large Scale Visual Recognition Challenge (ILSVRC). Это намного сложнее. Существует 1000 типов изображений, в том числе такие классические, как картофель фри или картофельное пюре. За последние шесть или семь лет этот набор данных использовался серьезными исследователями компьютерного зрения на ежегодных соревнованиях. Для NN большим был год

2012, когда в конкурсе победила программа глубокого обучения *Alexnet*, т.е. впервые победила программа NN. Программа Алекса Крижевски, Ильи Суцкевичера и Джейфри Хинтона [17] достигла *5 наивысших баллов* — в 15,5% случаев правильная метка не была одним из пяти лучших ответов, как было установлено оценкой программы вероятности того, что конкретная метка является правильной. Участник, занявший второе место, набрал 26,2%. С 2012 года все первые места имеют программы на базе NN.

2012	15,5
2013	11,2
2014	6,7
2015	3,6
Человек	5–10

Здесь запись “Человек” означает, что люди справляются с этой задачей в диапазоне от 5 до 10% в зависимости от обучения.

Таблицы и диаграммы с вышеприведенной информацией являются общими точками представления при объяснении влияния глубокого обучения на искусственный интеллект за последние 10 лет или около того.

3.6. Упражнения

Упражнение 3.1. а) Создайте ядро $3 * 3$, которое обнаруживает вертикальные линии на черно-белом изображении и возвращает значение 8, когда применяется к верхней левой части изображения на рис. 3.2. Оно должно возвращать нуль, если все пиксели в фрагменте имеют одинаковую интенсивность. б) Разработайте другое такое ядро.

Упражнение 3.2. В обсуждении уравнения 3.2 мы в комментарии сказали, что размер фильтра свертки не влияет на количество приложений при использовании дополнения *Same*. Объясните, как это может быть.

Упражнение 3.3. В нашем обсуждении дополнения мы говорили, что дополнение *Valid* всегда дает выходное изображение, имеющее меньшие двумерные размеры, чем входное. Строго говоря, это не так. Объясните (относительно неинтересный) случай, когда это утверждение неверно.

Упражнение 3.4. Предположим, что вход для сверточной NN является цветным изображением $32 * 32$. Мы хотим применить к нему восемь сверточных фильтров, все с формой $5 * 5$. Мы используем дополнение *Valid* и шаг 2 по вертикали и по горизонтали. а) Какова форма переменной, в которой мы храним значения фильтров? б) Какова форма вывода `tf.nn.conv2d`?

Упражнение 3.5. Объясните, что делает следующий код иначе, чем почти идентичный код в начале раздела 3.4.3?

```
convOut = tf.nn.conv2d(image, flts, [1,1,1,1], "Same")
convOut = tf.nn.maxpool(convOut, [1,2,2,1], [1,1,1,1] "Same").
```

В частности, необязательные значения `image` и `flts`, `convOut` имеют одинаковую форму, в обоих случаях? Обязательно ли им иметь одинаковые значения? Является ли один набор значений правильным подмножеством другого? Почему в каждом случае “да” или почему “нет”?

Упражнение 3.6. а) Сколько переменных создается, когда мы выполняем следующую команду `layers`?

```
layers.conv2d(image,10, [2,4], 2, "Same", use_bias=False)
```

Предположим, что `image` имеет форму `[100, 8, 8, 3]`. Какие из этих значений формы не имеют отношения к ответу? б) Не имеет значения, сколько, если для `use_bias` установлено значение `True` (стандартно)?

Глава 4

Векторное представление слов и рекуррентные NN

4.1. Векторное представление слова для языковых моделей

Языковая модель (language model) — это распределение вероятностей по всем строкам (strings) языка. На первый взгляд, это довольно сложная для понимания идея. Рассмотрим, например, предыдущее предложение “At first blush...” (На первый взгляд...). Существует большая вероятность, что вы никогда не видели это конкретное предложение, и если вы не перечитаете эту книгу, вы никогда не увидите его во второй раз. Какой бы ни была вероятность этого, она должна быть очень маленькой. И все же противопоставьте это предложение тому, в котором есть те же слова, но в обратном порядке. Это еще менее вероятно. Таким образом, цепочки слов могут быть более или менее разумными. Кроме того, программы, которые хотят, скажем, перевести польский на английский, должны иметь возможность различать предложения, которые звучат, как английский язык, и предложения, которые так не делают. Языковая модель является формализацией этой идеи.

Мы можем развить эту концепцию далее, разбив строки на отдельные слова, а затем выяснив вероятность следующего слова с учетом предыдущих? Итак, пусть $E_{1,n} = (E_1 \dots E_n)$ — это последовательность из n случайных величин, обозначающая строку из n слов, а $e_{1,n}$ — это значение одного кандидата. Например, если n равно 6, а $e_{1,6}$ составляет “We live in a small world” (Мы живем в очень маленьком мире), мы могли бы использовать цепное правило, чтобы получить следующее:

$$P(\text{We live in a small world}) = P(\text{We})P(\text{live}|\text{We})P(\text{in}|\text{We live})\dots \quad (4.1)$$

В более общем смысле

$$P(E_{1,n} = e_{1,n}) = \prod_{j=1}^{j=n} P(E_j = e_j | E_{1,j-1} = e_{1,j-1}). \quad (4.2)$$

Прежде чем продолжить, имеет смысл вернуться к тому, что мы упомянули как “разбиение строк на последовательность слов”. Это *лексический анализ* (*tokenization*), и если бы это была книга о понимании текста, мы могли бы посвятить этому отдельную главу. Тем не менее у нас есть и другие заботы, поэтому мы просто говорим, что “слово” для наших целей — это любая последовательность символов между двумя пробелами (символ новой строки мы также считаем пробелом). Обратите внимание, что это означает, что, например, “1066” — это слово в предложении “The Norman invasion happened in 1066.” (Нормандское вторжение произошло в 1066 году). На самом деле это неверно: согласно нашему определению “слова” это слово, которое встретилось в приведенном выше предложении, “1066.”, т.е. “1066” с точкой после него. Таким образом, мы также предполагаем, что знаки препинания (например, точки, запятые, двоеточия) отделены от слов, так что последняя точка становится самостоятельным словом, отдельным от слова “1066”, которое предшествовало ему. (Возможно, вы начинаете понимать, почему мы могли бы потратить на это целую главу.)

Кроме того, мы собираемся ограничить наш английский словарь неким фиксированным размером, скажем 10 000 разных слов. Для обозначения нашего словаря мы используем символы V (*vocabulary*) и $|V|$ как его размер. Это необходимо, потому что согласно приведенному выше определению “слова” мы должны ожидать увидеть слова в наших наборах для разработки и тестов, которые не встречаются в обучающем наборе, например “132 423” в предложении “Население Провиденса составляет 132 423 человек”. Мы делаем это, заменяя все слова, не входящие в V (так называемые *неизвестные слова* (*unknown word*)), специальным словом “*UNK*”. Таким образом, это предложение теперь встретилось бы в нашем корпусе как “Население Провиденса составляет *UNK* человек”.

Данные, которые мы используем в этой главе, известны как *Penn Treebank Corpus*, или сокращенно *PTB*. *PTB* состоит из порядка 1 000 000 слов новостных статей из *Wall Street Journal*. Он был подвергнут лексическому анализу, но не “помечен” (*unked*), поэтому объем словаря близок к 50 000 слов. Это “банк синтаксических деревьев” (“treebank”), поскольку все предложения были превращены в деревья, демонстрирующие их грамматическую структуру. Здесь мы игнорируем деревья, так как нас интересуют только слова. Мы также заменяем все слова, встречающиеся 10 или менее раз, словом **UNK**.

После этого вернемся к уравнению 4.2. Если бы у нас было очень большое количество текста на английском языке, мы могли бы оценить первые две или три вероятности с правой стороны, просто посчитав, как часто мы видим, например, “We live” (Мы живем) и как часто “in” (в) встречается следующим, а затем делим второе на первое (т.е. используем оценку максимального

правдоподобия), чтобы получить оценку, например, $P(\text{in}|\text{We live})$. Но по мере увеличения n это невозможно из-за отсутствия каких-либо конкретных примеров в обучающем корпусе, скажем, последовательности из 50 слов.

Одно из стандартных решений этой задачи заключается в предположении, что вероятность следующего слова зависит только от предыдущего одного или двух слов, так что при оценке вероятности следующего мы можем игнорировать все слова перед этим. Версия, в которой мы предполагаем, что слова зависят только от предыдущего слова, выглядит следующим образом:

$$P(E_{1,n} = e_{1,n}) = P(E_1 = e_1) \prod_{j=2}^{j=n} P(E_j = e_j | E_{j-1} = e_{j-1}) \quad (4.3)$$

Это модель *биграмм* (bigram model), где “биграмма” означает “два слова”. Это называется так потому, что каждая вероятность зависит только от последовательности из двух слов. Мы можем упростить это уравнение, если поместим воображаемое слово “STOP” (СТОП) в начале корпуса, а затем после каждого предложения. Это *дополнение предложений* (sentence padding). Поэтому, если первый “STOP” равен e_0 , уравнение 4.3 становится таким:

$$P(E_{1,n} = e_{1,n}) = \prod_{j=1}^{j=n} P(E_j = e_j | E_{j-1} = e_{j-1}) \quad (4.4)$$

Отныне мы предполагаем, что предложения всех наших языковых корпусов дополнены. Таким образом, за исключением первого слова “STOP”, наша языковая модель прогнозирует все слова “STOP”, а также все реальные слова.

С помещенными нами упрощениями должно быть ясно, что создание плохой модели языка вполне тривиально. Если, скажем, $|V| = 10\,000$, мы можем принять вероятность того, что любое слово будет следовать за любым другим как $1/10000$. Разумеется, мы хотим получить хороший вариант, в котором, если последним словом является “the”, распределение присваивает очень низкую вероятность “a”, а гораздо более высокую, — скажем, “cat”. Мы делаем это, используя глубокое обучение, т.е. мы даем глубокой сети слово w_i и ожидаем в качестве вывода разумного распределения вероятности по возможным следующим словам.

Для начала нам нужно как-то превратить слова в те вещи, которыми могут манипулировать глубокие сети, т.е. в числа с плавающей запятой. Теперь стандартное решение заключается в том, чтобы связать каждое слово с вектором чисел с плавающей запятой. Эти векторы являются *векторным представлением слова* (word embedding). Для каждого слова мы инициализируем его векторное представление как вектор из e чисел с плавающей запятой, где e — это системный гиперпараметр. Векторы, где e равно 20 редки, 100 встречаются часто,

а 1000 неизвестны. На самом деле мы делаем это в два этапа. Вначале каждое слово в словаре V имеет уникальный индекс (целое число) от 0 до $|V|-1$. Тогда у нас будет массив \mathbf{E} размерностью $|V| \times e$. \mathbf{E} содержит все векторные представления слов, так что, если, скажем, “the” имеет индекс 5, 5-я строка \mathbf{E} является векторным представлением слова “the”.

С учетом этого очень простая сеть прямой связи для оценки вероятности следующего слова показана на рис. 4.1. Небольшой квадрат слева является входом в сеть — это целочисленный индекс текущего слова, e_i . Справа — вероятности, присвоенные возможным следующим словам e_{i+1} , а функцией кросс-энтропийных потерь является $-\ln P(e_o)$, отрицательный натуральный логарифм вероятности, назначенный правильному следующему слову. Снова возвращаясь влево, текущее слово немедленно преобразуется в его векторное представление *слоем векторного представления* (embedding layer), который просматривает e_i -ю строку в \mathbf{E} . С этого момента все операции NN осуществляются над векторным представлением слова.

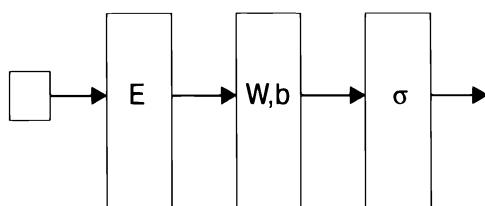


Рис. 4.1. Сеть прямой связи для языковой модели

Критически важным моментом является то, что \mathbf{E} — это параметр модели. Таким образом, первоначально числа в \mathbf{E} будут случайными со средним нулевым значением и небольшим стандартным отклонением, их значения модифицируются в соответствии со случайным градиентным спуском. В более общем смысле при обратном проходе Tensorflow начинает с функции потерь и работает в обратном направлении, ища все параметры, которые влияют на потери. \mathbf{E} является одним из таких параметров, поэтому TF его изменяет. Что удивительно в этом, кроме того факта, что процесс сходится к устойчивому решению, так это то, что решение обладает тем свойством, что слова, которые ведут себя аналогичным образом, заканчиваются векторными представлениями слов, которые “близки друг к другу”. Так что если e (размер вектора представления слова) равно, скажем, 30, то предлоги “near” (близко) и “about” (около) указывают примерно в одном и том же направлении в 30-мерном пространстве и ни один из них не очень близок, скажем, к слову “computer” (компьютер) (которое ближе к слову “machine” (машина)).

Но если немного подумать, это, возможно, не так удивительно. Давайте более внимательно рассмотрим, что происходит с векторными представлениями

слова, когда мы пытаемся минимизировать потери. Как уже говорилось, функция потерь — это кросс-энтропийные потери. Первоначально все значения логита примерно одинаковы, поскольку все параметры модели примерно равны (и близки к нулю).

Теперь предположим, что мы уже обучались на паре слов “says that” (говорит, что). Это может привести к тому, что параметры модели будут перемещаться таким образом, что векторное представление слова “says” приведет к более высокой вероятности того, что слово “that” будет следующим. Теперь рассмотрим случай, когда модель впервые встречает слово “recalls” (помните), и скажем, что, кроме того, за ним также следует слово “that”. Один из способов изменить параметры так, чтобы слово “recalls” прогнозировало слово “that” с более высокой вероятностью, — сделать его векторное представление более похожим на векторное представление, что и у слова “says”, поскольку оно тоже желает прогнозировать слово “that” в качестве следующего слова. Это действительно то, что происходит. В более общем смысле два слова, за которыми следуют одинаковые слова, получают одинаковые векторные представления слова.

Номера слов	Слово	Наибольшее косинусное сходство	Самые похожие
0	under		
1	above	0,362	0
2	the	-0,160	0
3	a	0,127	2
4	recalls	0,479	1
5	says	0,553	4
6	rules	-0,066	4
7	laws	0,523	6
8	computer	0,249	2
9	machine	0,333	8

Рис. 4.2. Десять слов, наибольшие косинусные сходства с предыдущими словами и индексы слов с наибольшим сходством

На рис. 4.2 показано, что происходит, когда мы запускаем нашу модель для миллиона слов текста, когда размер словаря составляет порядка 7 500 слов и размер векторного представления слова 30. *Косинусное сходство* (cosine similarity) двух векторов является стандартной мерой того, насколько два вектора близки один к другому. В случае двумерных векторов это стандартная функция косинуса, равная 1,0, если векторы указывают в одном и том же направлении, 0 — если они ортогональны, и -1,0 — если направлены противоположно. Вычисление косинусного сходства произвольной размерности осуществляется так:

$$\cos(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\left(\sqrt{\sum_{i=1}^{i=n} x_i^2} \right) \left(\sqrt{\sum_{i=1}^{i=n} y_i^2} \right)} \quad (4.5)$$

На рис. 4.2 показаны пять пар схожих слов, пронумерованных от нуля до девяти.

Для каждого слова мы вычисляем его косинусное сходство со всеми предшествующими ему словами. Таким образом, мы ожидаем, что все нечетные слова будут наиболее похожи на слова, которые расположены непосредственно перед ними, и это действительно так. Мы также ожидаем, что четные слова (первое из каждой пары схожих слов) не будут очень похожи ни на одно из предыдущих слов. По большей части это также верно.

Поскольку подобие векторных представлений в значительной степени отражает сходство, векторные представления слов были приняты как средство количественного определения “смысла”, и теперь мы знаем, как можно несколько улучшить этот результат. Основным фактором является то, сколько слов мы используем для обучения, хотя есть и другие архитектуры, которые также помогают. Однако большинство методов страдает от подобных ограничений. Например, они часто “слепнут”, когда пытаются различать синонимы и антонимы. (Возможно, слова “under” (под) и “above” (над) являются антонимами.) Помните, что языковая модель пытается угадать следующее слово, поэтому слова с похожими следующими словами имеют похожие векторные представления, и очень часто антонимы делают именно это. Кроме того, получить хорошие модели для векторных представлений фраз гораздо сложнее, чем для отдельных слов.

4.2. Создание языковых моделей с прямой связью

Теперь давайте создадим программу TF для вычисления вероятностей bigramm. Она очень похожа на модель распознавания цифр на рис. 2.2, так как в обоих случаях мы имеем единую полносвязную NN с прямой связью, заканчивающуюся слоем softmax, чтобы получить вероятности, необходимые для кросс-энтропийных потерь. Есть только несколько различий.

В первую очередь, вместо изображения NN получает в качестве ввода индекс слова i , где $0 \leq i < |V|$, и первым делом нужно заменить его на $\mathbf{E}[i]$, векторное представление слова:

```
inpt=tf.placeholder(tf.int32, shape=[batchSz])
answr=tf.placeholder(tf.int32, shape=[batchSz])
E = tf.Variable(tf.random_normal([vocabSz, embedSz],
                                 stddev = 0.1))
embed = tf.nn.embedding_lookup(E, inpt)
```

Мы предполагаем, что здесь не показан код, который читает слова и заменяет символы их уникальными индексами. Кроме того, этот код упаковывает batchSz из них в вектор. На этот вектор указывает inpt. Правильный ответ для каждого слова (следующего слова в тексте) — это похожий вектор answr. Далее мы создали встроенный поисковый массив E. Функция tf.nn.embedding_lookup создает необходимый код TF и помещает его в граф вычислений. Будущие манипуляции (например, tf.matmul) далее работают с embed. Естественно, TF может определить, как обновить E, чтобы снизить потери, равно как и другие параметры модели.

Переходя к другому концу сети прямой связи, мы используем встроенную функцию TF для вычисления кросс-энтропийных потерь:

```
xEnt = tf.nn.sparse_softmax_cross_entropy_with_logits
       (logits=logits, labels=answr)
loss = tf.reduce_sum(xEnt)
```

Функция TF tf.nn.sparse_softmax_cross_entropy_with_logits получает два именованных аргумента. Здесь аргумент logits представляет собой batchSz значений logit (т.е. массив batchSz на vocabSz логитов). Аргумент labels — это вектор правильных ответов. Функция передает логиты в softmax, чтобы получить вектор столбцов вероятностей batchSz на vocabSz, т.е. $s_{i,j}$ — это i,j -й элемент softmax, являющийся вероятностью слова j в i -ом примере данной партии. Затем функция определяет вероятность правильного ответа (из answr) для каждой строки, вычисляет ее натуральный логарифм и выводит массив batchSz на 1 (фактически вектор-столбец) этих логарифмических вероятностей. Стока 2 выше получает этот вектор-столбец и суммирует его, чтобы получить общую потерю для этой серии примеров.

Если вам любопытно, использование слова “sparse” (разреженный) здесь такое же, как, например, в словах “sparse matrix” (разреженная матрица), и предположительно взято из них. *Разреженная матрица* — это матрица с очень малыми ненулевыми значениями, поэтому эффективно хранить только место и значения ненулевых значений. Возвращаясь к нашему вычислению потерь в первой программе TF Mnist (с. 43), мы предположили, что правильные метки для цифровых изображений были предоставлены в виде унитарных векторов с положением правильного ответа, отличным от нуля. В tf.nn.sparse_softmax мы просто даем правильный ответ. Правильный ответ можно рассматривать как разреженную версию одного унитарного представления.

Возвращаясь к языковой модели с этим кодом, мы пройдем несколько эпох с нашими обучающими примерами и получим векторные представления слов, которые демонстрируют сходство слов, как на рис. 4.2. Кроме того, чтобы оценить языковую модель, мы можем выводить общие потери в обучающем наборе после каждой эпохи. Они должны уменьшаться с увеличением количества эпох.

В главе 1 (с. 30) мы предложили после эпохи обучения выводить среднюю потерю за пример, поскольку, если наши параметры улучшают модель, потери должны уменьшаться (числа, которые мы видим, должны уменьшаться). Здесь мы предлагаем небольшую вариацию этой идеи. В первую очередь, обратите внимание, что в языковом моделировании “пример” присваивает вероятности для возможных следующих слов, поэтому количество обучающих примеров — это количество слов в нашем учебном корпусе. Таким образом, вместо того чтобы говорить о средней потере за пример, мы говорим о *средней потере за слово* (average per-word loss). Далее, вместо того чтобы выводить среднюю потерю за слово, лучше вывести e , возведенное в эту степень, т.е. для корпуса d с $|d|$ словами, если общая потеря равна x_d , выводим

$$f(d) = e^{\frac{x_d}{|d|}}. \quad (4.6)$$

Это *перплексивность* (perplexity) корпуса d . Это хорошее число для размышлений, поскольку оно на самом деле имеет несколько интуитивное значение: в среднем угадывание следующего слова эквивалентно угадыванию результата броска честного кубика с таким количеством граней. Обратите внимание, что это значит угадывание второго слова нашего учебного корпуса по первому слову. Если у нашего корпуса есть словарь размером 10 000 слов и мы начинаем со всеми нашими параметрами близкими к нулю, то 10 000 логитов в первом примере равны нулю и все вероятности равны 10^{-4} . Читатели должны подтвердить, что это приводит к перплексивности, которая в точности соответствует размеру словаря. По мере того как мы обучаемся, перплексивность уменьшается, и для конкретного корпуса, который ваш автор использовал с размером словаря около 7 800 слов, после двух учебных эпох с обучающим набором из примерно 10^6 слов набор разработки имел перплексивность 180 или около того. На ноутбуке с четырьмя процессорами эпоха модели занимала около трех минут.

4.3. Улучшение языковых моделей с прямой связью

Существует множество способов улучшения языковой модели, которую мы только что разработали. Например, в главе 2 мы увидели, что добавление скрытого слоя (с функцией активации между двумя слоями) улучшило производительность Mnist с 92% до 98%. Добавление скрытого слоя улучшает перплексивность набора разработки с 180 до 177.

Но самый простой способ получить большую перплексивность — это перейти от модели биграмм языка к *модели триграмм* (trigram model). Помните, что при переходе от уравнения 4.2 к уравнению 4.4 мы предполагали, что

вероятность слова зависит только от предыдущего слова. Очевидно, это неверно. В общем, на выбор следующего слова могут влиять слова, находящиеся произвольно далеко сзади, и влияние слова “сзади” очень велико. Таким образом, правильно обученная модель, основывающая свои предположения на двух предыдущих словах (называемая моделью *триграмм* (trigram) потому, что вероятности основаны на последовательностях из трех слов), получает гораздо большую перплексивность, чем модели биграмм.

В модели биграмм у нас был один заполнитель для предыдущего индекса слова, `inpt`, и один — для прогнозируемого слова (при условии, что размер партии равен единице) `answr`. Теперь мы введем третий заполнитель, который имеет индекс на два слова назад, `inpt2`. В графе вычислений TF мы добавляем узел, который находит векторное представление `inpt2`,

```
embed2 = tf.nn.embedding_lookup(E, inpt2),
```

а затем — один для объединения двух:

```
both= tf.concat([embed, embed2], 1)
```

Здесь второй аргумент указывает, для какой оси тензора была выполнена конкатенация. (Помните, что на самом деле мы осуществляем векторное представление размера партии одновременно, поэтому каждый из результатов поиска представляет собой матрицу размера партии, умноженного на размер векторного представления. Мы хотим получить матрицу размера партии на (размер-векторного-представления * 2), поэтому конкатенация происходит вдоль оси 1, строки (помните, что столбцы — это ось 0). Наконец, нам нужно изменить размеры `W` с размер-векторного-представления * размер-словаря на (размер-векторного-представления * 2) * размер-словаря.

Другими словами, мы вводим векторные представления для двух предыдущих слов, а NN использует оба при оценке вероятности следующего слова. Кроме того, обратный проход обновляет векторные представления обоих слов. Это снижает перплексивность со 180 до 140. Добавление еще одного слова к входному слою еще больше снижает перплексивность, примерно до 120.

4.4. Переобучение

В разделе 1.6 мы обсудили предположение iid, скрывающееся за всеми гарантиями того, что методы обучения наших NN действительно приводят к хорошим весам. В частности, мы отметили, что как только мы используем учебные данные для более чем одной эпохи, все ставки отменяются.

Но кроме довольно надуманного примера, мы не предоставили никаких эмпирических доказательств по этому вопросу. Причина в том, что данные,

которые мы использовали для наших примеров в главе 1, Mnist, с точки зрения набора данных ведут себя очень хорошо. В конце концов, мы хотим получить такие учебные данные, которые охватывали бы все возможные варианты (и в правильных пропорциях), поэтому, когда мы смотрим на учебные данные, никаких сюрпризов не возникает. Имея всего десять цифр и 60 000 учебных примеров, набор Mnist вполне соответствует этому критерию.

К сожалению, большинство наборов данных не настолько полны, и письменные наборы данных на неком языке в целом (и национальный корпус американского английского Penn Treebank, в частности) далеки от идеала.

Даже если ограничить размер словаря 8 000 слов и рассмотреть только модели триграмм, в тестовом наборе есть большое количество триграмм, которые не отображаются в учебных данных. В то же время многократное рассмотрение одного и того же (относительно небольшого) набора примеров приводит к тому, что модель переоценивает их вероятность. На рис. 4.3 показаны результаты по перплексивности для двухслойной модели языка триграмм, обученной на корпусе Penn Treebank. В строках указаны количество эпох обучения и средняя перплексивность для каждого учебного примера в каждую эпоху, за которой следует среднее значение по всему корпусу разработки.

Эпоха	1	2	3	4	5	6	7	10	15	20	30
Учебный	197	122	100	87	78	72	67	56	45	41	35
Разработки	172	152	145	143	143	143	145	149	159	169	182

Рис. 4.3. Переобучение в языковой модели

В первую очередь, глядя на строку перплексивностей при обучении, мы видим, что оно монотонно уменьшается с увеличением количества эпох. Так и должно быть. Стока разработки рассказывает более сложную историю. Здесь перплексивность также начинает уменьшаться, со 172 в первую эпоху до 143 в четвертую, но затем остается устойчивой на протяжении двух эпох, а начиная с седьмой эпохи — увеличивается. К 20-й итерации значение достигает 169, а к 30-й — 182. Разница между результатами обучения и разработки на 30-й эпохе составляет 35 против 182, что является классическим примером переобучения на учебных данных.

Регуляризация (regularization) является общим термином для модификаций, позволяющих избежать переобучения. Самая простая методика регуляризации — это *ранняя остановка* (early stopping): мы просто прекращаем обучение в точке, где перплексивность разработки самая низкая. Но хотя ранняя остановка проста, это не лучший метод для решения проблемы переобучения. На рис. 4.4 показаны два гораздо лучших решения: *исключение* (dropout) и *регуляризация L2* (L2 regularization).

Эпоха	1	2	3	4	5	6	7	10	15	20	30
Исключение	213	182	166	155	150	144	139	131	122	118	114
Регуляризация L2	180	163	155	148	144	140	137	130	123	118	112

Рис. 4.4. Перплексивность языковой модели при использовании регуляризации

При исключении мы изменяем сеть так, чтобы случайно отбрасывать части нашего вычисления. Например, данные для исключения на рис. 4.4 получены при случайном исключении 50% вывода линейных блоков первого слоя. Таким образом, следующий слой видит нули в случайных местах своего входного вектора. Дело в том, что обучающие данные больше не идентичны на каждой эпохе, поскольку каждый раз отбрасываются разные единицы. Кроме того, обратите внимание, что классификатор не может зависеть от совпадения множества особенностей данных, выстраивающихся в линию определенным образом, а следовательно, он должен обобщать лучшее. Как видно из рис. 4.4, это действительно помогает предотвратить переобучение. Первая строка рис. 4.4 больше не демонстрирует изменений в перплексивности корпуса разработки. Даже на 30-й эпохе перплексивность уменьшается, хотя и со скоростью, сходной с ледниковой (около 0,1 единицы за эпоху). Кроме того, абсолютное значение при использовании исключения намного лучше, чем можно достичь при раннем останове — перплексивность 114 против 145.

Второй метод, который мы демонстрируем на рис. 4.4, — это регуляризация L2. Она начинается с наблюдения, что переобучение во многих видах машинного обучения сопровождается тем, что параметры обучения становятся достаточно большими (или довольно маленькими для весов ниже нуля). Ранее мы отмечали, что повторение одних и тех же данных приводит к тому, что NN переоценивает вероятности того, что он видел, за счет всех примеров, которые могут встречаться, но отсутствуют в обучающих данных. Эта переоценка достигается за счет весов с большими абсолютными значениями или, что почти эквивалентно, большими квадратами. В регуляризации L2 мы добавляем к функции потерь величину, пропорциональную сумме квадратов весов, т.е., если бы мы не использовали кросс-энтропийную потерю, наша новая функция потерь была бы такой:

$$\mathcal{L}(\Phi) = -\log(\Pr(c)) + \alpha \frac{1}{2} \sum_{\phi \in \Phi} \phi^2 \quad (4.7)$$

Здесь α — это действительное число, контролирующее то, как мы взвешиваем два члена. Оно обычно невелико; в приведенных выше экспериментах мы установили для него значение 0,01, вполне типичное значение. Когда мы дифференцируем функцию потерь по ϕ , второй член добавляет $\alpha\phi$ к сумме $\frac{\partial \mathcal{L}}{\partial \phi}$.

Это побуждает как положительные, так и отрицательные ϕ приближаться к нулю.

В этом примере обе формы регуляризации работают примерно одинаково, хотя в целом предпочтительным методом является исключение. Его легко добавить в сеть TF. Чтобы исключить, скажем, 50% значений, выходящих из первого слоя линейных блоков (например, $w1Out$), добавляем в программу код

```
keepP= tf.placeholder(tf.float32)
w1Out=tf.nn.dropout(w1Out, keepP)
```

Обратите внимание, что мы сделали вероятность заполнителем. Обычно мы хотим сделать это потому, что применяем исключение только во время обучения. При тестировании это не нужно, да и вообще вредно. Делая значение заполнителем, мы можем вводить значения 0,5 при обучении и 1,0 — при тестировании.

Использование регуляризации L2 так же просто. Если мы хотим, чтобы значения, например, W_1 , веса некоторых линейных блоков, не становились слишком большими, мы просто добавляем

```
.01*tf.nn.l2_loss(W1)
```

к функции потерь, используемой при обучении. Здесь 0,01 — это гиперпараметр для взвешивания того, насколько мы считаемся с регуляризацией по сравнению с исходной функцией кросс-энтропийных потерь. Если ваш код вычисляет перплексивность, возведя e до потери за слово, обязательно отделите объединенную потерю, используемую при обучении, от потери, используемой при вычислении перплексивности. Для последнего случая нам нужны только кросс-энтропийные потери.

4.5. Рекуррентные сети

Рекуррентная нейронная сеть (Recurrent Neural Network — RNN) в некотором смысле противоположна NN с прямой связью. Это сеть, в которой вывод также вносит свой вклад. В терминологии графов это направленный циклический граф, в отличие от направленного ациклического графа с прямой связью. Простейшая версия RNN показана на рис. 4.5. Прямоугольник с надписью $W_r b_r$ состоит из слоя линейных блоков с весами W_r и смещениями b_r , плюс функция активации. Ввод поступает в него снизу слева, а вывод o выходит справа и разделяется. Один экземпляр возвращается назад; именно этот круг делает такую сеть рекуррентной, а не прямой. Другая копия отправляется на второй слой линейных блоков с параметрами W_o , b_o , который отвечает за вычисление выходных данных RNN и потерь. Алгебраически мы можем выразить это следующим образом:

$$s_0 = 0 \quad (4.8)$$

$$s_{t+1} = \rho((e_{t+1} \cdot s_t)W_r + b_r) \quad (4.9)$$

$$o = s_{t+1}W_o + b_o \quad (4.10)$$

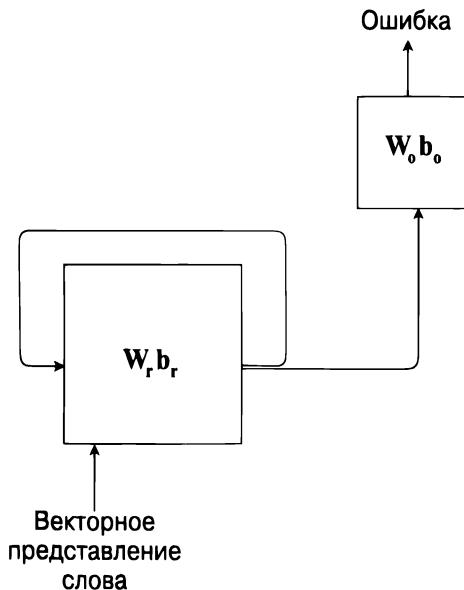


Рис. 4.5. Графическая иллюстрация рекуррентной NN

Мы начинаем рекуррентное отношение с состояния s_0 , инициализированного некоторым произвольным значением, обычно — вектором нулей. Размерность вектора состояний является гиперпараметром. Мы получаем следующее состояние, объединяя следующий вход (e_{t+1}) с предыдущим состоянием s_t и передавая результат через линейный блок W_r , b_r . Затем мы передаем вывод через функцию $\text{relu } \rho$. (Функция активации выбирается произвольно.) Наконец, выходной сигнал блока RNN o получается в результате подачи текущего состояния через второй линейный блок W_o , b_o . Функция потерь при обучении также выбирается произвольно, зачастую это кросс-энтропия, вычисляемая по o .

Рекуррентные сети целесообразны, когда мы хотим, чтобы предыдущие входы в сеть оказывали в будущем влияние произвольно. Поскольку язык работает таким же образом, сети RNN зачастую используются в задачах, связанных с языком в целом и в моделировании языка в частности. Таким образом, здесь мы предполагаем, что ввод — это векторное представление текущего слова w_i , прогноз — w_{i+1} , а потеря — наши стандартные кросс-энтропийные потери.

Вычисление прямого прохода NN работает почти так же, как и в NN с прямой связью, за исключением того, что мы помним об o из предыдущей итерации и объединяем его с текущим векторным представлением слова в начале прямого прохода. Обратный проход, однако, не так очевиден. Ранее, объясняя, почему параметры векторного представления слов также обновляются TF, мы говорили, что TF работает в обратном направлении от функции потерь, отслеживая ее и продолжая искать параметры, которые влияют на ошибку, а затем дифференцирует ошибку по отношению к этим параметрам. В главе 1 для NN Mnist это возвращало нас через слой с W и b , но затем прекратило, когда мы столкнулись только с пикселями изображения. То же самое верно для сверточных NN, хотя способы, которыми параметры вводятся в вычисление функции ошибок, несколько сложнее. Но теперь потенциально нет предела тому, как далеко нам нужно идти назад.

Предположим, мы читаем 500-е слово и хотим изменить параметры модели, потому что не спрогнозировали w_{501} с вероятностью 1. Возвращаясь назад, мы обнаруживаем, что часть ошибки связана с весами W_o сети в правом верхнем углу рис. 4.5. Но, конечно, одним из входов для этого слоя является выход рекуррентного блока o_{500} , когда он только что обработал слово w_{500} . И откуда взялось это значение? Ну, W_r , b_r , но также частично это связано с o_{499} . Короче говоря, чтобы сделать это “правильно”, нужно обратно проследить ошибку через 500 циклов рекуррентного слоя, корректируя веса W_r , b_r , снова и снова, благодаря вкладам всех ошибок, начиная со слова 1. Это не практично.

Мы решаем эту проблему грубой силой: просто отключаем вычисления после некоторого произвольного количества итераций в обратном направлении. Количество итераций называется *размером окна* (window size) и является системным гиперпараметром. Общая методика называется *обратным распространением во времени* (back propagation through time) и показана на рис. 4.6, где мы принимаем размер окна 3. (Более реалистичное значение для размера окна было бы, скажем, 20.) Более подробно на рис. 4.6 показано, что мы обрабатываем корпус, который начинается с фразы “It is a small world but I like it that way” (Это маленький мир, но мне он нравится) вместе с дополнением предложений. Обратное распространение во времени рассматривает рис. 4.6 так, как если бы это была сеть с прямой связью, получающая не одно слово, а размер окна (т.е. 3), а затем вычисляющая ошибку для трех. Для нашего короткого “корпуса” первый запуск обучения получает “STOP It is” в качестве входных слов и “it is a” — в качестве трех слов для прогнозирования.

На рис. 4.6 демонстрируется, что мы находимся на втором вызове, когда поступают слова “a small world”, и необходимо спрогнозировать “small world but”. В начале второго прямого прохода выход первого вызова поступает слева (O0), соединяется с векторным представлением “a” и передается через сеть RNN, где становится O1, передавая потери на E1.

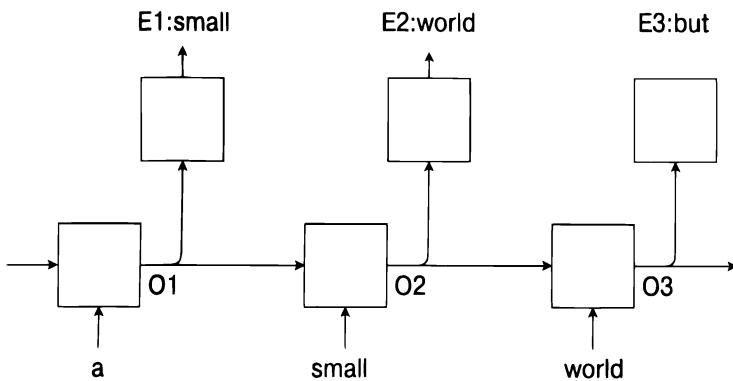


Рис. 4.6. Обратное распространение во времени с размером окна, равным 3

Однако, помимо перехода к E1, O1 объединяется со вторым словом “small”. Мы также вычислили ошибку. Здесь мы вычисляем влияние как \mathbf{W} , так и \mathbf{b} (не говоря уже векторных представлениях) на ошибку в прогнозировании слова “small”. Но \mathbf{W} и \mathbf{b} вызывают ошибку двумя различными способами — в основном непосредственно из-за ошибки, которая ведет от них к E2, а также из-за того, как они повлияли на O1. Естественно, когда мы в следующий раз вычисляем E3, \mathbf{W} и \mathbf{b} влияют на ошибку тремя способами: непосредственно из O3, из O1 и из O2. Таким образом, параметры в этих переменных изменяются шесть раз (или, что эквивалентно, программа сохраняет промежуточный итог и изменяет их один раз).

На рис. 4.6 проблема размера партии игнорируется. Как и следовало ожидать, функции TF RNN созданы для одновременного обучения и тестирования. Таким образом, каждый вызов RNN получает массив входных данных размером в размер партии на размер окна, чтобы вернуть массив спрогнозированных слов того же размера. Как отмечалось ранее, семантика этого заключается в том, что мы работаем с группами размера партии параллельно, поэтому последние слова, спрогнозированные из первого обучающего вызова, являются первыми входными словами для второго.

Чтобы это сработало, нам нужно быть осторожными при создании партии. На рис. 4.7 показано, что происходит с нашим корпусом: “STOP It is a small world but I like it that way STOP”. Мы подразумеваем размер партии 2 и размер окна 3. Основная идея заключается в том, чтобы сначала разделить корпус на две части, а затем заполнить каждую партию фрагментами из каждой части (в данном случае — половины) корпуса. В верхнем окне на рис. 4.7 демонстрируется корпус, разделенный на две части. Следующая пара окон демонстрирует две входные партии, которые созданы из него. Каждая партия имеет три сегмента слов из каждой половины. Также нужно объединить спрогнозированные

слова для передачи в сеть. Это то же самое, что и на приведенном выше рисунке, но каждое слово уже следующее в корпусе. Таким образом, верхняя строка диаграммы прогноза гласит: “It is a small world but”.

STOP	It	is	a	small	world
but	I	like	it	that	way

STOP	It	is
but	I	like

a	small	world
it	that	way

Рис. 4.7. Распределение слов, когда размер партии равен 2, а размер окна — 3

Поскольку “корпус” состоит из 14 слов, каждая половина состоит из шести слов. Чтобы понять, почему шесть, а не семь, сосредоточьтесь на прогнозах для второй партии. Внимательно прочитайте семь слов по половинам, и вы обнаружите, что нет прогноза слова для последнего ввода. Таким образом, корпус первоначально делится на S секций, где для корпуса размера x и размера партии b , $S = \lfloor (c - 1)/b \rfloor$ (где “ $\lfloor x \rfloor$ ” — это *функция пола* (floor function) — наибольшее целое число, меньше чем x). Здесь “минус 1” дает вводу последнее соответствующее прогнозируемое слово.

Мы пока ничего не сказали о том, что мы делаем в конце предложения. Самое простое — это сразу перейти к следующему. Это означает, что данный сегмент размера окна, передаваемый в RNN, может содержать фрагменты двух разных предложений. Однако мы поместили между ними слово “STOP”, поэтому RNN, в принципе, должен узнатъ, что это означает с точки зрения того, какие слова встречаются, например, с прописными буквами. Кроме того, могут быть хорошие подсказки о последующих словах из последних слов предыдущего предложения. Если мы просто занимаемся моделированием языка, то достаточно разделить предложения с помощью слова “STOP”, но не беспокоиться о разделении предложений при обучении или использовании RNN.

Давайте вернемся к RNN, снова рассмотрев рис. 4.5 и 4.6 и подумав, как мы программируем языковую модель RNN. Как мы только что заметили, код, перенесший нас из нашего корпуса слов в модель ввода, должен быть немного переработан. Ранее входные данные (и прогнозы) представляли собой вектор размера партии, теперь он представляет собой массив размера партии на размер окна. Нам также нужно превратить каждое слово в его векторное представление, но это не отличается от модели прямой связи.

Далее слово подается в RNN. Ключевой код TF для создания RNN таков:

```
rnn = tf.contrib.rnn.BasicRNNCell(rnnSz)
initialState = rnn.zero_state(batchSz, tf.float32)
outputs, nextState = tf.nn.dynamic_rnn(rnn, embeddings,
                                       initial_state=initialState)
```

В первой строке RNN добавляется к графу вычислений. Обратите внимание, что ширина массива весов RNN является свободным параметром, `rnnSz` (вы, может быть, помните, что, когда мы добавили дополнительный слой линейных блоков в модель Mnist в конце главы 2, у нас была похожая ситуация). Последняя строка — это вызов RNN. Он получает три аргумента и возвращает два. Входными данными являются, во-первых, собственно RNN, во-вторых, слова, которые RNN собирается обработать (размер партии на размер окна), и `initialState`, полученный из предыдущего вызова. Поскольку при первом вызове `dynamic_rnn` предыдущего состояния нет, мы создаем его фактично, вызовом функции в строке 2: `rnn.zero_state`.

Вызов `tf.nn.dynamic_rnn` имеет два вывода. Первый мы назвали `outputs`; это информация, которая передается для вычисления ошибок. На рис. 4.6 это выводы O1, O2, O3. Таким образом, `output` имеет форму [размер партии, размер окна, скрытый размер]. Первое измерение упаковывает примеры размера партии. Каждый пример состоит из O1, O2 и O3, поэтому второе измерение имеет размер окна. Наконец, O1, например, является вектором чисел с плавающей запятой размером `rnn`, которые составляют вывод RNN для одного слова.

```
[[[-0.077  0.022 -0.058 -0.229  0.145]
 [-0.167  0.062  0.192 -0.310 -0.156]
 [-0.069 -0.050  0.203  0.000 -0.092]]

 [[[ -0.073 -0.121 -0.094 -0.213 -0.031]
   [-0.077  0.022 -0.058 -0.229  0.145]]
 [[ 0.179  0.099 -0.042 -0.012  0.175]
   [-0.167  0.062  0.192 -0.310 -0.156]]
 [[ 0.103  0.050  0.160 -0.141 -0.027]
   [-0.069 -0.050  0.203  0.000 -0.092]]]
```

Рис. 4.8. Вывод RNN для next_state и outputs

Второй вывод от `tf.nn.dynamic_rnn` мы назвали `nextState`, и это последний вывод (O3) данного прохода RNN. В следующий раз, когда мы вызываем `tf.nn.dynamic_rnn`, у нас есть `initialState = nextState`. Обратите внимание, что `nextState` — это фактически информация, которая присутствует в `outputs`, поскольку она представляет собой набор O3 из примеров размера партии. Например, на рис. 4.8 показаны `next_state` и `outputs` для партии

размером 3, размером окна 2 и размером rnn 5. При размере окна 2 каждая вторая строка вывода является строкой следующего состояния. Нам, конечно, удобно упаковывать следующее состояние отдельно, но реальная причина этого повторения станет ясна в следующем разделе.

Последний элемент языковой модели — это потери. Они вычисляются в правой верхней части рис. 4.5. Как мы видим, вывод RNN сначала пропускается через слой линейных блоков, чтобы получить логиты для softmax, а затем мы вычисляем кросс-энтропийные потери из вероятностей. Как только что упоминалось, вывод RNN представляет собой трехмерный тензор с формой [размер партии, размер окна, размер rnn]. До сих пор мы передавали только двумерные тензоры, матрицы, через наши линейные блоки, и делали это с умножением матриц, например `tf.matmul(inpt, W)`.

Самый простой способ справиться с этим — изменить форму выходного тензора RNN так, чтобы сделать его матрицей с правильными свойствами:

```
output2 = tf.reshape(output, [batchSz*windowSz, rnnSz])
logits = matmul(output2, W)
```

Здесь **W** — это линейный слой (**W₀**), получающий выходные данные RNN и превращающий их в логиты на рис. 4.5. Затем мы можем передать это функции `tf.nn.sparse_softmax_cross_entropy_with_logits`, которая возвращает вектор столбца значений потерь, которые сводятся до одного значения с помощью `tf.reduce_mean`. Это окончательное значение может быть возведено в степень, чтобы дать нам нашу перплексивность.

Изменение формы вывода RNN здесь было удобно и по педагогическим причинам (это позволило бы повторно использовать `tf.matmul`), и по вычислительным (это приводит все в форму, требуемую `sparse_softmax`). В других ситуациях последующие вычисления могут потребовать исходной формы. Для этого мы можем обратиться к одной из многих функций TF, которые обрабатывают многомерные тензоры. Здесь мы использовали бы ту, которая описана в разделе 2.4.2. Ее вызов будет таков:

```
tf.tensordot(outputs, W, [[2], [0]])
```

Этот код выполняет повторное скалярное произведение (по сути, умножение матриц) между вторым компонентом (начинаясь с нуля) `outputs` и нулем `W`.

Еще один момент по поводу использования общей сети RNN: в коде Python, сопровождающем приведенный выше код TF для RNN, мы видим что-то вроде

```
inputSt = sess.run(initialSt)
for i in range(numExamps)
    'read in words and embed them'
```

```
logts, nxts=sess.run([logits,nextState],
                      {{input=wrds, nextState=inputSt}})
inputSt=nxts
```

Ничто здесь не должно восприниматься буквально, за исключением а) того, как инициализируется состояние ввода для RNN, б) как мы передаем его в TF с помощью фрагмента словаря nextState=inputState и в) как мы затем обновляем inputSt в последней строке выше. До сих пор мы использовали только feed_dict для передачи значений в заполнители TF. Здесь nextState указывает не на заполнитель, а на фрагмент кода, который генерирует нулевое состояние, с которого мы запускаем RNN. Это разрешено.

4.6. Долгая краткосрочная память

Долгая краткосрочная память (Long Short-Term Memory — LSTM) NN — это особая разновидность сети RNN, которая почти всегда превосходит простую RNN, представленную в последнем разделе.

Проблема со стандартными RNN заключается в том, что, хотя они должны помнить вещи из далекого прошлого, на практике они, кажется, быстро все забывают. На рис. 4.9 все, что заключено в пунктирную рамку, соответствует одному блоку RNN. Очевидно, что архитектуру LSTM разрабатывают довольно тщательно. В первую очередь, обратите внимание, что мы показали одну копию LSTM на *диаграмме обратного распространения во времени* (back-propagation diagram). Итак, слева есть информация, поступающая от обработки предыдущего слова (используя два тензора информации, а не один). Внизу появляется следующее слово. Справа есть два тензора, которые выводят информацию в следующий блок времени, и, как и в простых RNN, эта информация идет вверх по диаграмме, чтобы спрогнозировать следующее слово и потери (верхняя правая сторона).

Цель — улучшить память RNN о прошедших событиях, научив ее запоминать важные вещи и забывать все остальное. Для этого LSTM проходят две версии прошлого. “Официальная” селективная память находится вверху, а более локальная версия — внизу. Верхняя времененная шкала памяти называется *состоянием ячейки* (cell state), сокращенно — *c*. Нижняя линия называется *h*.

На рис. 4.9 представлено несколько новых функций подключения и активации. В первую очередь, мы видим, что линия памяти изменяется в двух местах, а затем передается в следующий блок времени. Они помечены как (X) и плюс (+). Идея в том, что память удаляется в *блоке времени* (time unit) и добавляется в *блоке “плюс”* (plus unit).

Почему мы говорим об этом? Посмотрите сейчас на текущее векторное представление слова, которое расположено внизу слева. Оно проходит через

слой линейных блоков, за которыми следует сигмовидная функция активации, на что указывают подписи \mathbf{W} , \mathbf{b} , S : \mathbf{W} и \mathbf{b} означают линейный блок, а S — сигмовидную функцию. Мы показали сигмовидную функцию на рис. 2.7. Возможно, вы захотите просмотреть ее, потому что некоторые из ее особенностей имеют значение в следующем обсуждении. В математической форме записи есть операция

$$\mathbf{h}' = \mathbf{h}_t \cdot \mathbf{e} \quad (4.11)$$

$$\mathbf{f} = S((\mathbf{h}' \mathbf{W}_f + \mathbf{b}_f)) \quad (4.12)$$

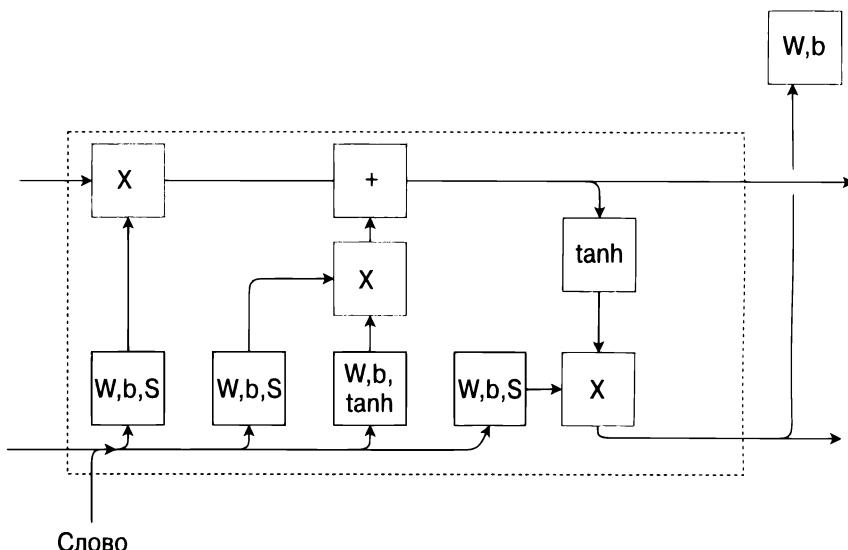


Рис. 4.9. Архитектура LSTM

Мы используем центральную точку · для обозначения конкатенации векторов. Повторим, что в нижнем левом углу мы объединяем предыдущую h -строку \mathbf{h}_t и текущее векторное представление слова \mathbf{e} , чтобы получить \mathbf{h}' , которое, в свою очередь, передается в линейный блок “забывания” (сопровождаемый сигмовидной функцией), чтобы получить \mathbf{f} , *сигнал забывания* (forgetting signal), который движется вверх по левой стороне рисунка.

Выход сигмовидной функции затем поэлементно умножается на с-линию памяти, идущую вверх слева. (Под “поэлементно” мы подразумеваем, что, например, $x[i, j]$ -й элемент одного массива умножается (или добавляется и т.д.) на $y[i, j]$ -й элемент другого.)

$$\mathbf{c}'_t = \mathbf{c}_t \odot \mathbf{f} \quad (4.13)$$

(Здесь “ \odot ” означает поэлементное умножение.) Учитывая, что сигмовидные функции ограничены нулем и единицей, результатом умножения должно быть уменьшение абсолютного значения в каждой точке основной памяти. Это соответствует “забыванию”. В целом эта конфигурация, которую сигмовидная функция передает на мультиплексивный слой, является обычной моделью, когда мы хотим “мягкого” управления (gating).

Сравните это с тем, что происходит в аддитивном блоке, с которым в следующий раз встречается память. Опять же, следующее векторное представление слова пришло снизу слева, и на этот раз оно проходит отдельно через два линейных слоя, один — с сигмовидной функцией активации, другой — с функцией активации \tanh (\tanh activation function), показанной на рис. 4.10. “Tanh” означает гиперболический тангенс (hyperbolic tangent).

$$\mathbf{a}_1 = S(\mathbf{h}' \mathbf{W}_{\mathbf{a}_1} + \mathbf{b}_{\mathbf{a}_1}) \quad (4.14)$$

$$\mathbf{a}_2 = \tanh((\mathbf{h}_t \cdot \mathbf{e}) \mathbf{W}_{\mathbf{a}_2} + \mathbf{b}_{\mathbf{a}_2}) \quad (4.15)$$

Важно, что, в отличие от сигмовидной функции, функция \tanh может выводить как положительные, так и отрицательные значения, поэтому она может выражать новый материал, а не только масштаб. Ее результат добавляется к состоянию ячейки в ячейке, помеченной знаком “+”:

$$\mathbf{c}_{t+1} = \mathbf{c}'_t \oplus (\mathbf{a}_1 \odot \mathbf{a}_2) \quad (4.16)$$

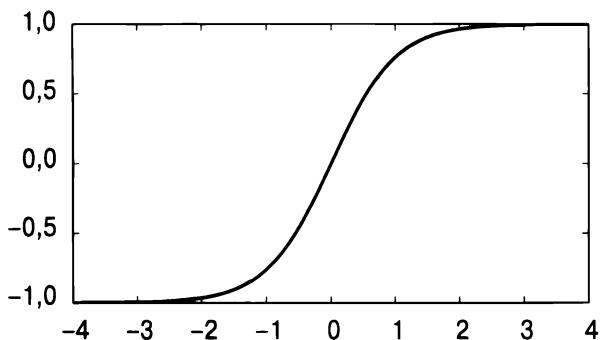


Рис. 4.10. Функция \tanh

После этого линия памяти ячейки разделяется. Одна копия выходит вправо, а другая копия проходит через функцию \tanh , а затем объединяется с линейным преобразованием более локальной истории/векторного представления, чтобы стать новой \mathbf{h} -линией внизу:

$$\mathbf{h}'' = \mathbf{h}' \mathbf{W}_h + b_h \quad (4.17)$$

$$\mathbf{h}_{t+1} = \mathbf{h}'' \odot \mathbf{a}_2 \quad (4.18)$$

Это должно быть объединено со следующим векторным представлением слова, и процесс повторяется. Здесь следует подчеркнуть, что линия ячейки памяти никогда не идет напрямую через линейные блоки. Все ослабляется (*de-emphasized*) (например, забывается) в блоке “X” и добавляется в блоке “+”, но это так. Такова логика механизма LSTM.

Что касается программы, требуется только одно небольшое изменение в версии TF:

```
tf.contrib.rnn.BasicRNNCell(rnnSz)
```

становится

```
tf.contrib.rnn.LSTMCell(rnnSz)
```

Обратите внимание, что это изменение влияет на состояние, которое передается из одного блока времени в другой. Ранее, как показано на рис. 4.8, состояния имели форму [batchSz, rnnSz]. Теперь это [2, batchSz, rnnSz], один тензор [batchSz, rnnSz] для с-линии и один — для h-линии.

По производительности версия LSTM намного лучше за счет траты на обучение большего времени. Возьмите модель RNN, подобную той, которую мы разработали в предыдущем разделе, выделите ей много ресурсов (векторы векторных представлений слов размером 128, скрытый размер 512), и получите респектабельную перплексивность 120 или около того. Измените один вызов функции с RNN на LSTM, и наша перплексивность снизится до 101.

4.7. Ссылки и что читать дальше

Документ, который представил то, что мы сейчас считаем стандартной языковой моделью с прямой связью, — это работа Бенджо и других [18]. Также возник термин “векторное представление слов”. Идея представления слов в непрерывном пространстве, в частности векторов чисел, появилась намного раньше. Однако именно в статье Бенджо показано, что векторное представление слов в том виде, в котором мы их теперь знаем, является почти автоматическим побочным продуктом языковых моделей NN.

Можно утверждать, что не ранее работы Миколова и других это векторное представление слов стало почти универсальным компонентом обработки естественного языка сетью NN. Они разработали несколько моделей, которые называются *word2vec*. Стандартная публикация — [19]. Самая популярная из моделей *word2vec* — это модель *skip-gram* (*skip-gram model*). В нашей презентации векторные представления слов были оптимизированы для прогнозирования следующего слова с учетом предыдущих слов. В модели скип-грамм каждое отдельное слово просят спрогнозировать все соседние слова. Одним

из поразительных результатов моделей word2vec является использование векторного представления слов для решения задачи *поиска аналогии слов* (word analogy problem), например мужчина может быть королем, а женщина — кем? Неожиданно ответы на эти вопросы выпали из созданного векторного представления слов. Кто-то просто взял векторное представление слова “король”, вычел его для “мужчины”, добавил “женский”, а затем искал слово как векторное представление, ближайшее к результату. Отличный блог о векторном представлении слов и их проблемах ведется Себастьяном Рудером [20].

Рекуррентные нейронные сети существуют примерно с середины 1980-х годов, но они не работали до тех пор, пока Зепп Хохрайтер и Юрген Шмидхубер не создали LSTM [21]. Блог Криса Кола дает хорошее объяснение LSTM [22], и мой рис. 4.9 представляет собой переработку одной из его диаграмм.

4.8. Упражнения

Упражнение 4.1. Предположим, что наш корпус начинается словами “*STOP* I like my cat and my cat likes me . *STOP*” (Я люблю мою кошку, и моя кошка любит меня.). Также предположим, что мы присваиваем отдельным словам уникальное целое число по мере их чтения в корпусе, начиная с 0. Если у нас есть размер партии 5, запишите значения, которые мы должны прочитать, чтобы заменить заполнители `inpu` и `answr` в первой обучающей партии.

Упражнение 4.2. Объясните, почему, если вы надеетесь получить хороший шанс изучить хорошую языковую модель на основе векторного представления слов, вы можете не устанавливать все `E` равными нулю. Убедитесь, что ваше объяснение также справедливо для установки всех `E` равными единице.

Упражнение 4.3. Объясните, почему при использовании регуляризации L2 вычисление фактических общих потерь является плохой идеей.

Упражнение 4.4. Подумайте о построении полносвязной модели триграмм языка. В нашей версии мы объединили векторные представления для двух предыдущих входов, чтобы сформировать вход модели. Влияет ли порядок, в котором мы объединяем, на способность модели обучаться? Объясните.

Упражнение 4.5. Рассмотрим модель униграмм NN. Может ли перплексивность этой модели быть лучше, чем при выборе слов из равномерного распределения? Почему “да” или почему “нет”? Объясните, какие части модели биграмм необходимы для оптимальной работы модели униграмм.

Упражнение 4.6. *Линейный управляемый блок (LGU — Linear Gated Unit) является вариантом LSTM. Возвращаясь к рис. 4.9, мы видим, что последний имеет один скрытый слой, который контролирует то, что удаляется из основной линии памяти, и второй, который контролирует то, что добавляется. В обоих случаях слои получают нижнюю линию управления в качестве входа и создают вектор чисел от 0 до 1, который умножается на линию памяти (забывается) или добавляется к ней (запоминается). LGU отличаются тем, что эти два слоя заменены одним слоем с одинаковым входом. Вывод умножается на линию управления, как и раньше. Однако он также вычитается из единицы, умножается на уровень управления и добавляется в линию памяти. В целом LGU работают так же, как LSTM, и, имея на порядок меньше линейных слоев, немного быстрее. Измените рис. 4.9 так, чтобы он представлял работу LGU.*

Глава 5

Обучение от последовательности к последовательности

Обучение от последовательности к последовательности (sequence-to-sequence learning или сокращенно seq2seq) представляет собой метод глубокого обучения для отображения одной последовательности символов на другую последовательность символов, когда невозможно (или по крайней мере мы не видим, как) выполнить отображение на основе отдельных символов. Прототипом приложения для seq2seq является *машинный перевод* (Machine Translation — MT) — компьютерный перевод между естественными языками, такими как французский и английский.

Примерно с 1990 года было признано, что выразить это отображение программно довольно сложно и что более косвенный подход работает намного лучше. Мы даем компьютеру *выровненный корпус* (aligned corpus) — множество примеров пар предложений, которые являются взаимными переводами одно другого, и требуем, чтобы машина сама выявила соответствие. Вот где начинается глубокое обучение. К сожалению, техники глубокого обучения, которую мы изучили для задач на естественном языке, например LSTM, самой по себе недостаточно для MT.

Важно отметить, что языковое моделирование (задача, на которой мы сконцентрировались в предыдущей главе) происходит поэтапно, т.е. мы передаем одно слово и прогнозируем следующее. MT работает не так. Рассмотрим некоторые примеры из *Canadian Hansard*'s, отчета обо всем, что было сказано в канадском парламенте, что по закону должно быть опубликовано на двух официальных языках Канады, французском и английском. Первая пара предложений, которая у меня под рукой, такова:

edited hansard number 1
hansard révisé numéro 1
(отредактированный hansard номер 1)

Первый урок для англоязычных людей, изучающих французский язык (и, вероятно, наоборот), заключается в том, что прилагательные обычно располагаются перед существительными, к которым они относятся в английском языке, и после них — во французском. Так что здесь прилагательные “edited” и “révisé” находятся в переводах не в одинаковых позициях. Дело в том, что мы не можем продвигаться слева направо по фразе на *исходном языке* (source language) (языке, с которого мы переводим), выдавая по одному слову за раз на *языке назначения* (target language). В данном случае мы могли бы ввести два слова и вывести два, но наблюдаемые последовательные несоответствия могут быть намного больше. Следующее происходит через несколько строк после предыдущего примера. Обратите внимание, что текст *лексически проанализирован* (tokenized) — французская пунктуация дважды отличается от слов, к которым она обычно прикрепляется:

this being the day on which parliament was convoked by proclamation of his excellency ...

parlement ayant été convoqué pour aujourd’hui, par proclamation de son excellence...

(это был день, когда парламент был создан по требованию его превосходительства...)

Пословный перевод с французского языка будет “parliament having been convoked for today, by proclamation of his excellency” (парламент был создан на сегодня по требованию его превосходительства), и, в частности, “this being the day” (это был день) переводится как “aujourd’hui” (сегодня). (Действительно, как правило, длина предложений в паре не одинакова.) Таким образом, требуется обучение от последовательности к последовательности, где последовательность обычно является целым предложением.

5.1. Парадигма Seq2Seq

Схема очень простой модели seq2seq показана на рис. 5.1. Она демонстрирует процесс во времени (время, как обычно, идет слева направо) перевода “hansard révisé numero 1” в “edited hansard number 1”. Модель состоит из двух сетей RNN. В отличие от LSTM, мы предполагаем модель RNN, которая проходит одну линию памяти. Мы могли бы использовать BasicRNNCell; однако лучший выбор — более новый конкурент LSTM, управляемый рекуррентный блок (Gated Recurrent Unit — GRU), который пропускает только одну строку памяти между блоками времени (time unit).

Модель работает в два прохода, каждый из которых имеет собственный GRU. Первый проход, *проход кодирования* (*encoding pass*), представлен в нижней половине рис. 5.1. Проход завершается, когда последний французский токен (всегда слово “STOP”) обрабатывается нижним GRU. Затем состояние GRU передается во вторую половину процесса. Цель этого прохода — создать вектор, который “суммирует” предложение. Это иногда называют *векторным представлением предложений* (*sentence embedding*) по аналогии с векторным представлением слов.

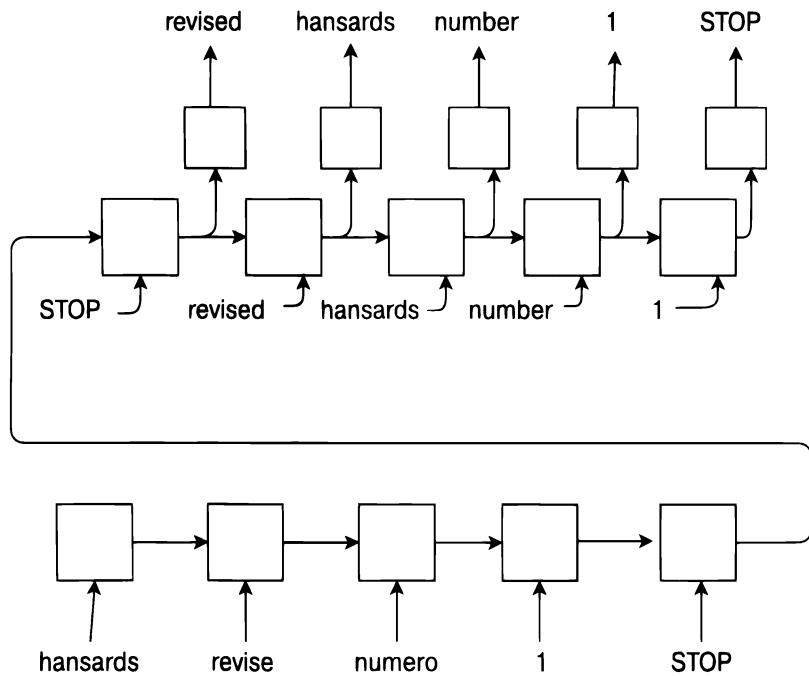


Рис. 5.1. Простая модель обучения от последовательности к последовательности

Вторая половина обучения seq2seq — это *проход декодирования* (*decoding pass*). Как видно из рисунка, это проход предложения языка назначения (здесь — английского). (Возможно, пришло время прямо указать, что мы говорим здесь о том, что происходит во время обучения, поскольку мы знаем английское предложение.) На этот раз задача заключается в том, чтобы спрогнозировать следующее английское слово после ввода каждого слова. Функция потерь — это обычная кросс-энтропийная потеря.

Термины *кодирование* и *декодирование* происходят из теории связи. Сообщение должно быть закодировано для формирования отправляемого сигнала, а затем снова декодировано из полученного сигнала. Если в связи есть шум,

полученный сигнал не обязательно будет идентичным отправленному. Представьте, что оригинальное сообщение было английским, а шум преобразовал его во французское. Тогда процесс его перевода (обратно) на английский язык будет декодированием.

Обратите внимание, что первое входное слово для прохода декодирования — это дополняющее слово “STOP”. Это также последнее слово вывода. Если бы мы использовали модель для настоящего французско-английского МТ, система не имела бы доступного английского языка. Но можно предположить, что мы начинаем обработку со слова “STOP”. Затем для генерации каждого последующего слова мы передаем LSTM предыдущее спрогнозированное слово. Обработка останавливается, когда LSTM прогнозирует, что следующее “слово” снова должно быть “STOP”. Естественно, мы должны так же работать при тестировании. Тогда мы знаем английский, но мы используем эту информацию только для оценки. Это означает, что при реальном переводе мы прогнозируем следующее слово перевода частично на основе последнего переведенного слова, *что вполне может быть ошибкой*. Если это так, то значительно возрастает вероятность того, что следующее слово будет неправильным, и т.д.

В этой главе мы игнорируем сложности реального МТ, оценивая способность нашей программы правильно прогнозировать следующее английское слово, учитывая правильное предыдущее слово. Конечно, в реальном МТ не существует единственно правильного перевода на английский язык для конкретного французского предложения, поэтому то, что наша программа не прогнозирует точное слово, используемое в переводе нашего параллельного корпуса, не означает, что это неправильно. Объективная оценка МТ является важной темой, но мы ее здесь игнорируем.

Последнее упрощение, прежде чем мы рассмотрим написание программы NN МТ. Рис. 5.1 представлен в виде схемы обратного распространения во времени, поэтому все блоки RNN в нижнем ряду фактически являются одинаковыми повторяющимися блоками, но в последовательные моменты времени. То же самое относится к блокам в верхнем ряду. Как вы, возможно, помните, модели обратного распространения во времени имеют гиперпараметр размера окна. В МТ мы хотим обработать все предложение сразу, но предложения бывают разных размеров. (В Penn Treebank они варьируются от одного слова до более чем 150.) Мы упростили нашу программу, работая только с предложениями, в которых французский и английский текст имеют длину менее 12 слов, или 13, включая слово “STOP”. Кроме того, мы делаем каждое предложение длиной 13, добавляя дополнительные слова “STOP”. Таким образом, программа может предполагать, что все предложения имеют одинаковую длину — 13 слов. Итак, рассмотрим короткие выровненные французско-английские предложения, с которых мы начали наше обсуждение МТ: “edited hansard number 1”

и “hansard révisé numéro 1”. Французское предложение, которое мы вводим, будет выглядеть так:

hansard révisé numéro 1 STOP STOP STOP STOP STOP STOP STOP STOP
STOP

А вот английское:

STOP edited hansard number 1 STOP STOP STOP STOP STOP STOP STOP
STOP

5.2. Написание программы Seq2Seq MT

Давайте начнем с обзора моделей RNN, которые мы рассмотрели в главе 4, но с небольшим отклонением. До сих пор мы не обращали особого внимания на хорошие методы разработки программного обеспечения. Однако здесь TF по причинам, которые должны быть объяснены, вынуждает нас навести порядок. Поскольку мы создаем две почти идентичные модели RNN, мы вводим конструкцию TF `variable_scope`. На рис. 5.2 показан код TF для двух RNN, которые нам нужны в нашей простой модели seq2seq.

```

1 with tf.variable_scope("enc") :
2     F = tf.Variable(tf.random_normal((vfSz, embedSz), stddev=.1))
3     embs = tf.nn.embedding_lookup(F, encIn)
4     embs = tf.nn.dropout(embs, keepPrb)
5     cell = tf.contrib.rnn.GRUCell(rnnSz)
6     initState = cell.zero_state(bSz, tf.float32)
7     encOut, encState = tf.nn.dynamic_rnn(cell, embs,
8                                         initial_state=initState)
9
10 with tf.variable_scope("dec") :
11    E = tf.Variable(tf.random_normal((veSz, embedSz), stddev=.1))
12    embs = tf.nn.embedding_lookup(E, decIn)
13    embs = tf.nn.dropout(embs, keepPrb)
14    cell = tf.contrib.rnn.GRUCell(rnnSz)
15    decOut, _ = tf.nn.dynamic_rnn(cell, embs, initial_state=encState)
```

Рис. 5.2. Код TF для двух RNN в модели MT

Мы делим код на две части; первая создает кодирование RNN, а вторая — декодирование. Каждый раздел заключен в команду TF `variable_scope`. Эта функция получает один аргумент — строку, служащую именем для области. Цель функции `variable_scope` — позволить нам собрать группу команд таким образом, чтобы избежать конфликтов имен переменных. Например, как

верхний, так и нижний сегменты используют имя переменной `cell` таким образом, что без двух отдельных областей они вступили бы в конфликт друг с другом, что очень плохо закончится.

Но даже если бы мы были осторожны и присвоили каждой из наших переменных уникальные имена, этот код *все равно* не работал бы правильно. По причинам, скрытым в деталях кода TF, когда функция `dynamic_rnn` создает материал для вставки в граф TF, она всегда использует одно и то же имя для указания на него. Если мы не поместим два вызова в отдельные области (или если код не настроен так, чтобы не учитывать, что эти два вызова, по сути, являются одним и тем же), мы получим сообщение об ошибке.

Теперь рассмотрим код в каждой области видимости переменной. Для кодировщика мы сначала создаем пространство для векторного представления французских слов `F`. Мы предполагаем заполнитель с именем `encIn`, который получает тензор индексов французских слов с формой “размер партии на размер окна”. Затем функция поиска возвращает трехмерный тензор с формой “размер партии на размер окна на размер векторного представления” (строка 3), к которому мы применяем исключение с вероятностью сохранения соединения, установленного в `keepProb` (строка 4). Затем мы создаем ячейку RNN, на этот раз используя вариант GRU для LSTM. Стока 7 далее использует ячейку для получения выходных данных и следующего состояния.

Второй GRU параллелен первому, за исключением того, что вызов функции `dynamic_rnn` получает в качестве входа вывод состояния кодировщика RNN, а не нулевое начальное состояние. Это `state=encState` в строке 15. Снова обратимся к рис. 5.1: пословный вывод декодера RNN подается на линейный слой. На рисунке не показано, но читатель должен представить, что вывод слоя (логиты) передается на вычисление потерь. Код будет выглядеть так, как показано на рис. 5.3. Единственным нововведением здесь является вызов `seq2seq_loss`, специализированной версии функции кросс-энтропийных потерь в случаях, когда логиты являются трехмерными тензорами. Она получает три аргумента. Первые два стандартны — логиты и двумерный тензор правильных ответов (размер партии на размер окна). Третий аргумент подразумевает взвешенную сумму — для ситуаций, когда одни ошибки должны иметь большее значение для общих потерь, чем другие. В нашем случае мы хотим, чтобы каждая ошибка учитывалась одинаково, поэтому третий аргумент имеет все веса, равные 1.

```
W = tf.Variable(tf.random_normal([rnnSz, veSz], stddev=.1))
b = tf.Variable(tf.random_normal([veSz, stddev=.1]))
logits = tf.tensordot(decOut, W, axes=[[2], [0]])+b
loss = tf.contrib.seq2seq.sequence_loss(logits, ans,
                                         tf.ones([bSz, wSz]))
```

Рис. 5.3. Код TF для декодера seq2seq

Как мы уже упоминали, вся идея простейшей модели seq2seq на рис. 5.1 заключается в том, что этап кодирования создает “сводку” французского предложения, передавая французский текст через GRU, а затем используя итоговый вывод состояния GRU в качестве сводки. Однако существует множество разных способов создания таких предложений, и значительная часть исследований была посвящена рассмотрению этих альтернатив. На рис. 5.4 приведена вторая модель, которая для МТ немного лучше. Разница в реализации небольшая: вместо того чтобы передавать конечное состояние кодера в декодер в качестве его начального состояния, мы скорее получаем сумму всех состояний кодера. Поскольку мы дополнили все французские и английские предложения до длины 13, мы получаем все 13 состояний и суммируем их. Остается надеяться, что эта сумма более информативна, чем один последний вектор, что на самом деле, кажется, и имеет место.

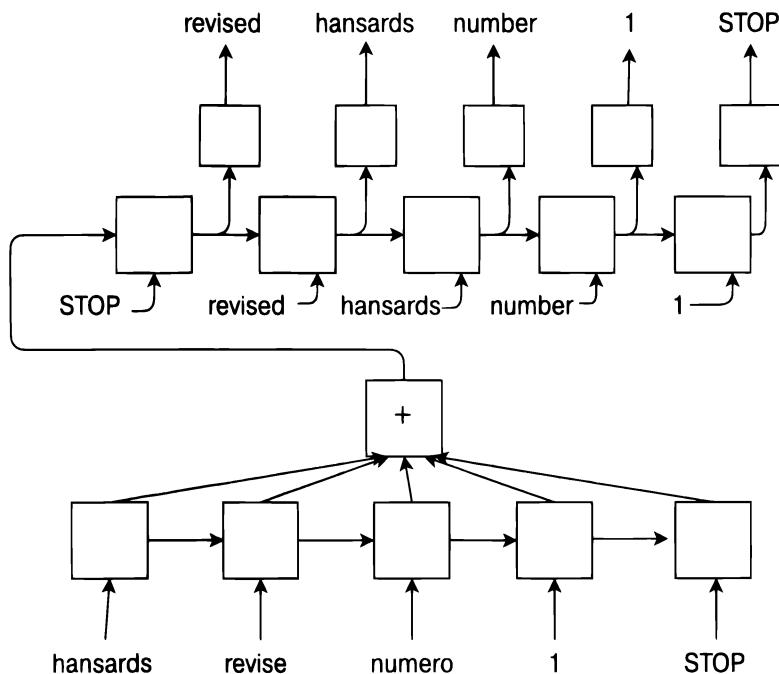


Рис. 5.4. Сводка предложений Seq2seq в ходе сложения

На самом деле ваш автор выбрал среднее значение векторов состояния, а не сумму. Вернувшись к главе 1 и посмотрев на вычисление прямого прохода, вы вспомните, что получение среднего значения, а не суммы не имеет различий в конечных вероятностях, поскольку softmax удаляет любые мультипликативные различия, а получение среднего значения просто соответствует делению на размер окна (13). Кроме того, направление градиента параметра также не

изменяется. То, что может изменяться (и действительно изменяется), является величиной изменения, которое мы вносим в значения параметров. В текущей ситуации получение суммы примерно эквивалентно умножению скорости обучения на 13. Как правило, в таких ситуациях лучше держать значения параметров около нуля и напрямую изменять скорость обучения.

5.3. Внимание в Seq2seq

Понятие *внимания* (attention) в моделях seq2seq вытекает из идеи, что, хотя в целом нам нужно понять все предложение, прежде чем мы сможем его перевести, на практике для данного фрагмента перевода одни части исходного предложения более важны, чем другие. В частности, гораздо чаще, чем наоборот, первые несколько французских слов соответствуют первым английским, середина французского предложения соответствует середине английского и т.д. Хотя это особенно верно для английского и французского языков, которые очень схожи, даже языки, у которых нет очевидных общностей, обладают этим свойством. Причиной является *данное новое различие* (given new distinction). Похоже, что это имеет место для всех языков, когда говорят что-то о новых вещах, о которых мы уже упоминали (а происходит это обычно в разговоре или письме): мы сначала говорим “дано” то, о чем мы говорили, и только потом упоминаем новый материал. Таким образом, в разговоре о Джеке мы могли бы сказать: “Jack ate a cookie” (Джек съел печенье), но если бы мы говорили о пачке печенья, то “One of the cookies was eaten by Jack” (Одно из печений съел Джек).

На рис. 5.5 показана небольшая вариация механизма суммирования seq2seq, показанного на рис. 5.4, в котором сводка, объединенная с векторным представлением английского слова, подается на ячейку декодера в каждой позиции окна. Это отличается от рис. 5.4, на котором вводится только английское слово. С нашей новой точки зрения, эта модель уделяет одинаковое внимание всем состояниям кодера при работе со всеми частями на английском языке. В моделях внимания мы модифицируем это так, чтобы разные состояния смешивались вместе в разных пропорциях перед передачей в декодер RNN. Мы называем это *вниманием только к позиции* (position-only attention). Истинные модели внимания куда сложнее, но мы оставляем их для самостоятельного изучения.

Итак, мы собираемся построить схему внимания, на которой, скажем, внимание на английское слово в позиции i соответствует состоянию французской кодировки в позиции j и зависит только от i и j . Как правило, чем ближе i и j , тем больше внимания. На рис. 5.6 представлена вымышленная матрица весов для модели, в которой размер окна кодера (французского) равен 4, а размер декодера равен 3. До этого мы предполагали, что оба размера окна равны 13, но нет никаких причин, по которым они должны быть одинаковыми. Кроме того, поскольку

во французском языке примерно на 10% больше слов, чем в соответствующем переводе на английский, есть веская причина соотнести немного большие размеры французских окон с меньшими английскими. Кроме того, по педагогическим соображениям асимметричная матрица помогает точно определить, какие числа относятся к позициям английского слова, а какие — французского.

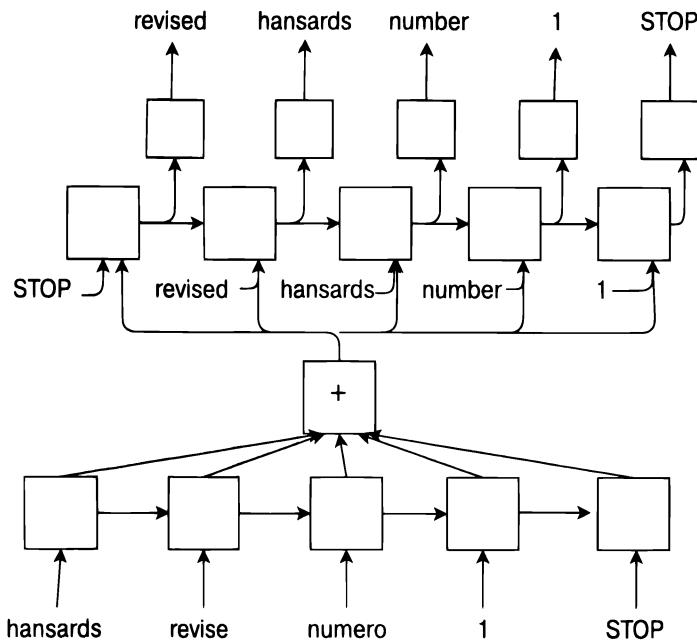


Рис. 5.5. Seq2seq, где сводка кодера подается непосредственно в каждую позицию окна декодера

1/2	1/6	1/6
1/6	1/3	1/6
1/6	1/3	1/3
1/6	1/6	1/3

Рис. 5.6. Возможная весовая матрица для взвешивания соответствующих французских/английских позиций

Здесь мы присваиваем $W[i, j]$ вес, приданый i -му французскому состоянию, когда он используется для прогнозирования j -го английского слова. Таким образом, общий вес для любого английского слова является суммой столбца, которую мы сделали 1. Например, в первом столбце приведены веса для первого английского слова. Глядя на первый столбец, мы видим, что первое французское состояние учитывает половину акцента (emphasis) (в нашем воображенном примере размер окна — 4), а остальные три французских состояния

совместно используют остальной акцент в равной степени. На данный момент мы предполагаем, что в нашей реальной программе на рис. 5.6 версии $13 * 13$ это константа TF.

Далее, учитывая весовую матрицу $13 * 13$, как мы можем использовать ее, чтобы варьировать внимание к конкретному английскому выводу? На рис. 5.7 показаны вычисления тензорного потока и числовые расчеты для ситуации, когда размер партии равен 2, размер окна — 3, а размер RNN — 4. В верхней части мы видим вывод воображаемого кодера `encOut`. Размер партии равен 2, и для простоты мы сделали две партии идентичными. В каждой партии у нас есть три вектора длиной 4, каждый из которых представляет собой вектор длиной 4 (размер RNN) для вывода RNN в этой позиции окна. Так, например, в партии 0 первый (отсчитываемый с 0) вектор состояния равен $(1, 1, 1, 1)$.

```

eo= ( (( 1,  2,  3,  4),
        ( 1,  1,  1,  1),
        ( 1,  1,  1,  1),
        ( -1, 0, -1, 0)),
      (( 1,  2,  3,  4),
        ( 1,  1,  1,  1),
        ( 1,  1,  1,  1),
        ( -1, 0, -1, 0)) )
encOut=tf.constant(eo, tf.float32)

AT = ( ( .6, .25, .25 ),
       ( .2, .25, .25 ),
       ( .1, .25, .25 ),
       ( .1, .25, .25 ) )
wAT = tf.constant(AT, tf.float32)

encAT = tf.tensordot(encOut,wAT,[[1],[0]])
sess= tf.Session()

print sess.run(encAT)
'''[[[ 0.80000001  0.5          0.5          ]
   [ 1.50000012  1.           1.           ]
   [ 2.           1.           1.           ]
   [ 2.70000005  1.5          1.5          ]
   ...]'''

decAT = tf.transpose(encAT,[0, 2, 1])
print sess.run(decAT)
'''[[[ 0.80000001  1.50000012  2.           2.70000005]
   [ 0.5          1.           1.           1.5          ]
   [ 0.5          1.           1.           1.5          ]
   ...]'''

```

Рис. 5.7. Упрощенные вычисления внимания при $bSz = 2$, $wSz = 3$ и $rnnSz = 4$

Далее у нас есть вектор готового веса, wAT , размером $4 * 3$ (размеры окна). Это позиция французского состояния по английскому слову. Таким образом, первый столбец в действительности говорит, что первому английскому слову должен быть задан вектор состояния, который составляет 60% первого французского вектора состояния и 20% каждого из двух других векторов состояния RNN. Это организовано так, чтобы веса для каждого английского слова составляли до 100%.

Далее следуют три команды TF, которые позволяют принимать невзвешенные состояния кодировщика и создавать взвешенные версии для каждого решения английского слова. Сначала мы изменим выходной тензор кодера с $[bSz, wSz, rnnSz]$ на $[bSz, rnnSz, wSz]$. Это осуществляется командой `tf.transpose`. Команда `transpose` получает два аргумента: тензор для транспонирования и заключенный в скобки вектор целых чисел, определяющих транспонирование, которое должно быть выполнено. Здесь мы запросили $[0, 2, 1]$, чтобы оставить 0-е измерение на месте, но “2” говорит, что измерение, которое изначально было вторым, стало измерением 1, а финальная единица делает последнее измерение тем, что раньше было первым. Мы выводим результат этого преобразования, когда выполняем вызов `print_sess.run(encOT)`.

Мы осуществили транспонирование, чтобы упростить умножение матриц на следующем этапе (`tensordot`). Фактически, если у нас нет усложнения размера партии, мы умножаем тензоры формы $[rnnSz, wSz] * [wSz, wSz]$ и мы могли бы использовать стандартное матричное умножение (`matmul`). Дополнительное измерение, обусловленное размером партии, исключает эту возможность, и мы возвращаемся к

```
encAT = tf.tensordot)(encOut.wAT, [[1], [0]])
```

Наконец, мы обратили транспонирование, которое сделали два этапа назад, чтобы вернуть выходные состояния кодера в их исходную форму.

Стоит немного остановиться, чтобы сравнить конечный результат в нижней части рис. 5.7 с выводом нашего воображаемого кодера в верхней части рисунка. В столбце $(0, 6, 0, 2, 0, 2)$ говорится, что первое английское слово означает вектор, состоящий из 60% “нулевого состояния”, т.е. $[1, 2, 3, 4]$. Поэтому мы ожидаем, что результирующее состояние будет увеличиваться при движении слева направо, что дает $(0, 6, 1, 4, 1, 8, 2, 6)$. Второе состояние, из всех 1, не имеет большого эффекта (оно добавляет 0,2 к каждой позиции). Но последнее состояние $(0, -1, 0, -1)$ должно давать результату модель вверх-вниз (up-and-down pattern), что оно и делает. (Все эти компоненты $(0, 6, 1, 4, 1, 8, 2, 6)$ увеличиваются, но первый и третий увеличиваются больше, чем второй.)

После того как у нас появляются повторно взвешенные состояния кодера для подачи в декодер, мы объединяем каждое с векторным представлением

английского слова, которое мы уже подали в декодер RNN. Это завершает нашу простую систему внимания МТ. Однако, чтобы завершить этот пример, важно одно: мы можем заставить нашу программу узнать веса внимания, просто сделав массив внимания 13×13 переменной TF, а не константой. Идея похожа на то, что мы делали, когда в главе 3 NN изучали ядра свертки. На рис. 5.8 показаны некоторые веса, изученные таким образом. Выделенные полужирным шрифтом цифры являются самыми большими числами в их ряду. Они демонстрируют ожидаемый сдвиг вправо — перевод слов в начале, середине или конце английского, как правило, должны уделять больше внимания началу, середине или концу французского.

-6.3	1.3	.37	.13	.06	.04	.11	.10	.02
-.66	-.44	.64	.26	.16	.02	.03	.04	.06
-.38	-.47	-.04	.63	.18	.10	.07	.06	.12
-.30	-.44	-.35	-.15	.48	.24	.06	.13	0
-.02	-.16	-.35	-.37	-.23	.12	.32	.22	.11
.05	-.11	-.11	-.35	-.04	-.22	.05	.26	.24
.10	.02	-.04	-.23	-.32	-.33	-.25	-.01	.28
0	.03	.01	-.18	-.21	-.26	-.30	-.1.1	-.17

Рис. 5.8. Веса внимания для верхнего левого угла 8×9 матрицы весов внимания 13×13

5.4. Seq2Seq с несколькими длинами окон

В последнем разделе мы ограничили наши пары перевода примерами, в которых оба предложения составляют 12 слов или менее (13, включая дополнительное слово “STOP”). В реальном МТ такие ограничения не допускаются. С другой стороны, наличие окна, скажем, 65, которое позволило бы переводить практически все предложения, которые появились на нашем пути, означало бы, что более обычное предложение будет дополнено 40 или 50 словами “STOP”. Решение, которое было принято в сообществе NN МТ, заключается в создании моделей с несколькими длинами окон. Теперь покажем, как это работает.

Рассмотрим снова строки 5–8 на рис. 5.2. С нашей текущей точки зрения, поразительно, что ни одна из двух команд TF, в первую очередь отвечающих за настройку кодера RNN, не упоминает размер окна. При создании GRU нужно знать размер RNN, поскольку, в конце концов, он отвечает за выделение пространства для переменных GRU. Но размер окна не входит в эту область. С другой стороны, `dynamic_rnn`, безусловно, должен знать размер окна, поскольку он отвечает за создание фрагментов графов TF, которые выполняют обратное распространение во времени. И он получает информацию, в данном случае — через переменную `embeds`, которая имеет размер `[bSz, wSz, embedSz]`. Итак,

предположим, что мы решили поддержать две разные комбинации размера окна. Первый, скажем, обрабатывает все предложения, в которых французский составляет 14 или меньше слов, а английский — 12 слов или меньше. Затем мы удваиваем их, 28 и 24. Если французское предложение содержит больше 28 слов или английское — больше 24, мы отбрасываем пример. Если французское или английское предложение больше 14 или 12 соответственно, но меньше пределов 28, 24, мы помещаем пару в большую группу. Затем мы создаем один GRU для использования в обоих dynamic_rnn следующим образом:

```
cell = tf.contrib.rnn.GRUCell(rnnSz)
encOutSmall, encStateS = tf.nn.dynamic_rnn(cell, smallerEmbs, ...)
encOutLarge, encStateL = tf.nn.dynamic_rnn(cell, largerEmbs, ...)
```

Обратите внимание, что, хотя у нас есть два (но возможно, пять или шесть) dynamic_rnn, в зависимости от диапазона размеров, которые мы хотим приспособить, все они используют одну и ту же ячейку GRU, поэтому они изучают и используют одни и те же знания французского. Подобным образом мы создали бы одну английскую ячейку GRU и т.д.

5.5. Упражнение по программированию

Эта глава посвящена технологии NN, используемой в МТ, поэтому здесь мы постараемся создать программу перевода. К сожалению, в нынешнем состоянии глубокого обучения это очень сложно. Хотя недавние успехи были впечатляющими, программы, которые могут похвастаться хорошими результатами, требуют около миллиарда учебных примеров и дней обучения, если не больше. Это не делает их хорошими обучающими примерами.

Вместо этого мы будем использовать около миллиона обучающих примеров — часть текста Canadian Hansard, ограниченную французско/английскими учебными примерами, в котором оба предложения состоят из 12 слов или менее (13, включая дополнительное слово “STOP”). Мы также установили гиперпараметры на малой стороне: размер векторного представления — 30, размер RNN — 64 и обучение на протяжении только одной эпохи. Мы устанавливаем скорость обучения 0,005.

Как отмечалось ранее, оценить программы машинного перевода сложно, за исключением того, что можно взять и оценить переводы самостоятельно. Мы принимаем особенно простую схему. Мы выполняем правильный английский перевод, пока не достигаем первого слова “STOP”. Созданное машиной слово считается правильным, если оно находится в той же позиции английского Hansard. Для повторения мы прекратим подсчет балов после первого слова “STOP”, например

the	law	is	very	clear	.	STOP
the	*UNK*	is	a	clear	.	STOP

будет насчитывать 5 правильных выводов из 7 слов. В конце мы делим общее количество правильных слов на сумму всех слов в наборе английских предложений.

С этим показателем наша авторская реализация набрала 59% правильности на тестовом наборе после одной эпохи (65% — после второй и 67% — после трех). Хорошо это или плохо, зависит от ваших предыдущих ожиданий. Учитывая наши предыдущие комментарии о необходимости миллиарда обучающих примеров, возможно, ваши ожидания были низкими; наши, конечно, были. Однако проверка выполненных переводов показывает, что даже 59% вводят в заблуждение своей оптимистичностью. Мы запустили программу, выводя первое предложение в каждой 400-й партии, используя размер партии 32. Вот первые два обучающих примера правильно переведенных входных данных:

Epoch 0 Batch 6401 Cor: 0.432627

* * *
* * *
* * *

Epoch 0 Batch 6801 Cor: 0.438996

le très hon. jean chrétien
right hon. jean chrétien
right hon. jean chrétien

Здесь “* * *” вставляется между сессиями парламента редактором. 14 410 строк файла из 351 846 строк состоят исключительно из этой маркировки. К этому моменту в первой эпохе (она уже на полпути) программа, без сомнения, запомнила соответствующий “английский” (который, конечно, идентичен). В том же духе имена следующего оратора всегда добавляются в Hansard перед тем, как он выступит. Жан Кретьен был премьер-министром Канады во время написания этого тома Hansard, и он, кажется, выступал 64 раза. Таким образом, перевод этого французского предложения также был запомнен. Действительно, можно спросить, не запомнен ли какой-либо из правильных переводов. Ответ — “да”, но это не так много. Вот список последних шести из 22 000 примеров тестовых наборов.

19154 the problem is very serious .

21191 hon. george s. baker :

21404 mr. bernard bigras (rosemont , bq) moved :
 21437 mr. bernard bigras (rosemont , bq) moved :
 21741 he is correct .
 21744 we will support the bill .

Это из серии с двойным вниманием к соответствующим словам, размером RNN 64, скоростью обучения 0,005 и одной эпохой. Описанная ранее метрика точности составила 68,6% для тестового набора. Мы вывели все тестовые примеры, которые были полностью правильны и не соответствовали ни одному английскому обучающему предложению.

Интересно и полезно получить представление о том, как меняется состояние между словами предложения. В частности, первая модель seq2seq использовала конечное состояние кодера для простейшего английского декодера. Поскольку мы только что приняли состояние в 13 слов, независимо от длины исходного французского предложения (максимум 12 слов), мы предполагаем, что не сильно потеряем, получив состояние после, скажем, восьми “STOP”, если оригинал на французском языке был в пять слов. Чтобы проверить это, мы рассмотрели 13 состояний, создаваемых кодером, и для каждого состояния вычислили косинусное сходство между последовательными состояниями. Ниже приводится учебное предложение, обрабатываемое на третьей эпохе.

Английский: that has already been dealt with.

Перевод: it is a . a .

Индексы французских слов: [18, 528, 65, 6476, 41, 0, 0, 0, 0, 0, 0, 0, 0]

Сходство состояний: .078 .57 .77 .70 .90 1 1 1 1 1 1 1 1 1 1

Сначала вы можете заметить ужасное качество “перевода” (два правильных слова из 8, “.” и “STOP”). Сходства состояний (state similarity), однако, выглядят вполне разумными. В частности, как только мы дойдем до конца предложения в слове f5 (на французском), все сходства состояний будут равны 1,0 — поэтому состояние, как мы и надеялись, совсем не меняется из-за дополнения.

Наименее похожие состояния — это первое по сравнению со вторым. Отсюда сходство увеличивается почти монотонно. Или, другими словами, по мере того как мы продвигаемся по предложению, появляется все больше информации о прошлом, которую стоит сохранить, так что вокруг остается больше старых состояний, делая следующее состояние похожим на текущее.

5.6. Упражнения

Упражнение 5.1. Предположим, что мы используем seq2seq с несколькими длинами окон для программы МТ и определились с двумя размерами предложений, один — до 7 слов (и “STOP”) для английского и 10 — для французского, а другой — до 10 слов для английского и 13 — для французского. Запишите входные данные, если французское предложение — “A B C D E F”, а английское — “M N O P Q R S T”.

Упражнение 5.2. Мы решили проиллюстрировать внимание в разделе 5.3 на особенно простой форме, в которой решение о внимании основывалось только на положении внимания во французском и английском предложениях. Более сложная версия основывается на выборе векторов входного состояния для английской позиции, над которой мы работаем, и предлагаемого вектора состояния, влияние которого мы определяем. Поскольку это может допускать более сложные суждения, требуется значительное усложнение модели. В частности, мы больше не можем использовать стандартное для декодера рекуррентной сети TF обратное распространение во времени. Объясните, почему.

Упражнение 5.3. Нередко отмечалось, что передача исходного языка в кодировщик seq2seq при *обратном* проходе (но оставление декодера работающим при прямом) улучшает производительность МТ незначительно, но постоянно. Придумайте правдоподобную историю, почему это может быть так.

Упражнение 5.4. В принципе, у нас может быть модель seq2seq с двумя потерями, которые мы суммируем, чтобы получить общую потерю. Одними из них могут быть текущие потери МТ, возникшие в результате неправильного прогноза следующего слова назначения с вероятностью 1. Вторыми могут быть потери в кодере после задействования кодировщика для прогнозирования следующего исходного (например, французского) слова, т.е. после потери языковой модели. а) Придумайте правдоподобную историю о том, почему это ухудшит производительность. б) Придумайте правдоподобную историю о том, почему это улучшит производительность.

5.7. Ссылки и что читать дальше

В 1980-х годах группа исследователей из IBM во главе с Фредом Йелинеком (Fred Jelinek) начала работу над проектом по созданию программы машинного перевода, научив машину переводить, замечая статистические закономерности. “Замечание” пришло из байесовского машинного обучения, а данные, как и в

этой главе, — из корпуса Canadian Hansard [23]. Этот подход после нескольких лет насмешек стал доминирующим и оставался таковым до недавнего времени. Теперь подходы к глубокому обучению быстро набирают популярность, и это только вопрос времени, когда все коммерческие системы МТ будут основаны на NN, если они еще этого не сделали. Ранним примером подхода NN является метод Кальчбреннера и Блансома [24].

Выравнивание в моделях seq2seq было представлено у Дмитрия Богданова и др. [25]. Эта группа также, по-видимому, первой применила для этого подхода стандартный термин *нейронный машинный перевод* (neural machine translation). В модели внимания только по позиции этой модели придаются только числовые значения для положений французских и английских слов при определении того, какой вес придать соответствующему французскому состоянию. В [25] модель также содержит информацию о состояниях LSTM для положений во французском и английском предложениях.

Что касается сетевых учебников по МТ, в сети появился таковой от Тада Луонга (Thad Luong) и других, причем только что, когда эта книга готовилась к изданию [26]. Предыдущий учебник Google seq2seq/МТ был не очень хорош для педагогических целей (представьте себе поваренную книгу, в которой рассказывается, как приготовить блины, словами “смешайте 100 000 галлонов (378541,18 литра) молока и миллион яиц...”), но этот выглядит довольно разумно и вполне может стать хорошей основой для дальнейшего изучения нейронных МТ в частности и seq2seq в целом.

Помимо МТ, есть много других задач, для которых предназначены модели seq2seq. Одна из самых “горячих” сейчас — это чат-боты (chatbot), программы, которые получают разговорные выражения и пытаются продолжить разговор. Они также являются одной из основ домашних помощников, таких как Alexa от Amazon. (“Горячие” — это, определенно, правильное слово здесь: есть сетевой журнал и статья по чат-ботам “Why Chatbots Are the Future of Marketing” (Почему чат-боты — это будущее маркетинга.) Возможный проект на эту тему описан в публикации Сурлиадипана Рама (Surlyadeepan Ram) [27].

Глава 6

Глубокое обучение с подкреплением

Обучение с подкреплением (Reinforcement Learning — RL) — это ветвь машинного обучения, связанная с изучением поведения *агента* (agent) в *среде* (environment) для максимизации *вознаграждения* (reward). Естественно, глубокое (deep) обучение с подкреплением ограничивает метод глубокого обучения.

Обычно среда определяется математически как *Марковский процесс принятия решений* (Markov decision process — MDP). MDP состоят из набора состояний $s \in S$, в котором может находиться агент (например, местоположений на карте), конечного набора действий ($a \in A$), функции $T(s, a, s') = \Pr(S_{t+1} = s' | S_t = s, A = a)$, которая переводит агента из одного состояния в другое, функции вознаграждения от состояния, действия и последующего состояния в действительные $R(s, a, s')$ и *дисконта* (discount) $\gamma \in [0, 1]$ (рассматривается далее). В целом действия являются вероятностными, поэтому T определяет распределение по возможному результирующему состоянию от выполнения действия в определенном состоянии. Модели называются Марковскими процессами принятия решений потому, что они делают *Марковское предположение* (Markov assumption, или Марковское свойство): если мы знаем текущее состояние, то история (как мы дошли до текущего состояния) не имеет значения.

В MDP время дискретно. В любой момент времени агент находится в каком-либо состоянии, предпринимает действие, которое переводит его в новое состояние, и получает некоторое вознаграждение, зачастую нулевое. Задача заключается в том, чтобы максимизировать *дисконтированное будущее вознаграждение* (discounted future reward), как определено

$$\sum_{t=0}^{t=\infty} \gamma^t R(s_t, a_t, s_{t+1}) \quad (6.1)$$

Если $\gamma < 1$, то эта сумма конечна. Если γ отсутствует (или, что эквивалентно, равно 1), то сумма может расти до бесконечности, что усложняет математику. Как правило, γ равно 0,9. Значение в уравнении 6.1 называется *дисконтированным* (discounted) будущим вознаграждением, поскольку повторяющее умножение на значение, меньшее единицы, заставляет модель “дисконтировать”

(т.е. снижать цену) вознаграждение в будущем по сравнению с вознаграждениями, которые мы получаем прямо сейчас. Это вполне резонно, поскольку никто не живет вечно.

Наша задача — *решить* (solve) вопрос MDP (мы также говорим о поиске оптимальной политики (policy)). Политика — это функция $\pi(s) = a$, которая для каждого состояния s определяет действие a , которое должен выполнить агент. Политика является оптимальной и обозначается как $\pi^*(s)$, если указанные действия приводят к максимальному ожидаемому дисконтированному будущему вознаграждению. *Ожидаемое* (expected) здесь означает поиск ожидаемого значения, как описано в разделе 2.4.3. Поскольку действия не являются детерминированными, одно и то же действие может в конечном итоге дать совершенно разные вознаграждения.

Итак, эта глава посвящена оптимальным политикам MDP: сначала с использованием так называемых *табличных* (tabular) методов, а затем — с использованием их партнеров глубокого обучения.

6.1. Итерация по значениям

Основной вопрос, на который мы должны ответить, прежде чем говорить о решении MDP, — предполагаем ли мы, что агент “знает” функции T и R , или он должен блуждать по среде, изучая их, а также создавая свою политику. Если мы знаем T и R , все значительно упрощается, поэтому мы начнем с этого случая. В этом разделе мы также предполагаем, что существует только конечное число состояний s .

Итерация по значениям (value iteration) примерно так же проста, как и изучение политики в MDP. (На самом деле это, возможно, не алгоритм обучения вообще, так как ему не нужно получать обучающие примеры или взаимодействовать со средой.) Алгоритм показан на рис. 6.1. V — это функция значения (value function) вектора размером $|s|$, в которой каждый элемент $V(s)$ является наилучшим ожидаемым дисконтированным вознаграждением, на которое мы можем надеяться, когда начинаем в состоянии s . Q (просто называемая *Q-функцией* (Q function)) представляет собой таблицу размером $|s| \times |a|$, в которой мы храним текущую оценку дисконтированного вознаграждения при выполнении действия a в состоянии s . Функция значения V имеет значение действительного числа для каждого состояния: чем больше число, тем легче достичь этого состояния. Q -функция более детализирована: она дает значения, которые мы можем ожидать для каждой пары “состояние–действие”. Если наши значения в V верны, то в строке 2(a)i $Q(s, a)$ установлено правильно. Это говорит о том, что значение $Q(s, a)$ состоит из немедленного вознаграждения $R(s, a, s')$ плюс значения для состояния, в котором мы заканчиваем, как указано V . Поскольку

действия не детерминированы, мы должны суммировать по всем возможным состояниям. Это дает нам ожидание.

1. Для всех s установить $V(s) = 0$

2. Повторять до сходжения:

(а) Для всех s :

- i. Для всех a установить $Q(s, a) = \sum_{s'} T(s, a, s')(R(s, a, s') + \gamma V(s'))$
- ii. $V(s) = \max_a Q(s, a)$

3. Возвратить Q

Рис. 6.1. Алгоритм итерации по значениям

Получив правильную Q , можем определить оптимальную политику p , всегда выбирая действие $a = \arg \max_{a'} Q(s, a')$. Здесь $\arg \max_x g(x)$ возвращает значение x , для которого $g(x)$ является максимальным.

Для большей конкретности мы рассмотрим очень простую MDP — задачу замерзшего озера (frozen-lake problem). Эта игра является одной из многих в составе *Open AI Gym* — группы компьютерных игр с унифицированными API, удобными для экспериментов по обучению с подкреплением. У нас есть сетка $4 * 4$ (озеро), показанная на рис. 6.2. Цель игры — добраться от начальной позиции (состояние 0 вверху слева) к цели (внизу справа), не провалившись сквозь лунку во льду. Мы получаем вознаграждение 1 каждый раз, когда предпринимаем действие и оказываемся в состоянии цели. Все остальные тройки “состояние–действие–состояние” не имеют вознаграждения. Если мы окажемся в состоянии лунки (или в состоянии цели), игра заканчивается, и если мы начинаем играть снова, мы возвращаемся в начальное состояние. В противном случае мы идем влево (l), вниз (d), вправо (r) или вверх (u) (числа от нуля до трех соответственно) с некоторой вероятностью “поскользнуться” и не пойти в намеченном направлении. Фактически способ, которым игра Open AI Gym запрограммирована на действие, например вправо, приводит нас с равной вероятностью к любому из непосредственно соседних состояний, за исключением прямо противоположного (например, влево), так что лед очень скользкий. Если действие заставляет нас покинуть озеро, оно оставляет нас в том состоянии, с которого мы начали.

0:S	1:F	2:F	3:F
4:F	5:H	6:F	7:H
8:F	9:F	10:F	11:H
12:H	13:F	14:F	15:G

S	начальное положение
F	замороженный участок
H	лунка
G	конечное положение

Рис. 6.2. Задача замерзшего озера

Чтобы вычислить V и Q для замерзшего озера, мы повторно просматриваем все состояния s и пересчитываем $V(s)$. Рассмотрим состояние 1. Это требует вычисления $Q(1, a)$ для всех четырех действий и последующей установки $V(1)$ на максимальное из четырех значений Q . Хорошо, давайте начнем с вычисления того, что произойдет, если мы решим двигаться влево, $Q(1, l)$. Для этого нужно сложить все s' — все игровые состояния. В игре 16 состояний, но начиная с состояния 1, мы можем достичь только трех из них с ненулевой вероятностью, состояний 0, 5 и 1 (попытка подняться вверх будет блокирована границей озера, и, таким образом, движения не будет совсем). Итак, рассматривая только конечные состояния s' , которые имеют ненулевые значения $T(1, l, s')$, мы вычисляем следующую сумму:

$$Q(1, l) = 0,33 \cdot (0 + 0,9 \cdot 0) + 0,33 \cdot (0 + 0,9 \cdot 0) + 0,33 \cdot (0 + 0,9 \cdot 0) \quad (6.2)$$

$$= 0 + 0 + 0 \quad (6.3)$$

$$= 0 \quad (6.4)$$

Первое из слагаемых говорит о том, что при попытке двигаться влево с вероятностью 0,33 мы в конечном итоге окажемся в состоянии 0. За это мы получаем нулевое вознаграждение, и наше предполагаемое будущее вознаграждение равно $0,9 \times 0$. Это нулевое значение, как и в случае, если вместо того, чтобы идти налево, мы поскользнулись (и оказались в состоянии 5) или остались в состоянии 1. Поэтому $Q(1, l) = 0$. Поскольку значения V для трех состояний достижимы из состояния 1 и все равны 0, $Q(1, d)$ и $Q(1, u)$ также равны 0, и в строке 2(a)ii устанавливается $V(1) = 0$.

Фактически на первой итерации V продолжает оставаться нулевым, пока мы не доберемся до состояния 14, в котором мы наконец получим ненулевые значения для $Q(14, d)$, $Q(14, r)$ и $Q(14, u)$:

$$Q(14, d) = 0,33 \times (0 + 0,9 \times 0) + 0,33 \times (0 + 0,9 \times 0) + 0,33 \times (1 + 0,9 \times 0) = 0,33$$

$$Q(14, r) = 0,33 \times (0 + 0,9 \times 0) + 0,33 \times (0 + 0,9 \times 0) + 0,33 \times (1 + 0,9 \times 0) = 0,33$$

$$Q(14, u) = 0,33 \times (0 + 0,9 \times 0) + 0,33 \times (0 + 0,9 \times 0) + 0,33 \times (1 + 0,9 \times 0) = 0,33$$

$$\text{и } V(14) = 0,33.$$

В левой части рис. 6.3 демонстрируется таблица значений V после первой итерации. Итерация по значению — это один из нескольких алгоритмов, которые работают в направлении оптимальной политики, сохраняя таблицы лучших оценок значений функций. Отсюда и название *табличные методы* (tabular method).

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0,33	0

0	0	0	0
0	0	0	0
0	0	0,1	0
0	0,1	0,46	0

Рис. 6.3. Значения состояний после первой и второй итераций по значению

На второй итерации снова большинство значений остаются равными 0, но на этот раз состояния 10 и 13 также получают ненулевые элементы Q и V , поскольку из них мы можем перейти в состояние 14, и, как только что наблюдалось, теперь $V(14) = 0,33$. Значения V после второго цикла итерации по значению показаны в правой части рис. 6.3.

Другая точка зрения на итерации по значениям состоит в том, что каждое изменение в V (и Q) включает в себя точную информацию о том, что произойдет за одно движение в будущем (мы получаем вознаграждение R), но затем возвращаемся к изначально неточной информации, уже включенной в эти функции. В конце концов функции включают в себя все больше и больше информации о состояниях, которых мы еще не достигли.

6.2. Q-обучение

Значение итерации предполагает, что обучаемый имеет доступ ко всем деталям среды модели. Теперь рассмотрим противоположный случай — *обучение без модели* (model-free learning). Агент может исследовать окружающую среду, делая ход, и возвращает информацию о вознаграждении и следующем состоянии, но не знает фактических вероятностей движения или функции вознаграждения T, R .

Если предположить, что наша среда представляет собой Марковский процесс принятия решений, наиболее очевидный способ планирования в среде без моделей — это случайное блуждание внутри среды, сбор статистики по T, R и последующее создание политики на основе Q -таблицы, как описано в последнем разделе. На рис. 6.4 показаны основные моменты программы для этого. Стока 1 создает игру замершего озера. Для запуска игры (из начального состояния) мы вызываем `reset()`. Одиночный запуск игры замершего озера заканчивается, когда мы либо падаем в яму, либо достигаем целевого состояния. Таким образом, внешний цикл (строка 2) указывает, что мы собираемся запустить игру 1000 раз. Внутренний цикл (строка 4) говорит, что для любой игры мы ограничиваем игру 99 ходами. (На практике этого никогда не произойдет — мы падаем в яму или достигаем цели задолго до этого.) Стока 5 говорит, что

на каждом ходе мы сначала случайным образом генерируем следующее действие. (Существует четыре возможных действия: влево, вниз, вправо и вверх: числа от 0 до 3 соответственно.) Стока 6 является критическим шагом. Функция `step(act)` получает один аргумент (действие, которое необходимо выполнить) и возвращает четыре значения. Первое — это состояние, в котором действие покинуло игру (в FL это целое число от 0 до 15), а второе — это значение вознаграждения, которое мы получаем (в FL обычно — 0, иногда — 1). Третье состояние, обозначенное на рисунке `dn`, является индикатором “истинно/ложно”, завершен ли раунд игры (т.е. мы упали в яму или достигли цели). Последний аргумент — это информация об истинных вероятностях перехода, которую мы игнорируем, если осуществляем обучение без моделей.

```

0 import gym
1 game = gym.make('FrozenLake-v0')
2 for i in range(1000):
3     st = game.reset()
4     for stps in range(99):
5         act=np.random.randint(0, 4)
6         nst,rwd,dn,_=game.step(act)
7         # обновление T и R
8         if dn: break
9

```

Рис. 6.4. Сбор статистики для игры Open AI Gym

Если хорошо подумать, случайные блуждающие в игре являются далеко не лучшим средством сбора нашей статистики. В основном происходит то, что мы попадаем в яму, возвращаемся к началу и продолжаем собирать статистику о том, что происходит в состояниях вблизи начального состояния. Гораздо лучшая идея — учиться и блуждать одновременно и позволить обучению влиять на то, куда мы идем. Если мы на самом деле собираем полезную информацию в процессе, то в ходе этого мы продвигаемся все дальше и дальше в игре, узнавая, таким образом, все больше и больше разных состояний. В этом разделе мы делаем это, осуществляя выбор в соответствии с вероятностью ϵ , либо а) выбирая случайный ход, либо (с вероятностью $(1 - \epsilon)$) б) основывая свое решение на знаниях, которые были получены до сих пор. Если ϵ фиксировано, это называется *эпсилон-жадной стратегией* (*epsilon-greedy strategy*).

Также нередко наблюдается уменьшение ϵ с течением времени (*стратегия снижения эпсилона* (*epsilon-decreasing strategy*)). Один из простейших способов сделать это подразумевает наличие связанного гиперпараметра E и установку $\epsilon = E/i+E$, где i — это количество раз, когда мы играли в игру. (Таким образом, E — это количество игр, в которых мы переходим от преимущественно случайных ходов к наиболее изученным.) Как вы можете ожидать, то, как мы решаем исследовать или основывать свой выбор на нашем нынешнем

понимании игры, может оказать большое влияние на скорость нашего изучения игры, и у этого есть собственное название — *компромисс между исследованием и использованием* (exploration-exploitation tradeoff) (применяя игровые знания, мы, как говорят, *используем* знания, которые уже приобрели). Другой популярный способ объединения исследования и использования — это всегда применять значения, заданные Q -функцией, но нужно сначала превратить их в распределение вероятностей, а затем выбрать действие в соответствии с этим распределением, а не всегда выбирать действие с наибольшим значением. (Последний называется *жадным алгоритмом* (greedy algorithm).) Поэтому, если бы у нас было три действия и их Q -значения были [4, 1, 1], мы выбрали бы первые две трети и т.д.

Q-обучение (Q-learning) — это один из первых и самых популярных алгоритмов обучения без моделей, сочетающий исследование и использование. Основная идея не в том, чтобы изучать R и T , а в том, чтобы изучать Q - и V -таблицы непосредственно. Теперь на рис. 6.4 нам нужно изменить строки 5 (мы больше не действуем совершенно случайно) и 7, в которых мы модифицируем Q и V , а не R и T .

Мы уже объяснили, что делать в строке 5, поэтому переходим к строке 7. Уравнения обновления Q-обучения таковы

$$Q(s, a) = (1 - \alpha)Q(s, a) + \alpha(R(s, a, n) + \gamma V(n)) \quad (6.5)$$

$$V(s) = \max_{a'} Q(s, a'), \quad (6.6)$$

Здесь s — это состояние, которое мы занимали, a — это действие, которое мы предприняли, а a' — это состояние, которое мы занимаем сейчас, только что сделав ход в игре в строке 6 на рис. 6.4.

Новое значение $Q(s, a)$ представляет собой смесь, контролируемую а его старого значения и новой информации, поэтому a является своего рода скоростью обучения. Обычно a невелико. Чтобы было понятно, зачем это нужно, имеет смысл сравнить эти уравнения со строками 2(a)i и 2(a)ii из алгоритма итерации по значениям на рис. 6.1. Таким образом, поскольку алгоритму были заданы R и T , мы могли бы суммировать все возможные результаты предпринятого нами действия. В Q-обучении мы так не можем. Все, что у нас есть, — это последний результат нашего хода. Новая информация основана только на одном ходе в нашем исследовании окружающей среды. Предположим, что мы находимся в состоянии 14 на рис. 6.2, но без нашего ведома существует очень малая вероятность (0,0001), что, если мы уйдем из этого состояния, мы получим “вознаграждение”, равное -10. Скорее всего, этого не произойдет, но если это произойдет, то это очень серьезно ухудшит ситуацию. Мораль в том, что алгоритм не должен

уделять слишком много внимания одному ходу. В итерации по значению мы знаем как T , так и R , и между ними алгоритм учитывает как вероятность отрицательного вознаграждения, так и низкую вероятность такого случая.

6.3. Основы глубокого Q-обучения

Овладев табличным Q-обучением, мы теперь сможем понять *глубокое Q-обучение* (deep-Q learning). Как и в табличной версии, мы начнем со схемы на рис. 6.4. На этот раз наибольшим отличием является то, что мы представляем Q -функцию не в виде таблицы, а с использованием модели NN.

В главе 1 мы кратко упомянули, что машинное обучение можно охарактеризовать как задачу *аппроксимации функций* (function-approximation problem) — поиск функции, которая лучше всего соответствует некой целевой функции; например, целевая функция может отображать пиксели в одно из десяти целых чисел, в которых пиксели взяты из изображения соответствующей цифры. Нам дано значение функции для некоторых входных данных, а цель заключается в том, чтобы создать функцию, которая точно соответствует ее выводу для всех этих значений, и тем самым заполнить значения функции в тех местах, где нам ее значение не было дано. В случае глубокого Q-обучения аналогия с аппроксимацией функций вполне уместна — мы собираемся аппроксимировать нашу (неизвестную) Q -функцию, используя NN, блуждая по Марковскому процессу принятия решений, обучаясь по пути.

Мы должны подчеркнуть, что переход от табличных моделей к глубокому обучению *не* мотивирован примером “замерзшего озера”, который является как раз той задачей, для которой подходит табличное Q-обучение. Глубокое обучение необходимо, когда существует слишком много состояний, чтобы создать для них таблицу.

Одним из событий в возрождении NN стало создание единой модели NN, которая могла бы применить глубокое обучение ко многим играм Atari. Эта программа была создана компанией DeepMind в 2014 году, приобретенной Google. Компания DeepMind смогла создать единую программу для изучения множества различных игр, представив игры в виде пикселей в изображениях, которые генерируют игры. Каждая комбинация пикселей является состоянием. Я так сразу и не вспомню размер изображения, которое они использовали, но даже если бы он был таким же маленьким, как изображения 28×28 , которые мы использовали для Mnist, и каждый пиксель был либо включен, либо выключен, это было бы 2^{784} возможных комбинаций значений пикселей — такое в принципе количество состояний будет необходимо в Q -таблице. В любом случае это слишком много для табличной схемы. (Я все же подсмотрел: окно игр Atari имеет размер 210×160 RGB, а программа DeepMind уменьшила его

до $84 * 84$ черно-белых пикселей.) Мы вернемся позже, чтобы обсудить случаи, более сложные, чем замерзшее озеро.

Замена Q -таблицы функцией NN сводится к следующему: чтобы получить рекомендацию по перемещению, а не смотреть в Q -таблицу, мы фактически обращаемся к Q -таблице, передавая состояние в однослойную NN, как показано на рис. 6.5. Код TF для создания только параметров модели Q -функции приведен на рис. 6.6. Мы передаем текущее состояние (скалярное `inptSt`), которое превращаем в унитарный вектор `oneH`, который преобразуется одним слоем линейных блоков Q . Q имеет форму $16 * 4$, где 16 — это размер одного унитарного вектора состояний, а 4 — это количество возможных действий. Вывод `qVals` — это элементы $Q(s)$, а `outAct` — максимум элементов в Q -таблице, являющийся рекомендацией политики.

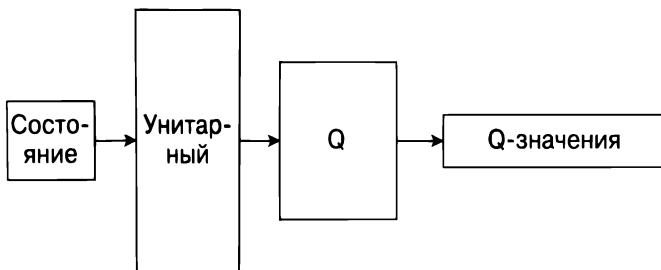


Рис. 6.5. NN глубокого Q-обучения для замерзшего озера

```

inptSt = tf.placeholder(dtype=tf.int32)
oneH=tf.one_hot(inptSt,16)
Q= tf.Variable(tf.random_uniform([16,4],0,0.01))
qVals= tf.matmul([oneH],Q)
outAct= tf.argmax(qVals,1)
  
```

Рис. 6.6. Параметры модели TF для функции Q-обучения

На рис. 6.6 подразумевается, что мы играем только по одной игре за раз, а потому, когда мы передаем исходное состояние (и получаем рекомендации политики), используется только одна из них. Исходя из нашей обычной практики NN, это соответствует размеру партии, равному единице. Например, исходное состояние, `inptSt`, является скаляром — номером состояния, в котором находится субъект. Из этого следует, что `oneH` является вектором. Тогда, поскольку `matmul` ожидает две матрицы, мы вызываем ее с `[oneH]`. Это, в свою очередь, означает, что `qVals` будет матрицей формы $[1, 4]$, т.е. будет иметь значения Q только для одного действия (вверх, вниз и т.д.). Наконец, `outAct` имеет форму $[1]$, поэтому рекомендуемое действие — `outAct[0]`. (Вы поймете, почему мы вдавались в эти подробности, когда мы рассмотрим остальную часть кода для глубокого обучения на рис. 6.7.)

Как и в табличном Q-обучении, алгоритм выбирает действие либо случайно (в начале процесса обучения), либо на основе рекомендации Q -таблицы (ближе к концу). При глубоком Q-обучении мы получаем рекомендацию Q -таблицы, вводя текущее состояние s в NN на рис. 6.5 и выбирая действие u, d, r или l . Когда у нас есть действие, мы вызываем `step`, чтобы получить результат, а затем учимся на нем. Естественно, чтобы сделать это в глубоком обучении, нам нужна функция потерь.

Но теперь, когда мы упомянули об этом, *какова функция потерь глубокого Q-обучения?* Это ключевой вопрос, поскольку, как было очевидно все время, когда мы делая ходы, особенно на ранних этапах обучения, мы не знаем, хороши они или плохи! Тем не менее мы знаем следующее: в среднем

$$R(s, a) + \gamma \max_{a'} Q(s', a') \quad (6.7)$$

(где, как и раньше, s' — это состояние, в котором мы заканчиваем ходы после a в s) — это более точная оценка $Q(s, a)$, чем текущее значение, потому что мы смотрим на один ход вперед. Итак, мы задаем потери

$$\left(Q(s, a) - \left(R(s, a) + \gamma \max_{a'} Q(s', a') \right) \right)^2, \quad (6.8)$$

квадрат разницы между тем, что только что произошло (когда мы сделали ход), и спрогнозированными значениями (из Q -таблицы/функции). Это *квадратичная функция потерь* (squared-error loss) или *квадратичная потеря* (quadratic loss). Разница между Q , рассчитанным сетью (первый член), и значением, которое мы можем вычислить, наблюдая фактическое вознаграждение за следующее действие плюс значение Q на один ход в будущем (второй член), является *ошибкой временных различий* (temporal difference error), или $TD(0)$. Если мы посмотрим на два хода в будущее, это будет $TD(1)$.

На рис. 6.7 приведен остаток кода TF (согласно приведенному на рис. 6.6). Первые пять строк создают оставшуюся часть графа TF. Теперь просмотрите оставшуюся часть кода с акцентом на строки 7, 11, 13, 14, 19 и 25. Они реализуют базовый AI Gym “wandering”, т.е. соответствуют всем строкам рис. 6.4. Мы создаем игру (строка 7) и играем 2000 отдельных игр (строка 11), каждая из которых начинается с `game.reset()` (строка 13). Каждый эпизод имеет максимум 99 ходов (строка 14). Фактический ход делается в строке 19. Игра завершается, как указано флагом `dn` (строка 25).

Остаются два пробела, строки 15–17 (выбор следующего действия) и 20–22. Стока 15 — это прямой проход, в котором мы присваиваем NN текущее состояние и возвращаем вектор длиной 1 (который следующая строка превращает в скаляр — номер действия). Мы также всегда даем программе небольшую

вероятность случайного действия (строка 18). Это гарантирует, что в конечном итоге мы исследуем все игровое пространство. Строки 20–22 относятся к вычислению потерь и выполнению обратного прохода для обновления параметров модели. Это также точка строк 1–5, которые создают граф TF для расчета и обновления потерь.

Производительность этой программы не так хороша, как у табличного Q-обучения, но, как мы уже говорили, табличные методы вполне подходят для MDP замершего озера.

```

1 nextQ = tf.placeholder(shape=[1, 4], dtype=tf.float32)
2 loss = tf.reduce_sum(tf.square(nextQ - qVals))
3 trainer = tf.train.GradientDescentOptimizer(learning_rate=0.1)
4 updateMod = trainer.minimize(loss)
5 init = tf.global_variables_initializer()
6 gamma = .99
7 game=gym.make('FrozenLake-v0')
8 rTot=0
9 with tf.Session() as sess:
10    sess.run(init)
11    for i in range(2000):
12        e = 50.0/(i + 50)
13        s=game.reset()
14        for j in range(99):
15            nActs,nxtQ=sess.run([outAct,qVals],feed_dict={inptSt: s})
16            nAct=nActs[0]
17            if np.random.rand(1)<e: nAct= game.action_space.sample()
18            s1,rwd,dn,_ = game.step(nAct)
19            Q1 = sess.run(qVals,feed_dict={inptSt: s1})
20            nxtQ[0,nAct] = rwd + gamma*(np.max(Q1))
21            sess.run(updateMod,feed_dict={inptSt:s, nextQ:nxtQ})
22            rTot+=rwd
23            if dn: break
24            s = s1
25    print "Percent games succesful: ", rTot/2000

```

Рис. 6.7. Остаток кода глубокого Q-обучения

6.4. Методы градиента политики

Теперь мы переходим к задаче Open AI Gym, которая не может быть решена стандартными табличными методами, а именно — к *тележке с шестом* (cart pole) и новому методу глубокого RL, *градиентам политики* (policy gradient). “Тележка с шестом” показана на рис. 6.8 и представляет собой тележку на одномерном пути. С помощью гибкого соединения к тележке прикреплен шест, и

когда тележка движется в том или ином направлении, верхушка шеста перемещается влево или вправо в соответствии с законами Ньютона. Состояния имеют четыре значения — позиция тележки и угол наклона шеста после предыдущего и текущего перемещений. Мы получаем значения в последовательные периоды времени, чтобы позволить программе выяснить направление движения. Игрок может выполнить два действия: переместить тележку вправо или влево. Импульс всегда имеет одинаковую величину. Если тележка переместится слишком далеко вправо или влево или вершина шеста отклонится слишком далеко от перпендикуляра, `step` дает сигнал о том, что текущая игра окончена, и необходимо будет перезапустить (`reset`) ее, чтобы начать новую. Мы получаем одну единицу вознаграждения за каждый ход, который делаем, прежде чем терпим неудачу. Естественно, цель в том, чтобы удерживать шест на тележке как можно дольше. Поскольку состояние соответствует четырем кортежам чисел, количество возможных состояний бесконечно, поэтому табличные методы исключены.

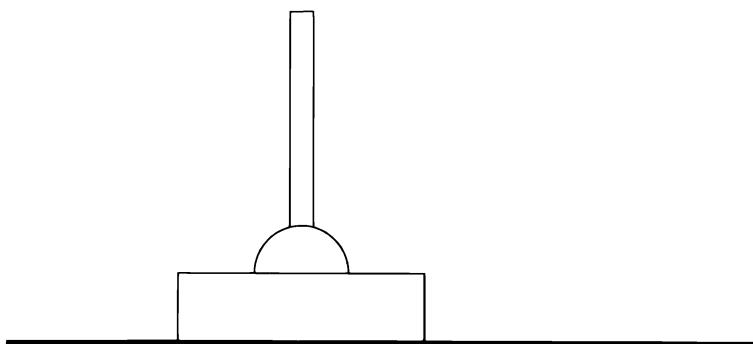


Рис. 6.8. Тележка с шестом

До сих пор мы использовали наши модели NN для аппроксимации Q -функции для нашего MDP. В этом разделе мы демонстрируем метод, в котором NN моделирует функцию политики непосредственно. Опять же, мы заинтересованы в обучении без моделей и снова принимаем парадигму блуждания по игровой среде, первоначально выбирая действия в основном случайным образом, но переходя к использованию рекомендаций NN. Как и везде в этой главе, главная проблема заключается в поиске подходящей функции потерь, так как мы не знаем, какие правильные действия мы должны предпринять.

В процессе глубокого Q-обучения мы делаем один ход за раз и зависим от того факта, что, сделав этот ход, получив вознаграждение и попав в новое состояние, наши знания о текущей локальной среде улучшились. Наши потери будут разницей между тем, что было спрогнозировано (например, Q -функцией) на основе прежних знаний, и тем, что фактически произошло.

Здесь мы попробуем нечто другое. Предположим, что мы играем всю итерацию текущей игры, не внося никаких изменений в нашу сеть, например мы делаем 20 ходов (направления движения тележки) до того, как шест упадет. На этот раз мы занимаемся исследованием/использованием, выбирая действия в соответствии с распределением вероятностей, полученным из Q -функции, вместо того чтобы брать максимум Q -функции.

В этом сценарии мы можем вычислить дисконтированное вознаграждение за первое состояние ($D_0(\mathbf{s}, \mathbf{a})$), когда за ним следуют все состояния и действия, которые мы только что опробовали:

$$D_0(\mathbf{s}, \mathbf{a}) = \sum_{t=0}^{n-1} \gamma^t R(s_t, a_t, s_{t+1}) \quad (6.9)$$

Если бы мы предприняли n ходов, мы могли бы вычислить будущее дисконтированное вознаграждение для любой комбинации “состояние–действие” s_i, a_i из рекуррентного отношения

$$D_n(\mathbf{s}, \mathbf{a}) = 0 \quad (6.10)$$

$$D_i(\mathbf{s}, \mathbf{a}) = R(s_i, a_i, s_{i+1}) + \gamma D_{i+1}(\mathbf{s}, \mathbf{a}) \quad (6.11)$$

Иначе говоря, будущее дисконтированное вознаграждение, например, за четвертое состояние в последовательности состояний, через которые мы переходим (при выполнении действия a), равно D_4 . Опять же, обратите внимание, что здесь мы получили информацию. Например, до того как мы попробовали первую случайную последовательность ходов, мы понятия не имели, что такое возможное вознаграждение. Теперь мы знаем, что, скажем, 10 возможно (и действительно разумно для случайной последовательности действий). Или, опять же, теперь мы знаем, уронили ли мы шест на 10-м ходу; тогда $Q(s_9, a_9) = 0$.

Вот хорошая функция потерь, которая фиксирует эти и многие другие факты:

$$L(\mathbf{s}, \mathbf{a}) = \sum_{t=0}^{n-1} D_t(\mathbf{s}, \mathbf{a})(-\log \Pr(a_t | s)) \quad (6.12)$$

Чтобы распаковать все это, сначала обратите внимание, что крайний справа член — это кросс-энтропийные потери, и он сам по себе побуждает сеть реагировать действиями a_t , когда она находится в состоянии s_t . Конечно, само по себе это довольно бесполезно, поскольку, особенно в начале обучения, мы выбирали действия случайным образом.

Далее рассмотрим влияние значений D_t . В частности, предположим, что a_0 была плохой реакцией на s_0 . Предположим, например, что тележка отцентрирована, шест с самого начала наклонен вправо и мы решаем пойти влево, заставив

шест наклониться еще дальше вправо. Читатель должен видеть, что при прочих равных значение D_0 в этом случае меньше, чем было бы, если бы мы выбрали правильный ход, — причина в том, что (при прочих равных), если первый ход хорош, шест и тележка должны оставаться в границах дольше (*и* больше), а значения D — больше. Таким образом, уравнение 6.12 дает более высокие потери плохому ходу a_0 , чем хорошему, обучая таким образом NN отдавать предпочтение хорошему ходу.

Эта комбинация “архитектура/функция потерь” известна как *REINFORCE*. На рис. 6.9 показана базовая архитектура. Важно отметить, что NN здесь используется двумя разными способами. В первую очередь, посмотрите на левую часть: мы передаем NN единственное состояние, которое, как упоминалось ранее, представляет собой четыре набора действительных значений, указывающих положение и скорость тележки, а также конца шеста. В этом режиме мы получаем вероятности для выполнения двух возможных действий, как показано в правом верхнем углу рисунка. В этом режиме мы не представляем заполнитель для вознаграждений или действий со значениями, поскольку а) мы их не знаем и б) они нам не нужны, потому что мы не вычисляем потери на данный момент. Сделав все ходы для всей игры, мы используем NN в другом режиме. На этот раз мы передаем ее последовательность действий и вознаграждений и просим ее вычислить потери и выполнить обратное распространение. В режиме обучения мы в некотором смысле вычисляем действия двумя разными способами. Во-первых, мы передаем NN состояния, которые проходим, и для каждого состояния слои вычисления политики вычисляют вероятности действий. Во-вторых, мы непосредственно поддерживаем действия, предпринятые в качестве заполнителя. Это связано с тем, что при выборе действий в режиме игры мы не обязательно выбираем действие

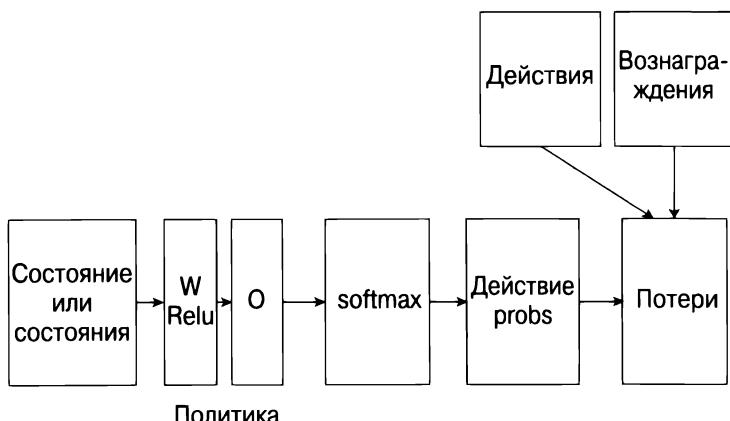


Рис. 6.9. Архитектура глубокого обучения для *REINFORCE*

с наибольшей вероятностью; мы выбираем его случайным образом на основе вероятностей действий. Чтобы вычислить потери в соответствии с уравнением 6.12, нам нужны оба.

На рис. 6.10 приведен код TF для создания градиента политики NN с использованием функции потерь в уравнении 6.12, а на рис 6.11 представлен псевдокод для использования NN при изучении политики и действий в игровой среде. Сначала рассмотрим псевдокод. Обратите внимание, что в самом внешнем цикле (строка 2) мы играем 3001 сеанс игры. Внутренний цикл (строка 2b) заставляет нас играть игровую сессию до тех пор, пока step не скажет, что игра закончена (строка D), или пока мы не сделаем 999 ходов. Мы выбираем случайное действие в соответствии с вероятностями, полученными из нашего NN (строки i, ii), а затем выполняем действие в игре. Мы сохраняем результаты в списке hist, чтобы у нас была запись о том, что произошло. Если действие приводит к окончательному состоянию, мы обновляем параметры модели.

```

state= tf.placeholder(shape=[None, 4], dtype=tf.float32)
W =tf.Variable(tf.random_uniform([4, 8], dtype=tf.float32))
hidden= tf.nn.relu(tf.matmul(state,W))
O= tf.Variable(tf.random_uniform([8,2], dtype=tf.float32))
output= tf.nn.softmax(tf.matmul(hidden,O))

rewards = tf.placeholder(shape=[None], dtype=tf.float32)
actions = tf.placeholder(shape=[None], dtype=tf.int32)
indices = tf.range(0, tf.shape(output)[0]) * 2 + actions
actProbs = tf.gather(tf.reshape(output, [-1]), indices)
aloss = -tf.reduce_mean(tf.log(actProbs) * rewards)
trainOp= tf.train.AdamOptimizer(.01).minimize(aloss)

```

Рис. 6.10. TF-инструкции графа для градиента политики NN случая шеста и тележки

Из рис. 6.10 мы видим, что output вычисляется из значений текущего состояния state, пропущенных через двухслойную NN с линейными блоками W и O, естественно разделенными tf.relu, а затем подается в softmax для преобразования логитов в вероятности. Как следует из предыдущих случаев применения многослойных NN, первый слой имеет размеры [input-size, hidden-size], а второй — [hidden-size, output-size], где hidden-size — это гиперпараметр (мы выбрали 8).

Поскольку здесь мы разработали новую функцию потерь и не использовали стандартную функцию из библиотеки TF, вычисление потерь должно было быть построено из более базовых функций TF (вторая половина рис. 6.10). Например, во всех наших предыдущих NN прямой и обратный проходы были неразрывно связаны, поскольку не было задействовано никаких вычислений извне TF. Здесь мы получаем значения reward извне: reward — это заполнитель,

который получается в соответствии со строками А, В и С на рис. 6.11. Точно так же `actions` является заполнителем.

1. `totRs=[]`

2. Для i в диапазоне (3001):

(a) `st=сбросить игру`

(b) для j в диапазоне(999):

i. `actDist = sess.run(output, feed dict=state:[st])`

ii. выбрать `act` случайно согласно `actDist`

iii. `st1,r,dn,_=game.step(act)`

iv. собрать `st,a,r` в `hist`

v. `st=st1`

vi. если `dn`:

A. `disRs = [Di(states, действия из hist) | i = 0 to j - 1]`

B. создать `feed_dict` с состоянием `state=st`, действия из `hist` и `rewards=disRs`.

C. `sess.run(trainOp,feed_dict=feed_dict)`

D. добавить j к концу `totRs`

E. прекратить

vii. если $i \% 100 = 0$: вывести среднее из последних 100 записей в `totRs`

Рис. 6.11. Псевдокод для NN “политика–градиент–обучение” случая шеста и тележки

$$\begin{array}{lll}
 \Pr(\mathbf{l} | s_1) & \Pr(\mathbf{r} | s_1) & \Pr(a_1 | s_1) \\
 \Pr(\mathbf{l} | s_2) & \Pr(\mathbf{r} | s_2) & \Pr(a_2 | s_2) \\
 & \rightarrow & \\
 \Pr(\mathbf{l} | s_n) & \Pr(\mathbf{r} | s_n) & \Pr(a_n | s_n)
 \end{array}$$

Рис. 6.12. Извлечение вероятностей действий из тензора всех вероятностей

В последних трех строках на рис. 6.10 все выглядит более знакомым. `loss` просто вычисляет величины из уравнения 6.12. Для `optimizer` мы использовали оптимизатор *Adam* (*Adam optimizer*). Мы могли бы использовать наш знакомый оптимизатор градиентного спуска, просто заменив его и удвоив скорость обучения, и мы достигли бы почти такой же хорошей производительности, но не совсем. Оптимизатор *Adam* немного сложнее и, как правило, считается лучшим. Он отличается от градиентного спуска несколькими способами, наиболее фундаментальным из которых является использование *импульса* (*momentum*).

Как следует из названия, оптимизатор, использующий импульс, имеет тенденцию продолжать перемещать значение параметра вверх/вниз, если оно перемещалось вверх/вниз в последнее время куда чаще, чем с градиентным спуском.

Остаются две средние строки на рис. 6.10, в которых устанавливаются `indices` и `actProbs`. Игнорируйте пока их работу и сосредоточьтесь на том, что им нужно делать. Для этого необходимо преобразование, показанное на рис. 6.12. Слева мы видим результат прямого прохода, вычисляющий вероятность того, что каждое из возможных действий `r` и `l` является наилучшим. Если бы это была глава 1 и у нас был бы полный контроль, мы умножили бы ее на тензор размером в одну партию для унитарных векторов, чтобы получить вероятности действий, которые мы должны предпринять в соответствии с контролем. На самом деле это то, что мы показываем справа на рис. 6.12.

Выполнение этого преобразования зависит от функции `gather`, получающей два аргумента,

```
tf.gather(tensor, indices)
```

и выдает элементы тензора, заданные числовыми индексами, а также объединяет их в новый тензор. Например, если `tensor` равен $((1,3), (4,6), (2,1), (3,3))$, а `indices` равен $(3,1,3)$, то на выходе получается $((3,3), (4,6), (3,3))$. В нашем случае мы превращаем матрицу вероятности действий слева на рис. 6.12 в вектор вероятностей и, в зависимости от предыдущей строки, устанавливаем `indices` в правильный список, поэтому `tf.gather` собирает вероятности только действий, указанных вектором `actions`. Доказательство того, что `indices` установлен правильно, оставим читателю в качестве упражнения (упражнение 6.5).

Имеет смысл вернуться и более внимательно посмотреть на то, как связаны Q-обучение и REINFORCE. В первую очередь, они различаются способом сбора информации об окружающей среде для информирования NN. Q-обучение продвигается на один шаг, а затем проверяет, близок ли прогноз NN к результату и тому, что фактически произошло. Оглядываясь на уравнение 6.8, функцию потерь Q-обучения, мы видим, что если прогноз и результат совпадают, то обновлять нечего. При REINFORCE, напротив, мы проигрываем весь эпизод перед тем, как изменить любые параметры NN, где эпизод (episode) — это полный ход игры, от начального состояния до тех пор, пока игра не сигнализирует о завершении. Обратите внимание, что мы могли бы сделать что-то вроде Q-обучения, но использовали расписание изменения параметров REINFORCE. Это замедляет обучение, поскольку мы вносим изменения в параметры гораздо реже, но в качестве компенсации мы вносим более качественные изменения, поскольку рассчитываем фактическое дисконтированное вознаграждение.

6.5. Методы актор-критик

Теперь, рассмотрев различия между Q-обучением и REINFORCE, мы сконцентрируемся на сходстве. В обоих случаях NN вычисляет либо политику, либо, при Q-обучении, функцию, которую можно просто использовать для создания политики (для любого состояния s всегда выполняйте действие a , которое максимизирует $Q(s, a)$). Таким образом, в обоих случаях наша NN аппроксимирует одну функцию, которая указывает нам, как действовать. Мы называем такие RL-программы *акторскими* (actor) методами. В этом разделе мы рассмотрим программы, которые имеют два подкомпоненты NN, каждый из которых имеет собственные функции потерь: один, как и прежде, является программой актора, а второй — программой *критика* (critic). Как вы можете догадаться, мы называем этот тип RL методом *актор-критик* (actor-critic method). В частности, в этом разделе мы рассмотрим улучшенный метод актор-критик (advantage actor-critic method или a2c). Это хороший выбор для нас, поскольку а) он работает достаточно хорошо и б) мы можем подходить к нему постепенно, начиная с REINFORCE. Мы называем первую версию (инкремент) a2c-. Мы снова применяем это к игре “Тележка с шестом”.

Это называется *улучшенным* (advantage) методом актора-критика потому, что в нем используется понятие “улучшения” (advantage). Улучшение пары “состояние–действие” заключается в разнице между значением Q состояния–действия и значением состояния:

$$A(s, a) = Q(s, a) - V(s) \quad (6.13)$$

Интуитивно мы ожидаем, что улучшение будет отрицательным числом, поскольку, скажем, в итерации по значениям $V(s)$ вычисляется в ходе обработки аргумента $\arg \max_a$ для возможных действий. Но для хороших действий A велико, поскольку идут отрицательные числа, поэтому A измеряет, насколько хорошим является действие в определенном состоянии по сравнению с состоянием в целом.

Затем мы определяем потери, которые a2c несет при изучении последовательности действий от начального состояния до конца игры, следующим образом:

$$L_A(s, a) = \sum_{t=0}^{n-1} A(s_t, a_t) (-\log \Pr(a_t | s_t)) \quad (6.14)$$

Это очень близко к потерям REINFORCE в уравнении 6.12, но мы заменили $D_t(s, a)$, дисконтированное вознаграждение, на $A_t(s, a)$. Мы назвали эту потерю L_A , чтобы отличать ее от общих потерь для a2c, которые, как мы увидим ниже, включают в себя вторую потерю L_C , связанную с критиком.

Мы помним, что потери REINFORCE предназначены для поощрения действий, которые приводят к большему вознаграждению. Теперь мы поощряем действия, которые лучше, чем альтернативные действия из того же состояния. Хотя это достаточно разумно, почему так должно быть лучше, чем поощрение действия с высоким вознаграждением действий напрямую?

Ответ связан с дисперсией $A(s, a)$. Как отмечено в разделе 2.4.3, дисперсия функции — это ожидание квадрата разности между значением функции и ее средним значением. Интуитивно понятно, что это означает, что сильно различающиеся функции имеют высокую дисперсию, и по сравнению с Q , A должна иметь гораздо меньшую дисперсию. Посмотрим на тележку. Предполагая, что игра дает нам разумную отдачу с точки зрения перемещения влево или вправо по сравнению с тем, как быстро движется шест, разница между движением вправо и влево будет небольшой, и, следовательно, A мало практически во всех частях пространства состояний. Сравните это с Q . Игра в тележку и шест после обучения на 100 играх в среднем дает около 20 ходов, прежде чем шест упадет, тогда как даже умеренно хорошая политика дает 200 или больше ходов.

Добавим теперь второй факт: при прочих равных легче аппроксимировать функцию с низкой дисперсией, чем функцию с высокой. Постоянная функция с нулевой дисперсией является самой простой из всех. Таким образом, если A намного проще оценить, это могло бы перевесить недостаток, связанный с максимизацией A , а не Q напрямую. Кажется, что так и есть. Конечно, мы не знаем, как вычислить A на этом этапе. Так что это следующий вопрос в нашей повестке дня.

Как вы помните, в REINFORCE мы следуем по пути, основанному на нашей текущей политике, до конца игры и используем дисконтированное вознаграждение $D_t(s.a)$ из уравнения 6.11 для оценки $Q(s, a)$. Теперь мы вычислим это как удвоенный долг в качестве оценки Q (уравнение 6.13). Что касается $V(s)$, мы встраиваем в нашу NN подсеть только для ее вычисления.

```

V1 =tf.Variable(tf.random_normal([4, 8], dtype=tf.float32, stddev=.1))
v1Out= tf.nn.relu(tf.matmul(state,V1))
V2 =tf.Variable(tf.random_normal([8, 1], dtype=tf.float32, stddev=.1))
v2Out= tf.matmul(v1Out,V2)
advantage = rewards-v2Out
aLoss = -tr.reduce_mean(tf.log(actProbs) * advantage)
cLoss= tf.reduce_mean(tf.square(rewards-vOut))
loss=aLoss + cLoss

```

Рис. 6.13. Код TF, добавленный к рис. 6.10 и 6.11 для a2c

На рис. 6.13 приведен дополнительный код построения сети TF, помимо того, который требуется для REINFORCE (рис. 6.10). Мы создали двухслойные

полносвязные NN v1Out и v2Out для вычисления V , функции значения “критика”. Она обучена для получения хороших оценок V с использованием квадратичных потерь на несоответствие между фактическим полученным вознаграждением и выводом аппроксимации NN ($cLoss$). Потери актора здесь следуют из уравнения 6.14, а потому используют функцию улучшения. Эти относительно небольшие изменения превращают наше REINFORCE в a2c–.

Фактический a2c, сверх a2c–, включает в себя два дополнительных улучшения. Одна проблема с REINFORCE (унаследованная a2c–) заключается в том, что ему нужно сыграть всю игру, прежде чем начнется какое-либо обучение. В начале, когда игра “Тележка с шестом” длится всего 10–20 ходов, это не является большим недостатком. Но игры REINFORCE заканчиваются одним или двумя сотнями ходов, а игры a2c– еще длиннее. A2c может улучшить положение, обновляя параметры модели гораздо раньше и чаще.

Трюк заключается в том, чтобы приостановить выполнение игры, скажем, через 50 (гиперпараметр) действий, чтобы обновить параметры модели. В REINFORCE мы не могли сделать это. В конце концов, смысл следовать целям всей игры заключался в том, чтобы получить хорошую оценку значений Q для действий, которые мы выполнили. Но a2c позволяет нам сделать оценку, просто сложив вместе а) фактические вознаграждения, которые мы накопили за последние 50 ходов, и б) значение V состояния, в котором мы оказались. Затем мы обнуляем значение $hist$ и перезапускаем его с нуля на 51-м ходу, только чтобы повторить все снова 50 ходов спустя. (В крайнем случае это также может ослабить требование REINFORCE о том, чтобы использовать его только в играх с явным перезапуском игры.)

Вторым улучшением в полной версии a2c является использование нескольких сред. Мы уже отмечали, что запуск серии обучающих примеров выгоден, поскольку позволяет лучше использовать способности быстрого умножения матриц. Проигрывание игры по одной за раз не позволяет этого при вычислении следующего игрового действия. Проигрывание нескольких игр эквивалентно партийным примерам в этом отношении.

6.6. Воспроизведение опыта

Ранее мы упоминали, что основным катализатором возрождения NN стал успех DeepMind с программой, которая может играть в несколько игр Atari на уровне экспертов. Используемая там технология NN известна как *DQN (Deep Q Network)*. Эта конкретная схема RL была в значительной степени заменена методами актор-критик, но программа представила также несколько улучшений, которые ортогональны использованию методов актор-критик. Одним из них является *воспроизведение опыта* (experience replay).

Как и следовало ожидать, RL является важной составляющей современного скачка в области беспилотных автомобилей. Одной из наибольших проблем в применении RL к этой области является получение обучающих данных. Текущему RL их требуется много, и по сравнению с компьютерами реальный мир, и в частности улицы и шоссе, движутся очень медленно. На самом деле, если вы начнете отсчитывать время в играх Open AI Gym, даже компьютерное моделирование может быть медленным — большая часть времени, проведенного в RL, тратится на ход игры. Если бы мы могли ускорить мир, мы могли бы учиться еще быстрее, но мы не можем.

При воспроизведении опыта мы используем одни и те же учебные данные несколько раз. Это проще всего объяснить в контексте игр Open AI Gym. Возвращаясь к REINFORCE, когда мы играли в игру, мы использовали переменную $hist$ для записи истории игры — каждое состояние, которое мы занимали, действие, которое мы предпринимали, состояние, в котором мы оказывались, и полученное вознаграждение. Это было необходимо в конце игры, чтобы вычислить D_s , но, вычислив их, мы отбрасывали историю. С воспроизведением опыта для каждого времени t мы сохраняем $\langle s_t, a_t, s_{t+1}, D_t \rangle$. С помощью этих чисел мы можем провести еще одну прямую и обратную передачи наших данных и получить из них больше “сока”. Есть и еще одно преимущество: мы можем играть, а затем воспроизводить каждый временной шаг в случайному порядке. Возможно, вы помните предположение iid, упомянутое в разделе 1.6, где мы отметили, что RL может быть особенно проблематичным, так как примеры обучения были соотнесены с самого начала. Выполнение случайных действий из нескольких различных игровых процессов значительно снижает эту проблему.

Конечно, за это мы платим цену. Старый учебный пример не так информативен, как новый. Кроме того, данные могут устареть. Предположим, у нас есть данные с самого начала нашего обучения, когда мы еще не умели, скажем, двигаться влево, когда шест наклоняется далеко вправо. И предположим, с тех пор мы научились действовать лучше. Это означает, что мы бесполезно переучиваемся со старых данных, что делать в состоянии s_{old} когда фактически наша текущая политика никогда не позволяет нам достигнуть этого состояния. Таким образом, вместо этого мы делаем нечто вроде следующего: сохраняем буфер из 50 игр, соответствующий, скажем, 5000 наборам “состояние–действие–состояние–вознаграждение” в четырех кортежах (перед неудачей мы усредняем 100 ходов). Теперь мы выбираем, например, 400 состояний наугад для обучения. Затем мы заменяем самую старую игру в буфере новой игрой, сыгранной с использованием новой политики, основанной на современных параметрах.

6.7. Ссылки и что читать дальше

Обучение с подкреплением имело богатую теорию и практику задолго до появления глубокого обучения, и глубокое обучение его не вытеснило. В конце концов, основная проблема в RL — как учиться, когда есть только косвенная информация о том, какие ходы хороши, а какие плохи, и глубокое обучение только сводит этот вопрос к вопросу о том, как определить функцию потерь. О природе решения не говорится ничего. Классический труд по RL принадлежит Ричарду Саттону и Эндрю Барто [28]. Я сам многое узнал из первого учебника Кельблинг и других [29], включая материал, предшествовавший глубокому обучению.

Что после глубокого обучения? Вначале своего изучения области обучения с подкреплением я наткнулся на блог Артура Джулиани на эту тему. Если вы заедете на его блог, особенно в частях 0 [30] и 2 [31], вы заметите, что мои презентации на темы тележки с шестом и REINFORCE значительно повлияли на его код, и его код стал отправной точкой для моего. Это напоминает мне оригинальную статью о REINFORCE, написанную Рональдом Вильямсом [41].

Алгоритм обучения с подкреплением a2c был предложен как вариант *асинхронного улучшенного алгоритма актор-критик* (Asynchronous Advantage Actor-Critic — a3c). На с. 87 мы отметили, что a2c позволяет некоторым средам лучше использовать программное и аппаратное обеспечение умножения матриц. В a3c эти среды обрабатываются асинхронно, по-видимому, для лучшего смешения комбинаций состояния и действия, которые ученик может наблюдать [32]. В этой же статье предложен алгоритм a2c в качестве подкомпоненты a3c. В конце концов, было показано, что он работает так же хорошо, и он намного проще.

6.8. Упражнения

Упражнение 6.1. Докажите, что V -таблица, показанная в правой части рис. 6.3, дает правильные значения (до двух значащих цифр) для значений состояния после второго прохода итерации по значениям.

Упражнение 6.2. Уравнение 6.5 имеет параметр a , но наша реализация TF на рис. 6.6 и 6.7, по-видимому, не упоминает a . Объясните, где он “прячется” и какое значение мы ему присвоили.

Упражнение 6.3. Предположим, что на фазе обучения алгоритма REINFORCE тележки с шестом потребовалось всего три действия (l,l,r), чтобы достичь конца, и $\Pr(l \mid s_1) = 0,2$; $\Pr(l \mid s_2) = 0,3$; $\Pr(r \mid s_3) = 0,9$. Покажите значения output, actions, indices и actProbs.

Упражнение 6.4. В REINFORCE мы сначала выбираем действия, которые ведут нас от начала игры “Тележка с шестом” до тех пор, пока (как правило) шест не упадет или тележка не выйдет за пределы окна. Мы делаем это без обновления параметров. Мы сохраняем эти действия и, по сути, снова и снова проходим весь сценарий, на этот раз вычисляя потери и обновляя параметры. Обратите внимание, что если бы мы сохранили действия *и их вероятности softmax*, мы могли бы вычислить потери, не выполняя все вычисления, которые передаются в функцию потерь во второй раз. Объясните, почему это, тем не менее, не работает — почему REINFORCE ничего не изучил бы, если бы мы делали это без дублированных вычислений.

Упражнение 6.5. Функция TF `tf.range` при наличии двух аргументов

```
tf.range(start, limit)
```

создает вектор целых чисел начиная со `start` и до (но не включая) `limit`. Если не указана именованная переменная `delta`, целые числа различаются на единицу. Таким образом, ее использование на рис. 6.10 создает список целых чисел в диапазоне от 0 до размера партии. Объясните, как они сочетаются со следующей строкой TF для выполнения преобразования на рис. 6.12.

Глава 7

Модели нейронных сетей, обучаемых без учителя

Эта книга пошла по большей части по непроверенному пути от задач *обучения с учителем* (supervised learning), таких как Mnist, к *обучению со слабым привлечением учителя* (weakly supervised learning), такому как обучение seq2seq и обучение с подкреплением. Говорят, что наша задача распознавания цифр относится к *обучению с полным привлечением учителя* (fully supervised learning), потому что каждый учебный пример сопровождается правильным ответом. В наших примерах для обучения с подкреплением учебные примеры не имеют маркировки. Вместо этого мы получаем форму слабой маркировки, поскольку процесс обучения направляет вознаграждения, которые мы получаем от Open AI Gym. В этой главе мы рассмотрим *обучение без учителя* (unsupervised learning), когда у нас нет меток или других форм контроля. Мы хотим узнать структуру наших данных только из самих данных. В частности, мы рассмотрим *автокодировщики* (autoencoder — AE) и *генеративно-состязательные сети* (Generative Adversarial Network — GAN).

7.1. Основы автокодировки

Автокодировщик — это функция, вывод которой при правильной работе практически идентичен вводу. Для нас эта функция — нейронная сеть. Чтобы сделать все нетривиальным, мы расположим на пути препятствия; наиболее распространенным методом является *снижение размерности* (dimensionality reduction). На рис. 7.1 показан простой двухслойный AE. Ввод (скажем, изображение Mnist размером $28 * 28$ пикселей) пропускается через слой линейных блоков и преобразуется в промежуточный вектор, который значительно меньше исходного ввода, например 256 по сравнению с исходным 784. Затем этот вектор сам проходит через второй слой, и цель заключается в том, чтобы выходные данные второго слоя были идентичны входным данным первого слоя. Рассуждения основаны на том, до какой степени мы можем уменьшить размеры среднего уровня по сравнению с входом, NN (нейронная сеть) закодировала

информацию о структуре изображений Mnist в среднем слое. Чтобы поднять это на более высокий уровень абстракции, мы имеем:

ввод → кодер → скрытый → декодер → ввод

Здесь *кодер* (encoder) выглядит, как ориентированная на задачу NN, а *декодер* (decoder) — как кодер наоборот. Процесс кодирования также известен как *субдискретизация* (downsampling) (поскольку он уменьшает размер изображения), а декодирования — как *передискретизация* (upsampling).

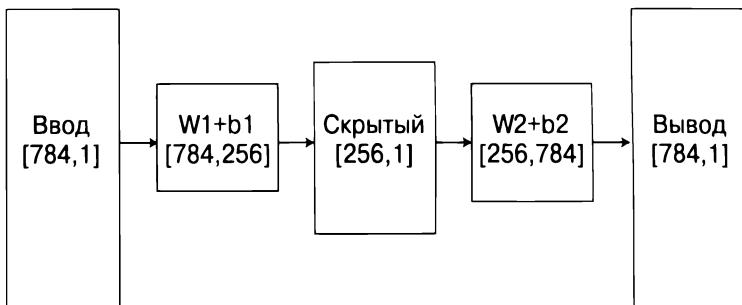


Рис. 7.1. Простой двухслойный AE

Мы заботимся об AE по нескольким причинам. Одна из них просто теоретическая и, возможно, психологическая. За исключением обучения в школе, люди мало учатся с учителем, и к тому времени, когда мы поступаем в школу, мы уже усвоили самые важные из наших навыков: наблюдение, разговорный язык, двигательные навыки и, в самом общем смысле, планирование. Предположительно, это возможно только через обучение без учителя.

Более практической причиной является использование *предобучения* (pre-training) или *сообучения* (co-training). Помеченных учебных данных обычно не достаточно, и наши модели почти всегда работают тем лучше, чем больше параметров им приходится обрабатывать. Предобучение — это метод изучения некоторых параметров для соответствующей задачи заранее, а затем начинается основной цикл обучения не с нашей стандартной случайной инициализации, а скорее со значений, достигнутых при обучении соответствующей задаче. Сообучение работает практически так же, за исключением того, что мы проводим обучение и “предобучение” одновременно, т.е. модель NN имеет две функции потерь, которые суммируются для получения полных потерь. Одним из них являются “действительные” потери, сводящие к минимуму количество ошибок в задаче, которую мы действительно хотим решить. Другим — потери для связанный задачи, которая может просто воспроизводить некоторые или все данные, которые мы используем, — это задача автокодирования.

Третья причина для изучения автокодирования — это вариант, называемый *вариационным автокодером* (*variational autoencoder*). Он похож на стандартный, за исключением того, что он предназначен для возврата случайных изображений в стиле тех, на которых обучался АЕ. Разработчик видеоигр может заставить действие игры происходить в городе, но не будет тратить время на отдельное проектирование, скажем, нескольких сотен зданий, в которых будет происходить действие. В наши дни эта задача может быть доверена хорошим вариационным автокодировщикам. В долгосрочной перспективе мы можем надеяться на гораздо лучшие автокодировщики, которые могли бы писать новые романы, которые читаются, как Хемингуэй или ваш любимый детектив.

	7	8	9	10	11	12	13	14	15	16	17	18	19	20
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	3	14	15	1	0	0	0	0	0	0	0	0	0	0
7	187	221	205	151	74	11	1	0	0	2	8	23	55	69
8	237	249	251	250	249	239	221	225	197	214	216	236	237	228
9	92	194	232	219	217	225	245	251	251	249	241	237	249	250
10	1	8	7	17	31	49	100	126	106	45	37	81	242	251
11	0	0	0	0	1	9	13	7	2	0	1	43	239	247
12	0	0	0	0	0	2	2	0	0	0	2	151	247	215
13	0	0	0	0	0	0	0	0	0	1	61	246	248	57
14	0	0	0	0	0	0	0	0	1	32	207	253	185	10
15	0	0	0	0	0	0	0	0	9	176	251	237	31	1
16	0	0	0	0	0	0	0	0	47	237	252	67	2	0
17	0	0	0	0	1	0	1	9	171	249	237	9	0	0
18	0	0	2	7	1	1	5	100	243	251	138	1	0	0
19	0	0	0	2	1	0	19	217	253	222	19	0	0	0
20	0	0	0	0	0	2	107	246	241	44	1	0	0	0
21	0	0	0	0	1	48	220	247	168	4	0	0	0	0
22	0	0	0	0	18	196	251	233	42	0	0	0	0	0
23	0	0	0	1	98	249	250	140	2	0	0	0	0	0
24	0	0	0	14	237	254	242	40	0	0	0	0	0	0
25	0	0	5	116	252	254	205	8	0	0	0	0	0	0
26	0	0	16	158	253	249	56	4	2	0	1	0	0	0
27	0	0	0	0	2	1	0	0	0	0	0	0	0	0

Рис. 7.2. Реконструкция тестового примера Mnist на рис. 1.1

Мы начнем с базового автокодирования, в ходе которого мы просто восстанавливаем ввод — цифры Mnist. На рис. 7.2 показан вывод четырехслойного АЕ, когда на входе находится изображение семерки в самом начале главы 1. Автокодер может быть выражен как

$$\mathbf{h} = S(S(\mathbf{x}\mathbf{E}_1 + \mathbf{e}_1)\mathbf{E}_2 + \mathbf{e}_2) \quad (7.1)$$

$$\mathbf{o} = S(\mathbf{h}\mathbf{D}_1 + \mathbf{d}_1)\mathbf{D}_2 + \mathbf{d}_2 \quad (7.2)$$

$$L = \sum_{i=1}^n (x_i - o_i)^2 \quad (7.3)$$

У нас есть два полносвязных слоя кодирования: первый — с весами \mathbf{E}_1 , \mathbf{e}_1 , второй — с \mathbf{E}_2 , \mathbf{e}_2 . При создании реконструкции цифры “7” на рис. 7.2 первый слой имеет форму [784, 256], а второй — [256, 128], поэтому конечное изображение имеет размер 128 “пикселей”, т.е. высота и ширина изображения равны $\sqrt{128}$. Уравнение 7.3 утверждает, что мы используем квадратичную функцию потерь, а это имеет смысл потому, что мы пытаемся прогнозировать значения пикселей, а не членство в классе. Для нашей нелинейности мы используем S , сигмовидную функцию.

Вы можете удивиться, почему мы использовали сигмовидную функцию активации, а не нашу почти стандартную `relu`. Причина в том, что до сих пор мы не особо интересовались фактическими значениями, которые передаются через сеть: в конце все они передаются через функцию `softmax`, и в основном все, что остается, — это их относительные значения. АЕ, напротив, сравнивает абсолютные значения входных данных со значениями выходных данных. Как вы, наверное, помните, обсуждая нормализацию данных наших изображений Mnist (19), мы разделили необработанные значения пикселей на 255, чтобы нормализовать их значения от 0 до 1. Как мы видели на рис. 2.7, сигмовидная функция колеблется от 0 внизу до 1 вверху. Поскольку это точно соответствует диапазону значений пикселей, NN не нужно учиться помещать значения в этот диапазон — диапазон скорее “встроен”. Естественно, это существенно упрощает получение данных значений при обучении, чем с функцией `relu`, которая имеет нижнее ограничение, но не верхнее.

Поскольку автокодировщики имеют несколько довольно похожих слоев (в этом случае все они полносвязны), модуль `layers`, обсуждаемый в разделе 2.4.4, здесь особенно полезен. В частности, обратите внимание, что модель, описанная в уравнениях 7.1–7.3, может быть кратко закодирована как

```
E1=layers.fully_connected(img,256,tf.sigmoid)
E2=layers.fully_connected(E1,128,tf.sigmoid)
```

```
D2=layers.fully_connected(E2, 256, tf.sigmoid)
D1=layers.fully_connected(D2, 784, tf.softmax)
```

Здесь мы предполагаем, что наше изображение, `img`, поступает как вектор на 784.

Еще один способ предотвратить простое копирование автокодировщиком входа на выход — это добавление *шума* (*noise*). Здесь мы используем “шум” в техническом смысле случайных событий, которые портят исходное изображение, которое в этом контексте будет называться *сигналом* (*signal*). В *бесшумном автокодировщике* (*de-noising autoencoder*) мы добавляем шум к изображению в виде случайно обнуленных пикселей. Обычно около 50% пикселей могут быть ухудшены таким образом. Функция потери AE снова является квадратичной функцией потерь, на этот раз — между пикселями в неповрежденном изображении и выходным изображением декодера.

7.2. Сверточное автокодирование

В последнем разделе был создан AE для чисел Mnist с использованием кодера, чтобы уменьшить начальное изображение с 784 пикселей сначала до 256, а затем — до 128. Сделав это, декодер обращает процесс в том смысле, что, начав со 128 “пикселей”, мы воссоздаем сначала изображение обратно до 256, а затем — до 784 пикселей. Все это было сделано с помощью полно связанных словес. Первый имел весовую матрицу формы [784, 256], второй — [256, 128], а затем, для декодера, — [128, 256] и [256, 784]. Однако, как мы узнали из нашего предыдущего исследования глубокого обучения для компьютерного зрения, наилучшие результаты достигаются при использовании свертки. В этом разделе мы строим AE, используя сверточные методы.

				0	0	0	0	0	0	0	0	0
				0	1	0	2	0	3	0	4	
1	2	3	4	0	0	0	0	0	0	0	0	
4	3	2	1	→	0	4	0	3	0	2	0	1
2	1	4	3	0	0	0	0	0	0	0	0	0
3	4	1	2	0	2	0	1	0	4	0	3	
				0	0	0	0	0	0	0	0	0
				0	3	0	4	0	1	0	2	

Рис. 7.3. Дополнение изображения для декодирования в сверточном AE

Сверточный кодер для снижения размерности изображения не вызывает проблем. В главе 3 мы отметили, как, скажем, два горизонтальных и вертикальных шага уменьшают размер изображения в два раза по каждому измерению. В главе 3 мы не занимались сжатием изображения, поэтому, считая

размер канала (сколько фильтров мы применили к каждому фрагменту изображения), мы фактически получили больше чисел, описывающих изображение в конце процесса свертки, чем в начале (изображение $7 * 7 * 32$ различных фильтра дает 1568). Здесь мы, определенно, хотим, чтобы кодированный промежуточный слой имел намного меньше значений, чем исходное изображение, поэтому мы могли бы, скажем, сделать три слоя свертки; первый уровень довел бы нас до $14 * 14 * 10$, второй — до $7 * 7 * 10$ и третий — до $4 * 4 * 10$ (точные цифры, конечно, гиперпараметры).

Декодирование со сверткой гораздо менее очевидно. Свертка никогда не увеличивает размер изображения, поэтому неясно, как может работать передискретизация. Решение буквально состоит в том, чтобы расширить исходное изображение, прежде чем свернуть его набором фильтров. На рис. 7.3 мы рассматриваем случай, когда скрытый слой AE представляет собой изображение $4 * 4$, и мы хотим расширить его до изображения $8 * 8$. Мы делаем это, окружая каждый “действительный” пиксель достаточным количеством нулей, чтобы получить размер $8 * 8$. (Реальные значения пикселей приведены только для иллюстрации.) Это требует добавления к каждому действительному пикселию нулевых значений слева, по диагонали влево и вверх. Как и следовало ожидать, добавив достаточное количество нулей, мы можем расширить изображение до любого желаемого размера. Затем, свернув это новое изображение с использованием conv2d при шаге 1 и таком же дополнении, мы получим новое изображение $8 * 8$.

```

mnist = input_data.read_data_sets("MNIST_data")

orgI = tf.placeholder(tf.float32, shape=[None, 784])
I = tf.reshape(orgI, [-1, 28, 28, 1])
smallI = tf.nn.max_pool(I, [1, 2, 2, 1], [1, 2, 2, 1], "SAME")
smallerI = tf.nn.max_pool(smallI, [1, 2, 2, 1], [1, 2, 2, 1], "SAME")
feat = tf.Variable(tf.random_normal([2, 2, 1, 1], stddev=.1))
recon = tf.nn.conv2d_transpose(smallerI, feat, [100, 14, 14, 1],
                               [1, 2, 2, 1], "SAME")
loss = tf.reduce_sum(tf.square(recon-smallI))
trainop = tf.train.AdamOptimizer(.0003).minimize(loss)

sess = tf.Session()
sess.run(tf.global_variables_initializer())

for i in range(8001):
    batch = mnist.train.next_batch(100)
    fd={orgI:batch[0]}
    oo,ls,ii,_=sess.run([smallI,loss,recon,trainop],fd)

```

Рис. 7.4. Транспонированная свертка для цифр Mnist

Таким образом, мы получаем изображение подходящего размера, но имеет ли оно соответствующие значения? В доказательство этого на рис. 7.4 приведен код TF, иллюстрирующий передискретизацию с помощью свертки. Здесь мы сначала субдискретизируем (уменьшаем) изображение Mnist, а затем передискретизируем (увеличиваем), используя свертку. Субдискретизация выполняется в два этапа:

```
smallI=tf.nn.max_pool(I,[1,2,2,1],[1,2,2,1],"SAME")
smallerI=tf.nn.max_pool(smallI,[1,2,2,1],[1,2,2,1],"SAME")
```

Первая команда создает версию изображения 14×14 с каждым отдельным оригинальным фрагментом 2×2 , представленным наибольшим значением пикселя в этом фрагменте. Пример изображения 14×14 цифры “7” приведен на рис. 7.5, слева. Вторая команда создает еще меньшую версию — 7×7 . Следующие две строки рис. 7.4 (feat, recon) создают recon, реконструированное изображение 7×7 с передискретизацией до 14×14 , как показано на рис. 7.5, справа. На рис. 7.4 не показано, как мы обычно используем сверточную передискретизацию, которая должна следовать за сверточной субдискретизацией. Мы скорее хотели бы начать с понятного изображения, чтобы лучше видеть, что с ним происходит. Из рис. 7.5 мы видим, что хотя реконструкция едва ли идеальна, она в основном работает. (На рисунке мы заменили нули пробелами, чтобы сделать контур семерки более четким.)

					6	6	2	2	5	5	6	6
9	9	4	5	9	9	5	6	2	2	5	5	5
9	9	9	9	9	9	9	6	6	6	6	6	6
9	9	9	9	9	9	9	7	5	6	5	6	5
8	8		9	9	9		5	5	5	5	6	6
			9	9	5		5	5	5	5	6	5
			6	9	9				6	6	6	6
			8	9	9				5	6	5	6
			9	9	7				1	1	6	4
1	9	9	1						1	1	5	4
4	9	9							2	2	6	6
1	9	4							2	2	5	6

Рис. 7.5. Цифра “7” Mnist размером 14×14 и версия, восстановленная из версии 7×7

Ключевой строкой на рис. 7.4 является вызов conv2d_transpose. Как мы только что упомянули, общим случаем является использование conv2d_transpose для “отмены” использования стандартного conv2d, как в

```
tf.nn.conv2d(img,feat,[1,2,2,1],"SAME")
```

Этот вызов субдискретизирует (`уменьшит`) изображение, которое `conv2d_transpose` может передискретизировать (`увеличить`). Если мы проигнорируем третий аргумент в транспонированной версии, аргументы двух функций будут абсолютно одинаковыми. Однако не все они имеют одинаковый импорт. Да, в обоих случаях первый аргумент — это четырехмерный тензор для манипуляции, а второй — это банк сверточных фильтров для использования. Но `conv2d_transpose` независимо от того, что говорят аргументы шага и дополнения, будет использовать шаг 1 и то же дополнение. Цель этих аргументов состоит скорее в том, чтобы определить, как добавить все лишние нули, что показано на рис. 7.3. Например, чтобы отменить сжатие из-за шага 2, мы обычно хотели бы дополнить каждый действительный пиксель тремя дополнительными нулевыми пикселями, как на рис. 7.3.

К сожалению, невозможно полностью определить размер выходного изображения `conv2d_transpose` только из этой информации. Таким образом, третий аргумент `conv2d_transpose` — это размер желаемого выходного изображения. На рис. 7.4 это [100, 14, 14, 1]: 100 — размер партии, мы хотим получить выходное изображение 14×14 и только один канал. Ситуация, которая вызывает неоднозначность, возникает из-за шагов, больших, чем с таковыми у дополнений *Same*. Рассмотрим, например, два изображения, одно — 7×7 и одно — 8×8 . В обоих случаях при свертке с шагом 2 и дополнением *Same* мы получим изображение размером 4. Таким образом, идя в другую сторону, `conv2d_transpose` с шагом 2 и дополнением *Same* не может знать, какое из этих выходных изображений желает пользователь.

На этом этапе разговора о транспонированной свертке мы сконцентрировались на том, чтобы удостовериться, что входные данные дополнены таким образом, чтобы получить желаемый эффект передискретизации. Мы не заботились о том, как фильтры справляются с этой задачей. Действительно, в начале обучения их нет. На рис. 7.6 показано изображение с передискретизацией из нулевого учебного примера. Преобладающим видимым эффектом на рисунке являются чередования 0 и -1, что, без сомнения, является артефактом, возникающим в результате чередования значений заполнения нулями и ненулевых значений действительных пикселей в изображении, переданном в `conv2d_transpose`. Существует математическая теория о том, как транспонированная свертка может найти правильные значения ядра, но для наших целей нам нужно только знать, что переменные значения ядра и обратное распространение делают это вместо нас.

Что касается использования `conv2d_transpose` для автоматического кодирования, оно в точности аналогично полносвязному автокодированию. У нас есть один или несколько уровней субдискретизации с использованием `conv2d` (возможно, с использованием `max_pool` или `avg_pool`), за которым обычно

следует равное количество уровней передискретизации с использованием соп-
v2d_transpose.

0	0	0	0	1	-1	1	-1	1	-1
0	0	0	0	-1	0	-1	0	-1	0
0	-1	1	-1	1	-1	0	-1	0	0
-1	0	-1	0	-1	0	-1	0	0	0
0	0	1	-1	1	-1	1	-1	0	-1
0	0	-1	0	-1	0	-1	0	-1	0
0	-1	0	-1	0	0	1	-1	1	-1
-1	0	-1	0	0	0	-1	0	-1	0
1	-1	1	-1	1	-1	1	-1	0	0
-1	0	-1	0	-1	0	-1	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

Рис. 7.6. Передискретизированная маленькая цифра Mnist из нулевого учебного примера

7.3. Вариационное автокодирование

Вариационный автокодировщик (Variational Autoencoder — VAE) — это вариант АЕ, задача которого заключается не в том, чтобы точно воспроизвести исходное изображение, а в том, чтобы создать новое изображение, являющееся членом того же класса изображений, но узнаваемое как новое. Он получил свое название от *вариационных методов* (variational method), темы в байесовском машинном обучении. Опять же, мы возвращаемся к нашему дружественному набору данных Mnist. Наша первоначальная цель для VAE подразумевает ввод цифрового изображения Mnist и получение нового изображения, которое будет выглядеть, как несколько другое.

Если нас не заботит то, что новое изображение будет заметно отличаться, стандартное автоматическое кодирование в значительной степени решит проблему. В конце концов, посмотрите еще раз на рис. 7.2: наша АЕ-реконструкция “7” с начала книги (см. рис. 1.1). В то время мы гордились тем, насколько они похожи, но, конечно, не идентичны. Однако они настолько близки, что если бы мы распечатали эту более позднюю версию в оттенках серого, было бы очень трудно отличить ее от рис. 1.2. Кроме того, некоторые размышления должны указывать на то, что стандартная АЕ с квадратичной функцией потерь — не совсем то, чего мы хотим. Рассмотрим три изображения 7×7 на рис 7.7. Верхнее — это небольшое изображение 1, а два других должны быть реконструкциями первого. Первое из них выглядит аналогично оригиналу, но, поскольку цифра сдвинута на два пикселя, перекрывающихся значений нет, а ее среднеквадратичная ошибка по сравнению с верхним изображением на

рис. 7.7 равна 10. Кроме того, правое нижнее изображение имеет вид, больше похожий на нашу функцию потерь, поскольку она отличается только на два пикселя. Таким образом, стандартное автокодирование и квадратичная функция потерь не очень подходят для нашей задачи. Тем не менее на время отложите это возражение, чтобы сконцентрироваться на том, как работает VAE. Мы вернемся к этому вопросу позже и покажем, как VAE решают эту проблему.

Рис. 7.7. Оригинальное изображение и две “реконструкции”

Заглянем под капот того, как наша программа собирается ввести изображение, а затем вызвать в воображении вектор случайных чисел. И именно эти случайные числа определяют разницу между исходным и новым изображениями: помещают в одно и то же изображение пары случайных чисел — и мы получаем точно такой же вариант исходного изображения. Позже мы увидим, что мы можем опустить изображение, и в этом случае мы получим не вариацию конкретного изображения, а скорее совершенно новое изображение в общем стиле всех изображений. Обычно это выглядит как одна из возможных цифр, но в зависимости от того, насколько хорошо VAE выполняет свою работу, это может быть и не так. Если вы перейдете к рис. 7.10, то увидите несколько примеров.

На рис. 7.8 показана схема архитектуры VAE. Изображение поступает в нижней части схемы, и мы можем проследить вычисления кодировщика. Затем кодированная информация используется для создания не одного векторного представления изображения, а двух векторов действительных значений, σ и μ . Затем мы строим векторное представление нового изображения, генерируя вектор случайных чисел r и вычисляя $\mu + \sigma r$. Здесь μ можно рассматривать как исходную версию векторного представления изображения, и мы варьируем его с помощью σr , чтобы получить другое векторное представление, которое близко, но не слишком, к оригиналу. Затем это пересмотренное

векторное представление проходит через стандартный декодер для создания нового изображения. Если мы используем программу (в отличие от обучения), то это новое изображение выводится пользователю. Если мы обучаемся, новое изображение поступает в слой “потерь изображения” (image loss). Потери изображения — это просто квадратичная потеря разности значений пикселей между исходным и новым изображениями. Таким образом, источником вариации выходного изображения в VAE является r . Введя одно и то же изображение и один и тот же вектор случайных чисел, мы получим одно и то же изображение.

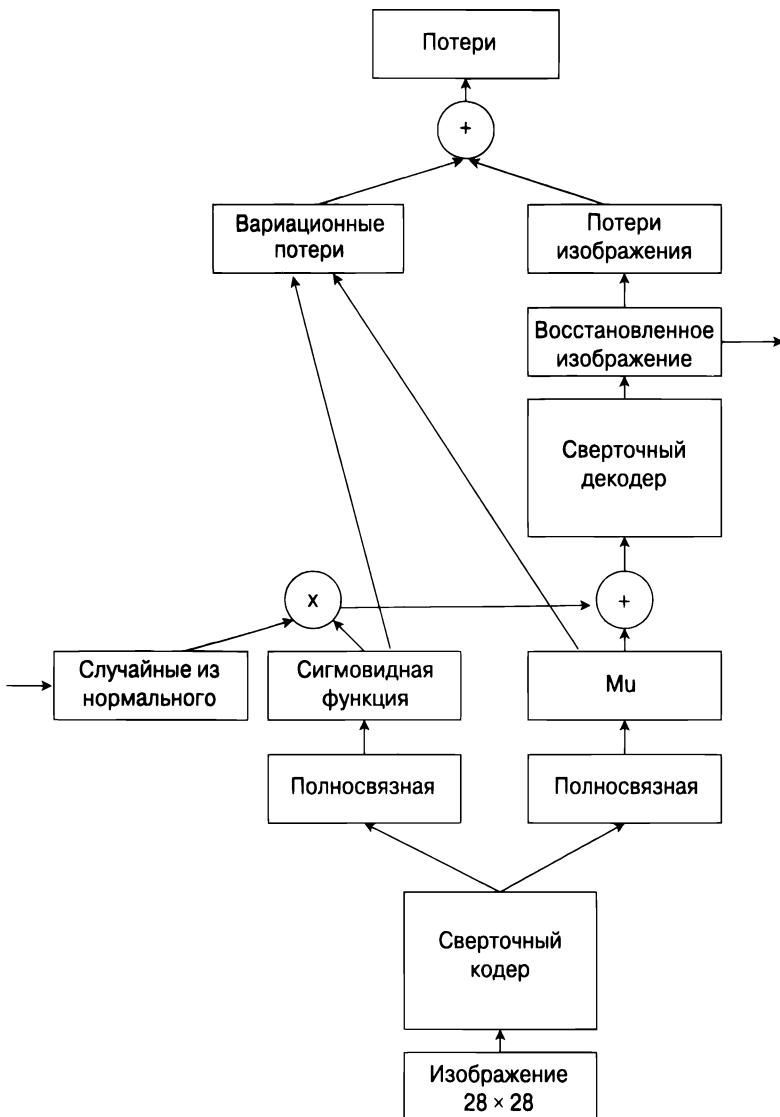


Рис. 7.8. Структурная модель вариационного автокодировщика

Повторим это немного другими словами: μ — это базовая кодированная версия исходного изображения, как и всегда в этой главе, а σ определяет допустимые границы того, насколько мы можем варьировать исходное изображение и при этом иметь “узнаваемую” версию. Помните, что и σ , и μ являются векторами действительных чисел размером, скажем, 10. Если $\mu[0] = 1,12$, это означает, что кодированное изображение должно иметь в качестве первого действительного числа нечто более или менее близкое к 1,12 с вариацией около 1,12, контролируемой $\sigma[0]$. Если $\sigma[0]$ велико (и наша NN работает должным образом), это означает, что можно немного изменить первое измерение кодированного изображения, не делая выходную версию неузнаваемой. Если $\sigma[0]$ мало, то это не так.

Но здесь есть большая проблема. Мы предположили, что каким-то образом σ и μ , которые мы получим от кодировщика, будут такими числами, что $\mu + r * \sigma$ является разумной кодировкой изображения. Если потери — это просто квадратичная функция потерь, мы *не получим* желаемого результата.

Именно здесь VAE склоняются к глубокой математике, но вместо этого мы попросим читателя принять некоторые маловероятные изменения в нашей функции потерь. VAE имеют две потери, которые складываются вместе, чтобы дать общие потери. Одна из них, которую мы уже видели, — это квадратичная функция потерь между пикселями исходного изображения и пикселями реконструкции. Мы называем это *потерями изображения* (image loss).

Вторая из потерь — это *вариационные потери* (variational loss):

$$L_v(\mu, \sigma) = -\sum_i \frac{1}{2} (1 + 2\sigma[i] - \mu[i]^2 - e^{2\sigma[i]}) \quad (7.4)$$

Это для одного примера, и это точечные вычисления над σ и μ . (Помните: они оба — векторы.) Чтобы сделать это более понятным, рассмотрим, что произойдет, если сначала σ , а затем μ установить равными нулю:

$$L_v(\mu, 0) = \mu^2 \quad (7.5)$$

$$L_v(0, \sigma) = -\frac{1}{2} - \sigma + \frac{e^{2\sigma}}{2} \quad (7.6)$$

Из первого из этих уравнений мы видим, что NN поощряется к сохранению средних значений $\mu = 0$. Конечно, потери изображения будут противодействовать этому поощрению, но при прочих равных L_v хочет $\mu \approx 0$.

Второе из приведенных выше уравнений будет поддерживать σ около 1. Когда σ меньше 1, доминирует второй член $L_v(0, \sigma)$, и мы уменьшаем потери, увеличивая σ , тогда как когда σ больше 1, третий член начинает расти довольно быстро.

То, что это выглядит как стандартное нормальное распределение, не является совпадением. Если бы мы рассмотрели математику, мы бы поняли а), почему имеет смысл, чтобы кодирование изображения выглядело больше как стандартное нормальное распределение, и б) почему минимальная вариационная потеря (σ) не совсем равна 1. Вместо этого мы просим читателя согласиться с тем, что добавление этой второй функции потерь к нашему общему вычислению дает нам то, что мы хотим: NN, которая, учитывая многомерные σ и μ , вычисленные нашими NN на основе действительного изображения Mnist, производит кодировку нового, немного другого изображения I' из

$$I' = \mu + r\sigma \quad (7.7)$$

Здесь r — вектор случайных чисел, полученных из стандартного нормального распределения. Как только мы получим эту гарантию, у нас будет наш VAE.

Давайте вернемся к проблеме, которую мы упомянули ранее (с. 154). Стандартный AE с квадратичной функцией потерь не соответствует цели, которую мы поставили для VAE: создавать заметно различающиеся версии изображения, которые, тем не менее, похожи на оригинал. Нашим примером проблемы на рис. 7.7 были две мнимые реконструкции 1; одна перевела два пикселя по горизонтали, одна пропустила два пикселя в середине вертикального штриха. В соответствии с квадратичной функцией потерь вторая была более похожа, но первая будет хорошей реконструкцией VAE, а вторая — нет. Утверждается, что, работая должным образом, VAE преодолевают эту трудность.

В первую очередь, обратите внимание на то, что при обучении VAE мы используем квадратичную функцию потерь, но по своей природе VAE должны получать большие квадратичные потери, поскольку они могут восстанавливать только исходное изображение до некоторой заранее заданной случайности. Далее отметим, что эта случайность кроется в кодировке изображения в середине архитектуры VAE. Утверждается, что это лучшее место, если VAE должны работать надлежащим образом.

В основном неявным, но иногда явным в нашем обсуждении AE является наблюдение, что они достигают снижения размерности, отмечая общие черты между входными данными. AE приспособливают векторное представление, чтобы понять общие черты, только “упоминая” различия. Предположим, что, как и следовало, цифры Mnist немного различаются по своему положению на странице. Тогда один из способов использования этого факта для уменьшения кодировки заключается в том, чтобы, скажем, одно из действительных чисел в кодировке указывало общую горизонтальную позицию цифры. (Или, возможно, имея только двадцать действительных значений, в которых можно закодировать среднее значение для цифры 1, мы не можем позволить себе посвятить действительное число этой работе, и наша AE “решает”, что некоторые другие

варианты более важны для хорошей реконструкции нашего изображения. Это только иллюстрация.) Дело в том, что если существует действительное число, которое кодирует общую горизонтальную позицию, то нижнее левое изображение на рис. 7.7 фактически очень близко к верхнему оригиналу — оно отличается только в одной кодированной позиции.

В любом случае VAE работают. На рис. 7.9 показаны оригинал 4 Mnist и новая версия. Даже нескольких секунд обучения достаточно, чтобы убедить вас, что они разные. Кроме того, они различаются довольно типичным для VAE способом, правильное изображение и реконструкция, являются менее различными. Это размытие 4, если хотите. Наиболее заметно, что левая четверка (оригинал) имеет хвостик внизу основного вертикального штриха, который полностью отсутствует справа.

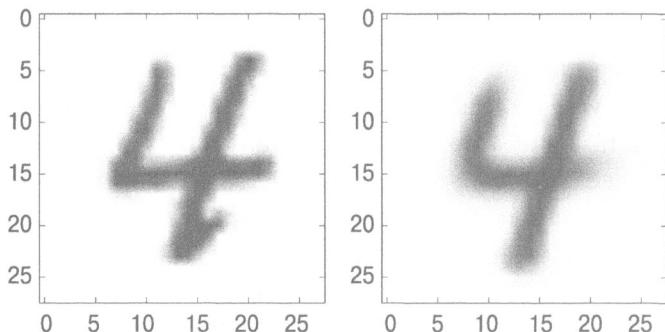


Рис. 7.9. Оригинальное изображение и реконструкция VAE

До сих пор мы рассматривали проблему генерации изображения, которое похоже на исходное изображение, но заметно отличается от него. Однако ранее мы отмечали, что VAE также могут работать более свободно — учитывая общий класс изображений, создают другой член этого класса. Это гораздо более сложная проблема, но VAE практически не меняется. Обучение, по сути, точно такое же. Разница лишь в том, как мы используем VAE. Чтобы создать новое изображение, не основанное на существующем, VAE генерирует случайное число (опять-таки, из стандартного нормального распределения), но на этот раз вставляет его (через `feed_dict`) для использования в качестве всего кодирования изображения. На рис. 7.10 показаны некоторые примеры результатов той же программы, которая генерировала примеры на рис. 7.9, но на этот раз не было предоставлено изображение для эмуляции. Четыре из изображений являются узнаваемыми цифрами, подобными цифрам Mnist, но нижнее правое изображение выглядит как “3” или “8”, а то, которое перед ним, может быть очень плохим подобием “8”. Более мощная модель с намного большим количеством эпох обучения и большим вниманием к гиперпараметрам даст гораздо лучший результат.

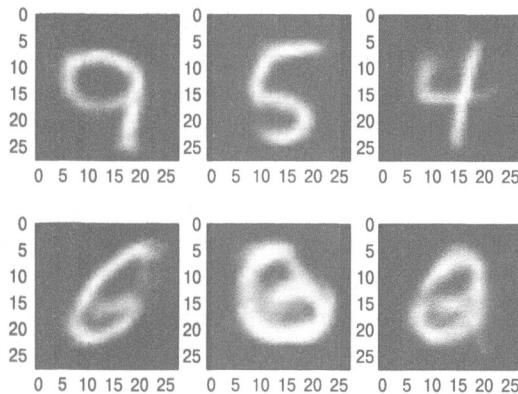


Рис. 7.10. Цифры Mnist, созданные VAE “с нуля”

Прежде чем закончить рассмотрение VAE, несколько слов для тех, кто хотел бы лучше понять вариационные потери уравнения 7.4. В нашей первоначальной формулировке мы генерируем новое изображение I' из оригинала I . Для этого мы используем сверточный кодер C , чтобы создать уменьшенное представление $z = C(I)$. Здесь мы сосредоточимся на двух вероятностных распределениях, $\Pr(z | I)$ и $\Pr(z)$. В частности, сначала мы предполагаем, что оба являются нормальными распределениями.

Если вы только что спокойно перешли от последнего предложения к этому, либо вы не слишком задумывались о том, что здесь написано, либо вы намного умнее меня. Как мы можем просто предположить, что распределение чего-то столь же сложного, как образы реального мира или даже цифры MNIST, — это так же просто, как нормальное распределение? Следует признать, что это будет n -мерная нормаль, где n — это размерность представления изображения, создаваемого сверточным кодером C , но даже n -мерные нормали — это довольно простые вещи (двумерная нормаль имеет форму колокола).

Ключевая идея, которую нужно иметь в виду, заключается в том, что $\Pr(z)$ — это распределение вероятностей, основанное не на исходном изображении, а на представлении $C(I)$. Возможно, вы помните, что раньше у нас мог быть элемент представления, показывающий, как далеко центр цифры находится от центра изображения, так что цифра “1” в верхней части рис. 7.7 имеет представление, очень похожее на нижний правый вариант. Предположим, что мы можем довести эту программу до логического конца, чтобы каждый параметр в векторе представления был числом, подобным этому, некоторые спецификации основных способов, которыми одно изображение числа может отличаться от другого. Другими примерами могут быть “ положение главного штриха ” и “ диаметр верхнего круга ” (для восьмерок) и т.д. Как вы можете себе представить, диаметр верхнего круга у восьмерок может, как правило, составлять,

скажем, 12 пикселей, но может быть расположен значительно выше или ниже. В таком случае было бы разумно описать изменение этого параметра как нормальное со средним значением 12 и стандартным отклонением, скажем, 3. Что касается $\text{Pr}(z | I)$, которое является нормальным, применимы аналогичные рассуждения. Когда дано изображение, мы хотим, чтобы VAE производил много разных похожих, но отличающихся изображений. Вероятность любого из них может быть нормальным распределением по ключевым факторам, таким как диаметры окружностей в восьмерках.

Предстоит сделать еще много шагов, прежде чем мы доберемся до вариационных потерь в уравнении 7.4, но мы просто собираемся взять один или два из них. Далее мы предполагаем, что $\text{Pr}(z)$ является стандартным нормальным — нормальным распределением с $\mu = 0$ и $\sigma = 1$, записанными как $N(0, 1)$. С другой стороны, мы предполагаем, что выход кодера является нормальным распределением, среднее значение и стандартное отклонение которого зависят от самого изображения, например $N(\mu(I), \sigma(I))$. Это объясняет, а) почему у нас есть кодировщик на рис. 7.8, приводящий к двум значениям, помеченным как μ и σ , и б) почему, чтобы получить случайное изменение, мы выбрали числа в соответствии со стандартным нормальным распределением.

Мы делаем последний шаг. Предположение о том, что одно из нормальных распределений является стандартным, а другое не может быть таковым, в целом удовлетворительно. Скорее, мы просто пытаемся минимизировать расхождение. Мы можем более точно смоделировать конкретное изображение, если мы можем выбирать подходящие μ и σ и нас подталкивают потери изображения в этом направлении. С другой стороны, мы хотим, чтобы эти значения были как можно ближе к 0 и 1, насколько это возможно, и вот именно здесь и возникают вариационные потери.

Другими словами, мы хотим минимизировать разницу между двумя вероятностными распределениями $N(0, 1)$ и $N(\mu(I), \sigma(I))$. Стандартной мерой разницы между двумя распределениями является *расстояние Кульбака–Лейблера* (Kullback-Leibler divergence):

$$D_{KL}(P \| Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)} \quad (7.8)$$

Например, если $P(i) = Q(i)$ для всех i , то соотношение двух всегда равно 1, а $\log 1 = 0$. Теперь мы можем охарактеризовать цель VAE как минимизацию потерь изображения при одновременной минимизации $D_{KL}(N(\mu(I), \sigma(I)) \| N(0, 1))$. К счастью, существует решение в закрытой форме для минимизации последних, и с некоторой дополнительной алгеброй (и умной идеей или двумя) это приводит к функции вариационных потерь, представленной выше.

7.4. Генеративно-состязательные сети

Генеративно-состязательные сети (Generative Adversarial Network — GAN) представляют собой модели NN без учителя, которые работают за счет установки двух конкурирующих одна с другой моделей NN. В случае Mnist первая сеть (называемая *генератором* (generator)) генерирует цифру Mnist “с нуля”. Вторая, *дискриминатор* (discriminator), получает выходной сигнал генератора или пример действительной цифры Mnist. Его выходные данные — это оценка вероятности того, что данные, которые он получил, взяты из действительного примера, а не созданы генератором. Впоследствии решение дискриминатора действует как сигнал об ошибке для обеих моделей, но в “противоположных направлениях” в том смысле, что, если он уверен в правильности решения, это означает большую ошибку для генератора (он не обманул дискриминатор), тогда как, если дискриминатор одурачен, это является большой потерей для дискриминатора.

Как и в случае с AE, GAN можно использовать для изучения структуры входных данных без меток. Кроме того, как и в случае с VAE, GAN может генерировать новые варианты для класса, включая классы изображений, и, в принципе, почти для всех. GAN являются горячей темой в глубоком обучении, поскольку в некотором смысле они являются *универсальной функцией потерь* (universal loss function). Любой набор картинок, текстов, решений по планированию и тому подобного, для которого у нас есть данные, может использоваться для обучения GAN без учителя с теми же базовыми потерями.

Базовая архитектура GAN показана на рис. 7.11. Чтобы понять, как это работает, мы возьмем тривиальный пример: дискриминатору даются два числа, по одному за раз. Одно из них — “действительные” данные, генерируемые нормальным распределением, например, со средним значением 5 и стандартным отклонением 1, поэтому производимые ими числа в основном находятся между 3 и 7. “Поддельное” число создается генератором, являющимся однослойной NN. (В этом разделе только противоположность “действительному” числу является поддельное, а не комплексное.) Генератору присваивается случайное число, которое с равной вероятностью может находиться в диапазоне от -8 до +8. Чтобы обмануть дискриминатор, он должен научиться изменять случайное число, предположительно, чтобы оно оказалось между 3 и 7. Первоначально генератор NN имеет параметры, близкие к нулю, так что фактически он в основном генерирует числа, близкие к нулю.

GAN используют несколько аспектов TF, которые мы еще не рассмотрели, и на рис. 7.12 мы приводим полный код простой GAN. Мы пронумеровали каждый небольшой раздел кода для справки.

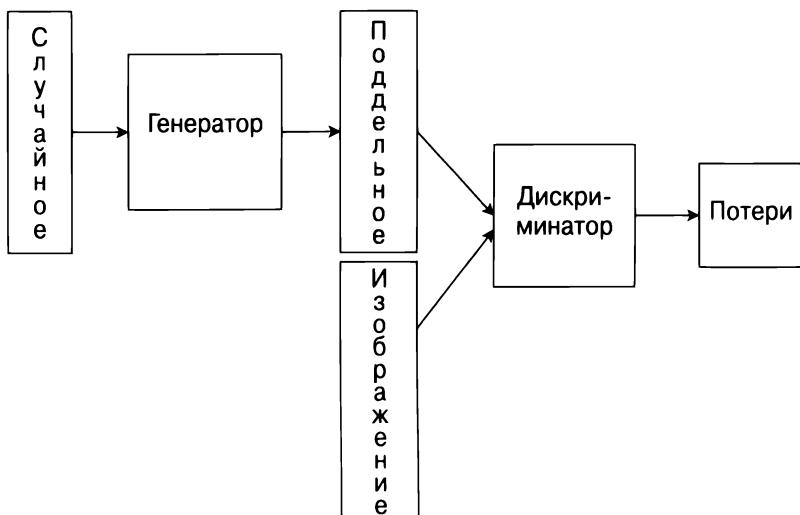


Рис. 7.11. Структура генеративно-состязательной сети

Сначала рассмотрим раздел 7, в котором мы обучаем GAN. Выберите основной обучающий цикл, который мы установили на 5001-й итерации. Первое, что мы делаем, — генерируем некоторые действительные данные: случайное число около 5 и случайное число от –8 до 8 для передачи генератору. Затем мы по отдельности обновляем сначала дискриминатор, а затем генератор. Наконец каждые 500 итераций мы выводим данные для отслеживания.

Мы скоро вернемся к этому разделу кода, но сначала рассмотрим, как все должно работать. Мы хотим, чтобы дискриминатор вывел единственное число (o), которое предназначено для вероятности того, что число, которое он только что видел, получено из действительного распределения. Обратите внимание на раздел 3 кода. Когда мы выполняем функцию `discriminator`, она устанавливает четырехслойную полносвязную NN с прямой связью. Первые три слоя имеют функции активации `relu`, а последний использует сигмовидную функцию. Поскольку сигмовидные функции выводят числа от 0 до 1, как мы отмечали ранее, они хороши для получения чисел, предназначенных для вероятностей (вероятность того, что входные данные получены из действительного распределения). Более конкретно, выходные данные интерпретируются как вероятность того, что входные данные являются членами класса (класса действительных входных данных). Это определяет наш выбор функции потерь — мы используем кросс-энтропийные потери.

```

bsz, hSz, numStps, logEvery, genRange = 8, 4, 5000, 500, 8

1 def log(x): return tf.log(tf.maximum(x, 1e-5))

2 with tf.variable_scope('GEN'):
    gIn = tf.placeholder(tf.float32, shape=(bsz, 1))
    g0=layers.fully_connected(gIn, hSz, tf.nn.softplus)
    G=layers.fully_connected(g0,1,None)
    gParams =tf.trainable_variables()

3 def discriminator(input):
    h0 = layers.fully_connected(input, hSz*2,tf.nn.relu)
    h1=layers.fully_connected(h0,hSz*2, tf.nn.relu)
    h2=layers.fully_connected(h1,hSz*2, tf.nn.relu)
    h3=layers.fully_connected(h2,1, tf.sigmoid)
    return h3

4 dIn = tf.placeholder(tf.float32, shape=(bsz, 1))
with tf.variable_scope('DIS'):
    D1 = discriminator(dIn)
    with tf.variable_scope('DIS', reuse=True):
        D2 = discriminator(G)
    dParams = [v for v in tf.trainable_variables()
               if v.name.startswith('DIS')]

5 gLoss=tf.reduce_mean(-log(D2))
dLoss=0.5*tf.reduce_mean(-log(D1) -log(1-D2))
gTrain=tf.train.AdamOptimizer(.001).minimize(gLoss, var_list=gParams)
dTrain=tf.train.AdamOptimizer(.001).minimize(dLoss, var_list=dParams)

6 sess = tf.Session()
sess.run(tf.global_variables_initializer())

7 gmus,gstds=[], []
for i in range(numStps+1):
    real=np.random.normal(5, 0.5, (bsz,1))
    fakeRnd= np.random.uniform(-genRange,genRange, (bsz,1))
    # обновить дискриминатор
    lossd,gout, _ = sess.run([dLoss,G,dTrain], {gIn:fakeRnd, dIn:real})
    gmus.append(np.mean(gout))
    gstds.append(np.std(gout))
    # обновить генератор
    fakeRnd= np.random.uniform(-genRange,genRange, (bsz,1))
    lossg, _ = sess.run([gLoss, gTrain], {gIn:fakeRnd})
    if i % logEvery == 0:
        frm=np.max(i-5,0)
        cmu=np.mean(gmus[frm:(i+1)])
        cstd=np.mean(gstds[frm:(i+1)])
        print('{:}\t{:.\.3f}\t{:.\.3f}\t{:.\.3f}\t{:.\.3f}\t'.
              format(i, lossd, lossg, cmu, cstd))

```

Рис. 7.12. GAN для изучения среднего значения нормального распределения

Когда дискриминатору передают число из действительного нормального распределения (n_r), потери являются отрицательным логарифмом вывода NN дискриминатора (o_r). Когда дискриминатору передается созданное поддельное число, потери составляют $\ln 1 - o_f$. В разделе 7 на рис. 7.12 показано, что в обучающем цикле мы сначала создаем несколько действительных чисел, а затем несколько случайных чисел, которые передаются генератору. Мы делаем это сразу после комментария “обновить дискриминатор”, предоставляя оптимизатору Adam и то, и другое. Функция потерь дискриминатора L_d определяется выражением

$$L_d = \frac{1}{2}(-\ln(o_r) - \ln(1 - o_f)) \quad (7.9)$$

Взглянув на раздел 5, в котором мы выкладываем код потерь и обучения, вы увидите потери дискриминатора, определенные, как в уравнении 7.9. Здесь o_r — это степень, в которой дискриминатор считает реальный вклад действительно реальным. И наоборот, потери также включают в себя член, который штрафует дискриминатор в той степени, в которой он считает, что поддельное (созданное генератором) число (o_f) является действительным.

Итак, подумайте, что может произойти на первом примере обучения. Как мы уже отмечали, генератор выдает что-то близкое к нулю, скажем, 0,01. Действительная выборка из нашего нормального распределения может быть 3,8. Тем не менее параметры дискриминатора также инициализируются значением, близким к нулю, так что o_r и o_f оба близки к нулю. Первоначально в L_d будет доминировать член $-\ln(o_r)$, который уходит в отрицательную бесконечность, когда o_r стремится к нулю, но программа быстро учится не назначать какие-либо вероятности слишком близко к нулю.

Более важным является то, как производная потерь влияет на параметры дискриминатора. Оглядываясь назад на однослойную сеть из главы 1, мы видим, что производная весовых параметров пропорциональна входным данным для рассматриваемого веса (уравнение 1.22). Если мы рассмотрим математику, то обнаружим, что, когда мы даем дискриминатору действительную выборку, вес перемещается вверх пропорционально значению выборки (например, 3,8). И наоборот, тот же вес перемещается вниз, когда мы обучаемся на поддельном значении генератора, но это всего лишь 0,01, поэтому дискриминатор слегка движется к тому, чтобы сказать, что более высокие входные значения являются действительными.

Далее мы рассмотрим, как программа должна оценивать производительность генератора. Что ж, генератор хочет обмануть дискриминатор в том смысле, что, когда последний получает ввод от первого, генератор хочет, чтобы дискриминатор думал, что он только что увидел действительные, а не поддельные числа. Таким образом, генератор потерь L_g должен быть

$$L_g = -\ln(o_f) \quad (7.10)$$

Читатель может убедиться, что первая строка раздела 5 кода GAN определяет потери генератора именно таким образом.

Итак, подведем итог сказанному до сих пор: раздел 7, основной цикл обучения, сначала обучается дискриминатор, выдавая некоторые действительные и некоторые ложные числа с функцией потерь, которая штрафует ошибки в обоих направлениях, когда действительные оцениваются как ложные, и наоборот. Затем обучается генератор с потерями, основываясь только на том, насколько правильно дискриминатор идентифицирует подделки генератора.

К сожалению, все немного сложнее. В первую очередь, обратите внимание на то, что в разделе 4 есть два вызова кода для настройки дискриминатора. Причина в том, что в TF нельзя передавать данные в одну сеть раздельно из двух разных источников (в отличие от передачи в сеть объединения двух тензоров). Таким образом, мы создаем в некотором смысле два дискриминатора. Один, D1, получает ввод от действительного распределения, тогда как второй получает на вход фальшивку от генератора D2.

Конечно, нам не нужны две отдельные сети, поэтому, чтобы объединить их, мы настаиваем на том, чтобы две сети использовали одни и те же параметры, выполняя, таким образом, одну и ту же функцию и не занимая гораздо больше места, чем одна версия. Это цель вызовов функции `tf.variable_scope`, которые мы видим в разделе 4. Сначала мы использовали эту функцию TF в главе 5 (см. с. 165). Там нам потребовались две модели LSTM и нам нужно было избегать конфликтов имен, поэтому мы определили одну — в переменной “enc” (кодировщик) и одну — в переменной “dec” (декодер). В первой из них все определенные переменные должны иметь часть “enc” перед именами, а во втором — “dec”. Здесь у нас противоположная проблема. Мы хотим избежать двух отдельных наборов, поэтому даем обеим областям одинаковые имена и во втором вызове явно указываем TF повторно использовать одни и те же переменные, добавляя `reuse = True`.

Есть одно последнее осложнение, с которым нужно справиться, прежде чем двигаться дальше. В разделе 7 на рис. 7.12 демонстрируется, что мы сначала обучаем дискриминатор, а затем — генератор на каждом этапе обучения. Кроме того, чтобы обучить дискриминатор, мы передаем TF два случайных числа, действительную выборку и случайное число для передачи в генератор. Затем мы запускаем генератор, чтобы создать поддельное число, и вероятность, которую дискриминатор назначает этому числу как полученному из действительного распределения o_f . Последнее является частью потерь, определенных в уравнении 7.9. Однако в отсутствие какой-либо специальной обработки на обратном проходе параметры генератора также модифицируются, причем, что еще хуже, модифицируются таким образом, чтобы уменьшить потери дискриминатора.

Как отмечалось выше, это *не то*, чего мы хотим. Мы хотим, чтобы параметры генератора двигались так, чтобы сделать задачу дискриминатора более сложной. Таким образом, запуская обратное распространение и изменяя параметры дискриминатора в разделе 5 (см. строку, задающую `dTrain`), мы инструктируем TF, только чтобы изменить один из двух наборов параметров. В частности, обратите внимание на именованный аргумент для `AdamOptimizer`. Этот аргумент указывает оптимизатору (при данном конкретном вызове) изменять только переменные TF в списке.

Что касается определения классов параметров, `gParams` и `dParams`, обратитесь к концу разделов 2 и 4. Функция TF `trainable_variables` возвращает тензор всех переменных в графе TF, определенных до этой точки.

7.5. Ссылки и что читать дальше

Истоки автокодирования в NN, наверное, потеряны в глубине веков. В учебнике Гудфеллоу и других [7] цитируется кандидатская диссертация Янна Лекуна [33] как самое раннее упоминание этого предмета.

Из тем, упомянутых в этой главе, первое, что, кажется, достигло определенной идентичности в сообществе NN, — это вариационные автокодировщики. Стандартной ссылкой здесь является статья Дидерика Кингмы и Макса Веллинга [34]. Я нашел блог Феликса Мора [35] очень полезным, и мой код основан на его коде. Для тех, у кого есть статистические предпосылки и кто хочет изучить математику VAE, я нахожу учебник Карла Доерша по VAE хорошим справочным материалом [36].

Однако история GAN совершенно ясна: основная идея возникла в ее нынешнем виде в статье Айана Гудфеллоу и других [37]. Я многому научился из блога Джона Гловера [38]. Мой код для GAN по обучению нормального распределения основан на его коде, и он, в свою очередь, зачитывает [39].

Комментарий о том, что GAN являются “универсальными” функциями потерь (с. 161), взят из выступления Филиппа Изолы [40], в котором представлено много интересных способов организации обучения без учителя визуальной обработке, с особым акцентом на использование GAN.

7.6. Упражнения

Упражнение 7.1. Подумайте о субдискретизации с полносвязными слоями. Цифры Mnist всегда имеют ряды нулей по краям. Учитывая наш комментарий о том, что AE работают, когда кодирование изображения игнорирует общие черты, что это означает (при прочих равных) относительно значений обученных весов первого уровня?

Упражнение 7.2. Та же ситуация и вопрос, что и в упражнении 7.1, но теперь об обученных весах для последнего слоя перед выводом.

Упражнение 7.3. Как вызов `conv2d_transpose` выполняет транспонирование, показанное на рис. 7.3?

Упражнение 7.4. С учетом того, что `img` является массивом $2 * 2$ пикселей из чисел от 1 до 4, покажите дополненную версию следующего вызова `conv2d_transpose`:

```
tf.nn.conv2d transpose(img, flts, [1, 6, 6, 1], [1, 3, 3, 1], "SAME")
```

Можете игнорировать первый и последний компоненты формы.

Упражнение 7.5. Это не важно, но почему мы установили цикл обучения GAN на 5001-ю итерацию на рис. 7.12, а не на 5000-ю?

Упражнение 7.6. GAN на рис. 7.12 выводит как среднее значение созданных данных, которое мы использовали для оценки точности GAN, так и стандартное отклонение созданных данных, которое мы игнорировали выше. На самом деле, несмотря на то что мы установили стандартное отклонение действительных чисел равным 0,5 (укажите, где это делается!), фактическая σ , выводимая после каждого 500 итераций, начиналась выше, быстро уменьшаясь ниже 0,5, и, кажется, направлялась еще ниже. Объясните, почему эта модель GAN не оказывает давления на нее, чтобы узнать правильную σ , и почему разумно, что фактическое значение, которое она предлагает, ниже, чем действительное.

Приложение

Ответы к некоторым упражнениям

Глава 1

Упражнение 1.1. Если a — это значение цифры в первом примере обучения, то после обучения на этом одном примере должно увеличиваться только значение b_a , а для всех значений цифр $a' \neq a$, $b_{a'}$ должно уменьшаться.

Упражнение 1.2. а) Логиты прямого прохода: $(0 * 0,2) + (1 * (-0,1)) = -0,1$ и $(0 * (-0,3)) + (1 * 0,4 * 1) = 0,4$ соответственно. Для вычисления вероятностей мы сначала вычисляем знаменатель softmax $e^{-0,1} + e^{0,4} = 0,90 + 1,50 = 2,40$. Тогда вероятности составляют $0,9/2,4 = 0,38$ и $1,5/2,4 = 0,62$. б) потеря составляет $-\ln 0,62 = -1,9$. Из уравнения 1.22 мы видим, что $\Delta(0, 0)$ будет произведением членов, включающих $x_0 = 0$, поэтому $\Delta(0, 0) = 0$.

Упражнение 1.5. Вычисление произведения матриц дает

$$\begin{pmatrix} 4 & 7 \\ 8 & 15 \end{pmatrix} \quad (\text{A.1})$$

Затем, добавив правый вектор к обеим строкам, получаем

$$\begin{pmatrix} 8 & 12 \\ 12 & 19 \end{pmatrix} \quad (\text{A.2})$$

Упражнение 1.6. Производная квадратичной потери по b_j вычисляется почти так же, как для кросс-энтропийных потерь, за исключением

$$\frac{\partial L}{\partial l_j} = \frac{\partial}{\partial l_j} (l_j - t_j)^2 = 2(l_j - t_j) \quad (\text{A.3})$$

Поскольку производная от l_j как функция от b_j равна 1, получаем

$$\frac{\partial L}{\partial b_j} = 2(l_j - t_j) \quad (\text{A.4})$$

Глава 2

Упражнение 2.1. Если мы не указываем индекс снижения, он считается равным нулю, и в этом случае мы добавляем столбцы. Это дает $[0, 3, 2, 9]$.

Упражнение 2.2. Эта новая версия значительно медленнее, чем оригинальная, поскольку на каждой итерации основного цикла обучения она создает новый оптимизатор градиентного спуска, а не использует один и тот же каждый раз.

Упражнение 2.4. Вы не можете взять `tensordot`, поскольку размеры не равны. Если можно, первый тензорный аргумент имеет форму $[4, 3]$, а второй — $[2, 4, 4]$. Просто объединив их, вы получите $[4, 3, 2, 4, 4]$. Поскольку мы берем скалярное произведение 0-го компонента первого тензора и 1-го компонента второго, они выпадут, дав форму результата $[3, 2, 4]$.

Глава 3

Упражнение 3.1. а) Одним из примеров будет

$$\begin{matrix} -2 & 1 & 1 \\ -2 & 1 & 1 \\ -2 & 1 & 1 \end{matrix}$$

Что касается части б), то дело в том, что таких ядер бесконечно много. Умножение чисел в вышеупомянутом ядре на любое положительное число будет примером.

Упражнение 3.5. С точки зрения синтаксиса единственное отличие — это шаг $[1, 1, 1, 1]$, а не более ранние $[1, 2, 2, 1]$. Таким образом, мы применяем `maxpool` к каждому фрагменту 2×2 , а не к любому другому. Когда шаг `maxpool` равен 1, форма `convOut` такая же, как у исходного изображения, тогда как при шаге 2 оно примерно вдвое меньше по высоте и ширине. Таким образом, ответ на первый вопрос — “нет”. Это также не тот случай, когда они имеют одинаковый набор значений, поскольку, например, если бы у нас было два фрагмента небольших значений рядом друг с другом, но окруженных большими значениями, вывод с однопиксельным шагом включал бы большее из этих двух небольших значений, тогда как в двухпиксельной строке эти значения будут “утоплены” соседними значениями пикселей. Наконец, третий ответ — “да” в том смысле, что каждое значение пула в первом случае повторяется во втором случае, но не наоборот.

Упражнение 3.6, а. Каждое ядро, которое мы создаем, имеет форму [2, 2, 3], что подразумевает 12 переменных на ядро. Так как мы создаем 10 из них, создаются 120 переменных. Поскольку количество случаев применения ядра не имеет никакого отношения к его размеру/форме, ни размер партии (100), ни высота/ширина (8/8) не влияют на этот ответ.

Глава 4

Упражнение 4.2. Важное различие между установкой E равным 0 (или всем 1) заключается в том, что NN никогда не видит фактический ввод, видит только его векторное представление. Таким образом, установка всех векторных представлений равными одному и тому же значению приводит к тому, что все слова идентичны. Очевидно, что это устраивает любой шанс чему-нибудь обучить.

Упражнение 4.3. На самом деле вычисление общих потерь при использовании регуляризации L_2 требует вычисления суммы квадратов весов для всех весов в модели. Эта величина не нужна где-либо еще в графе вычислений. Например, чтобы вычислить производную от общих потерь по $w_{i,j}$ требуется только добавление $w_{i,j}$ к обычным потерям.

Упражнение 4.5. Да, это может быть лучше, чем выбор из равномерного распределения. Следует научиться назначать более высокую вероятность более распространенным словам (например, “the”). Тем не менее обратите внимание, что униграммная модель не имеет “входных данных”, а следовательно, не нуждается в векторных представлениях или линейных блоках входных весов. Однако ей нужны смещения, так как, изменяя их, происходит обучение назначению вероятностей в соответствии с частотой слов.

Глава 5

Упражнение 5.2. Более сложный механизм внимания использует состояние декодера после времени t , чтобы определить внимание, используемое при декодировании в момент времени $t+1$. Но мы, очевидно, не узнаем это значение, пока не обработаем время t . При обратном распространении во времени мы обрабатываем окно слов одновременно. За исключением первой позиции в окне декодера, у нас нет значения исходного состояния, которое нам нужно для вычисления. По сути, нам нужно написать новый механизм обратного распространения во времени.

Упражнение 5.4. а) Мы хотим хороший машинный перевод. В той степени, в которой другая функция потерь влияет на вещи, она, по-видимому, заключается в перемещении весов, поэтому они справляются с этой задачей хуже. Таким образом, производительность ухудшается. б) Но часть а) верна только для обучающих данных. В других примерах ничего не говорится о производительности. Добавление второй функции потерь должно помочь программе узнать больше о структуре французского языка. Это должно улучшить его производительность на новых примерах, которые составляют основную часть тестовых данных.

Глава 6

Упражнение 6.3. В итерации по значению единственный способ, которым состояние может получить ненулевое значение, — если а) для действия есть способ получить немедленное вознаграждение из этого состояния (только состояние 14) или б) можно выполнить действие из этого состояния, которое приводит к состоянию, которое уже имеет ненулевое значение (состояния 10, 13 и 14). Таким образом, все другие состояния должны сохранять свои нулевые значения. Состояние 10 имеет максимальное значение Q для $Q(14, l)$, $Q(14, d)$ и $Q(14, r)$. В каждом случае значения заканчиваются в состоянии 15 и составляют $0,33 \times 0,9 \times 0,33 = 0,1$, поэтому $V(10) = 0,1$. Точно так же $V(14) = 0,1$. Наконец $V(15)$ получает свое значение от $Q(15, d)$ или $Q(15, r)$. В обоих случаях вычисление составляет $0,33 \times 0,9 \times 0,1 + 0,33 \times 0,9 \times 0,33 + 0,33 \times 1 = 0,03 + 0,1 + 0,33 = 0,47$.

Упражнение 6.4. Мы заявили, что, вычислив возможные действия в REINFORCE и сохранив их вероятности, мы сможем вычислить потери на втором проходе, не начиная с состояния, а затем вычисляя вероятности. Однако, если бы мы сделали это и *не* переделали вычисления, ведущие из состояния к вероятностям действий, обратный проход TF не смог бы проследить вычисления обратно через полно связанные уровни, которые вычисляют действия из состояния. Таким образом, этот слой (или слои) не будут обновлять свои значения, и программа не научится вычислять лучшие рекомендации действий для состояний.

Глава 7

Упражнение 7.1. Чтобы игнорировать граничные значения, очевидная вещь, которую нужно сделать, — установить равными нулю значения, которые соединяют пиксели с первым слоем. Пусть x будут значениями пикселей в одномерной версии изображения $0 < x < 783$. Если i, j располагаются в диапазоне пикселей на изображении $28*28$, то для всех x , таких как $x = j + 28i$, $i < 2$ или $i > 25$ и $j < 2$ или $j > 25$ для всех y , $0 \leq y \leq 256$:

$$E_1(x, y) = 0$$

Упражнение 7.3.

```
tf.nn.conv2d_transpose(smallerI, feat, [1, 8, 8, 1], [1, 2, 2, 1], "SAME")
```

Упражнение 7.5. Мы хотели, чтобы отслеживаемые значения выводились после последней итерации. Если бы мы установили диапазон 5000, последней была бы итерация 4999 и она не была бы выведена.

Предметный указатель

A

Activation function, 49
Actor, 138
Actor-critic method, 138
Adam optimizer, 136
AE, 145
Agent, 121
Aligned corpus, 103
Attention, 110
Autoencoder, 145
Axon, 14

B

Back propagation, 24
 through time, 92
Backward pass, 24
Batch size, 27
Bias, 15
Bigram model, 81
Binary classification, 13
Broadcasting, 33

C

Cart pole, 131
Ceiling function, 65
Cell state, 97
Channel, 66
Chatbot, 119
Classification, 13
Convolutional
 filter, 62

kernel, 62
neural network, 61
Convolution function, 66
Convolve, 62
Corpora, 16
Cosine similarity, 83
Co-training, 146
CPU, 32
Critic, 138
Cross-entropy, 23
 loss, 21

D

Data normalization, 28
Decoder, 146
Decoding pass, 105
Deep
 learning, 11
 reinforcement learning, 34
 RL, 34
Deep-Q learning, 128
Dendrite, 14
De-noising autoencoder, 149
Development set, 16
Dimensionality reduction, 145
Discount, 121
Discounted future reward, 121
Discretize, 11
Discriminator, 161
Dot product, 15
Downsampling, 146
Dropout, 88

E

Early stopping, 88
 Embedding layer, 82
 Emphasis, 111
 Encoder, 146
 Encoding pass, 105
 Environment, 121
 Episode, 137
 Epoch, 18
 Epsilon-decreasing strategy, 126
 Epsilon-greedy strategy, 126
 Expected value, 55
 Experience replay, 140
 Exploration-exploitation tradeoff, 127

F

Feature, 13
 Feed-forward neural network, 19
 Floor function, 65; 94
 Forgetting signal, 98
 Forward pass, 24
 Frozen-lake problem, 123
 Fully
 connected, 57
 Fully supervised learning, 13; 145
 Function approximation, 16
 problem, 128

G

GAN, 145; 161
 Gated Recurrent Unit, 104
 Gaussian distribution, 42
 Generative Adversarial Network, 145;
 161
 Generator, 161
 Given new distinction, 110

GPU, 32

Gradient, 24
 descent, 20; 27
 operator, 31
 Graphics Processing Unit, 32
 Greedy algorithm, 127
 GRU, 104

H

Held-out set, 16
 Heuristic, 13
 Hidden size, 51
 Hyperbolic tangent, 99
 Hyperparameter, 16

I

Image loss, 155
 Instability, 34

K

Kernel function, 62
 Kullback-Leibler divergence, 160

L

L2 regularization, 88
 Label, 13
 Language model, 79
 Layer, 19
 Leaky relu, 50
 Learning rate, 20
 LGU, 102
 Light intensity, 11
 Linear
 Gated Unit, 102
 unit, 15
 Logit, 22

Long Short-Term Memory, 97

Loss function, 20

LSTM, 97

M

Machine learning, 13

Machine Translation, 103

Markov

assumption, 121

decision process, 121

Matrix, 30

MDP, 121

Model-free learning, 125

Momentum, 136

MT, 103

Multilevel perceptron, 20

N

Neural net, 11

Noise, 149

O

One-hot vector, 44

P

Padding, 64; 69

Parameter, 15

Parameterized class, 16

Patch, 63

Penn Treebank Corpus, 80

Perceptron, 14

algorithm, 15

Perplexity, 86

Pixel, 11

value, 11

Placeholder, 40

Plus unit, 97

Policy, 122

gradient, 131

Position-only attention, 110

Pre-training, 146

Probability distribution, 21

PTB, 80

Q

Q function, 122

Q-learning, 127

Quadratic loss, 37; 130

Q-обучение, 127

Q-функция, 122

R

Rectified Linear Unit, 49

Recurrent Neural Network, 90

REINFORCE, 134

Reinforcement Learning, 121

ReLU, 49

с утечкой, 50

Reward, 121

RL, 121

RNN, 90

S

Saver object, 51

Semi-supervised, 13

Sentence

embedding, 105

padding, 81

Seq2seq, 103

Sequence-to-sequence learning, 103

Sigmoid function, 49

Signal, 149

Skip-gram model, 100
 Sparse matrix, 85
 Squared-error loss, 130
 Standard deviation, 42
 State similarity, 117
 Step function, 21
 Stochastic gradient descent, 24; 27
 Stride, 64; 69
 Supervised learning, 13; 145

T

Tanh activation function, 99
 Temporal difference error, 130
 Tensor, 41
 Tensordot, 53
 Test set, 16
 Time unit, 97; 104
 Tokenization, 80; 104
 Training example, 15
 Training set, 16
 Trigram model, 86

U

Universal loss function, 161
 Unsupervised learning, 13; 145
 Upsampling, 146

V

VAE, 153
 Validation set, 16
 Value
 function, 122
 iteration, 122
 Vanishing gradient, 50
 Variational
 autoencoder, 147; 153

loss, 156
 method, 153

W

Weakly supervised learning, 145
 Weight, 15
 Window size, 92
 Word embedding, 81

X

Xavier initialization, 55

Z

Zero-one loss, 20

A

Автокодировщик, 145
 Агент, 121
 Аксон, 14
 Актор, 138
 Акцент, 111
 Алгоритм перцептрона, 15
 Аппроксимация функций, 16

Б

Бесшумный автокодировщик, 149
 Блок
 времени, 97; 104
 линейной ректификации, 49
 плюс, 97

В

Вариационные потери, 156
 Вариационный
 автокодер, 147
 автокодировщик, 153
 метод, 153

Векторное представление предложений, 105
слова, 81
Вес, 15
Внимание, 110
только к позиции, 110
Вознаграждение, 121
Воспроизведение опыта, 140
Выбранный набор, 16
Выровненный корпус, 103

Г

Гауссово распределение, 42
Генеративно-состязательная сеть, 145; 161
Генератор, 161
Гиперболический тангенс, 99
Гиперпараметр, 16
Глубокое
Q-обучение, 128
обучение, 11
с подкреплением, 34
Градиент, 24
политики, 131
Градиентный спуск, 20; 27
Графический процессор, 32

Д

Данное новое различие, 110
Двоичная
классификация, 13
функция потерь, 20
Декодер, 146
Дендрит, 14
Дисконт, 121

Дисконтируемое будущее
вознаграждение, 121
Дискриминатор, 161
Долгая краткосрочная память, 97
Дополнение, 64; 69
предложений, 81

Ж

Жадный алгоритм, 127

З

Задача
аппроксимации функций, 128
замерзшего озера, 123
Заполнитель, 40
Значение пикселя, 11

И

Импульс, 136
Инициализация Ксавье, 55
Интенсивность света, 11
Исключение, 88
Исчезающий градиент, 50
Итерация по значению, 122; 124

К

Канал, 66
Квадратичная
потеря, 37; 130
функция потерь, 130
Классификация, 13
Кодер, 146
Компромисс между исследованием
и использованием, 127
Корпус, 16

Косинусное сходство, 83

Критик, 138

Кросс-энтропия, 23

Л

Лексический анализ, 80; 104

Линейный

блок, 15

выпрямитель, 49

управляемый блок, 102

Логит, 22

М

Марковский процесс принятия решений, 121

Марковское

предположение, 121

свойство, 121

Матрица, 30

Машинное обучение, 13

Машинный перевод, 103

Метка, 13

Метод актор-критик, 138

Многослойный перцепtron, 20

Модель

skip-gram, 100

биграмм, 81

триграмм, 86

Н

Набор разработки, 16

Нейрон, 14

Нейронная сеть, 11

с прямой связью, 19

Нестабильность, 34

Нормализация данных, 28

Нормальное распределение, 42

О

Обратное распространение, 24

во времени, 92

Обратный проход, 24

Обучающий

набор, 16

пример, 15

Обучение

без модели, 125

без учителя, 13; 145

от последовательности к

последовательности, 103

со слабым привлечением

учителя, 145

с подкреплением, 121

с полным привлечением

учителя, 145

с учителем, 13; 145

с частичным привлечением

учителя, 13

Ожидаемое значение, 55

Оператор градиента, 31

Оптимизатор Adam, 136

Оцифровка, 11

Ошибка временных различий, 130

П

Параметр, 15

Параметризованный класс, 16

Передискретизация, 146

Перплексивность, 86

Перцептрон, 14

Пиксель, 11
 Политика, 122
 Полное обучение с учителем, 13
 Полносвязный слой, 57
 Потери изображения, 155
 Потеря, 22
 Пошаговая функция, 21
 Предобучение, 146
 Предположение iid, 33
 Признак, 13
 Проверочный набор, 16
 Проход
 декодирования, 105
 кодирования, 105
 Прямой проход, 24

Р

Размер
 окна, 92
 партии, 27
 Разреженная матрица, 85
 Ранняя остановка, 88
 Распределение вероятностей, 21
 Расстояние Кульбака–Лейблера, 160
 Регуляризация, 88
 L2, 88
 Рекуррентная нейронная сеть, 90

С

Свертка, 62
 Сверточная нейронная сеть, 61
 Сверточное ядро, 62
 Сверточный фильтр, 62
 Сигмовидная функция, 49
 Сигнал, 149
 забывания, 98

Скалярное произведение, 15
 Скорость обучения, 20
 Скрытый размер, 51
 Слой, 19
 векторного представления, 82
 Смещение, 15
 Снижение размерности, 145
 Сообучение, 146
 Состояние ячейки, 97
 Среда, 121
 Стандартное отклонение, 42
 Стохастический градиентный
 спуск, 24; 27
 Стратегия снижения эпсилона, 126
 Субдискретизация, 146
 Сходства состояний, 117

Т

Тележка с шестом, 131
 Тензор, 41
 Тестовый набор, 16
 Универсальная функция потерь, 161

У

Унитарный вектор, 44
 Управляемый рекуррентный
 блок, 104

Ф

Фрагмент, 63
 Функция
 argmax, 46
 softmax, 21
 активации, 49
 активации tanh, 99
 значения, 122

кросс-энтропийных потерь, 21
поля, 65; 94
потерь, 20; 22
потолка, 65
свертки, 66
ядра, 62

Х

Хранимые объекты, 51

Ц

Центральный процессор, 32
Циклический граф, 90

Ч

Чат-бот, 119

Ш

Шаг, 64; 69
Широковещание, 33
Шум, 149

Э

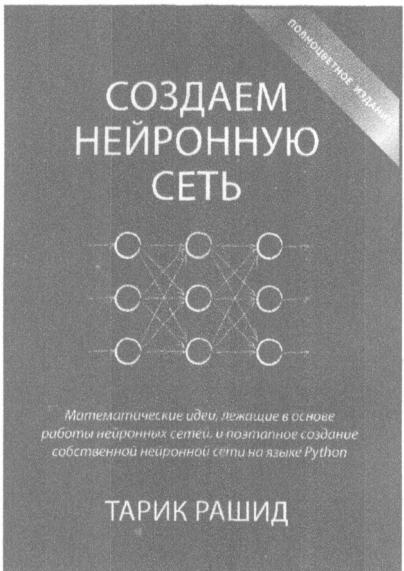
Эвристика, 13
Эпизод, 137
Эпоха, 18
Эпсилон-жадная стратегия, 126

Я

Язык Tensorflow, 39
Языковая модель, 79

СОЗДАЕМ НЕЙРОННУЮ СЕТЬ

Тарик Рашид



Эта книга представляет собой введение в теорию и практику создания нейронных сетей. Она предназначена для тех, кто хочет узнать, что такое нейронные сети, где они применяются и как самому создать такую сеть, не имея опыта работы в данной области. Изложение материала сопровождается подробным описанием процедуры поэтапного создания полностью функционального кода, который реализует нейронную сеть на языке Python и способен выполняться даже на таком миниатюрном компьютере, как Raspberry Pi Zero.

Основные темы книги:

- нейронные сети и системы искусственного интеллекта;
- структура нейронных сетей;
- сглаживание сигналов, распространяющихся по нейронной сети, с помощью функции активации;
- тренировка и тестирование нейронных сетей;
- интерактивная среда программирования IPython;
- распознавание образов с помощью нейронных сетей.

www.williamspublishing.com

ISBN 978-5-9909445-7-2 в продаже

ГЛУБОКОЕ ОБУЧЕНИЕ

ГОТОВЫЕ РЕШЕНИЯ

Давид Осинга



www.williamspublishing.com

ISBN: 978-5-907144-50-7

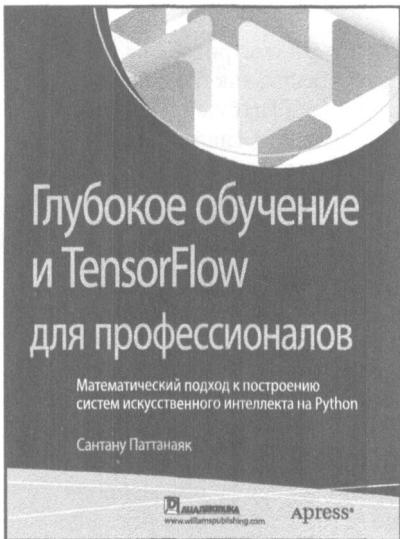
Благодаря готовым примерам, приведенным в книге, вы научитесь решать задачи, связанные с классификацией и генерированием текста, изображений и музыки. В каждой главе описывается несколько решений, объединяемых в единый проект, например приложение, реализующее тренировку музыкальной рекомендательной системы. Также имеется глава с описанием методик, которые в случае необходимости помогут выполнить отладку нейронной сети. Основные темы книги:

- использование векторных представлений слов для вычисления схожести текстов;
- построение рекомендательной системы фильмов на основе ссылок в Википедии;
- визуализация внутренних состояний нейронной сети;
- создание модели, рекомендующей эмодзи для фрагментов текста;
- повторное использование предварительно обученных сетей для создания службы обратного поиска изображений;
- генерирование пиктограмм с помощью генеративно-состязательных сетей (GAN), автокодировщиков и рекуррентных сетей (RNN);
- распознавание музыкальных жанров и индексирование коллекций песен.

в продаже

ГЛУБОКОЕ ОБУЧЕНИЕ И TENSORFLOW ДЛЯ ПРОФЕССИОНАЛОВ

Сантану Паттанаяк



www.williamspublishing.com

Данная книга представляет собой углубленное практическое руководство, которое позволит читателям освоить методы глубокого обучения на уровне, достаточном для развертывания готовых решений. Прочитав книгу, вы сможете быстро приступить к работе с библиотекой TensorFlow и заняться оптимизацией архитектур глубокого обучения. Весь программный код доступен в виде блокнотов iPython и сценариев, позволяющих с легкостью воспроизводить примеры и экспериментировать с ними. Благодаря этой книге вы:

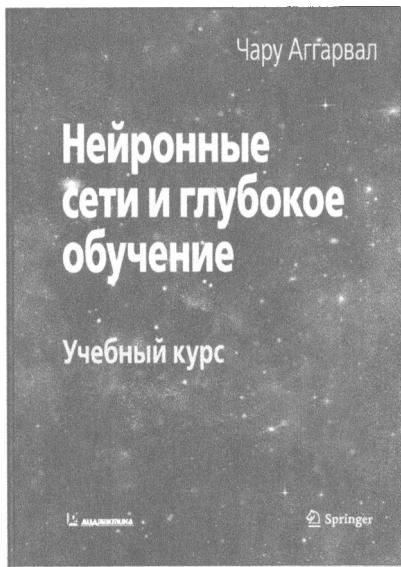
- овладеете полным стеком технологий глубокого обучения с использованием TensorFlow и получите необходимую для этого математическую подготовку;
- научитесь развертывать сложные приложения глубокого обучения в производственной среде с помощью TensorFlow;
- сможете проводить исследования в области глубокого обучения и выполнять самостоятельные эксперименты в TensorFlow.

ISBN: 978-5-907144-25-5

в продаже

НЕЙРОННЫЕ СЕТИ И ГЛУБОКОЕ ОБУЧЕНИЕ УЧЕБНЫЙ КУРС

Чару Аггарвал



www.dialektika.com

В книге рассматриваются как классические, так и современные модели глубокого обучения. В первых двух главах основной упор сделан на понимании взаимосвязи традиционного машинного обучения и нейронных сетей. Главы 3 и 4 посвящены подробному обсуждению процессов тренировки и регуляризации нейронных сетей. В главах 5 и 6 рассмотрены сети радиально-базисных функций (RBF) и ограниченные машины Больцмана. В главах 7 и 8 обсуждаются рекуррентные и сверточные нейронные сети. Главы 9 и 10 посвящены более сложным темам, таким как глубокое обучение с подкреплением, нейронные машины Тьюринга, самоорганизующиеся карты Кохонена и генеративно-состязательные сети. Книга предназначена для студентов старших курсов, исследователей и специалистов-практиков.

ISBN: 978-5-907203-01-3 в продаже

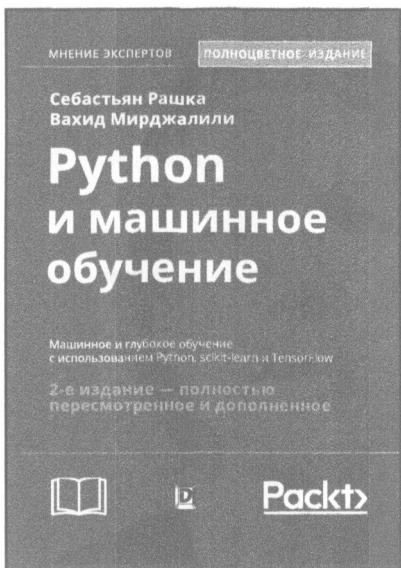
PYTHON И МАШИННОЕ ОБУЧЕНИЕ

МАШИННОЕ И ГЛУБОКОЕ ОБУЧЕНИЕ

С ИСПОЛЬЗОВАНИЕМ PYTHON, SCIKIT-LEARN И

TENSORFLOW, 2-Е ИЗДАНИЕ

**Себастьян Рашка
и Вахид Мирджалили**



www.dialektika.com

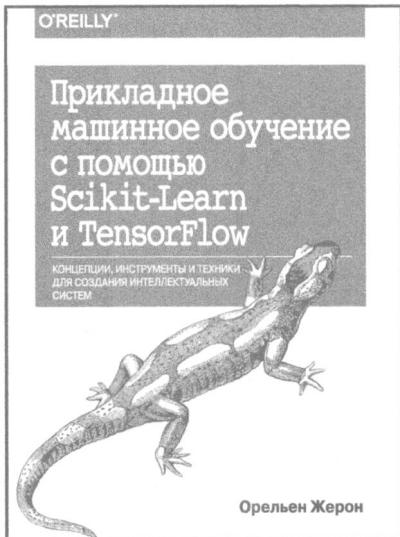
Машинное обучение поглощает мир программного обеспечения, и теперь глубокое обучение расширяет машинное обучение. Освойте и работайте с передовыми технологиями машинного обучения, нейронных сетей и глубокого обучения с помощью 2-го издания бестселлера Себастьяна Рашки. Будучи основательно обновленной с учетом самых последних библиотек Python с открытым кодом, эта книга предлагает практические знания и приемы, которые необходимы для создания и содействия машинному обучению, глубокому обучению и современному анализу данных. Если вы читали 1-е издание книги, то вам доставит удовольствие найти новый баланс классических идей и современных знаний в машинном обучении. Каждая глава была серьезно обновлена, и появились новые главы по ключевым технологиям. У вас будет возможность изучить и поработать с TensorFlow более вдумчиво, нежели ранее, а также получить важнейший охват библиотеки для нейронных сетей Keras наряду с самыми свежими обновлениями библиотеки scikit-learn.

ISBN 978-5-907114-52-4

в продаже

ПРИКЛАДНОЕ МАШИННОЕ ОБУЧЕНИЕ С ПОМОЩЬЮ SCIKIT-LEARN И TENSORFLOW

Орельен Жерон



www.dialektika.com

Благодаря серии недавних достижений глубокое обучение значительно усилило всю область машинного обучения. В наше время даже программисты, почти ничего не знающие об этой технологии, могут использовать простые и эффективные инструменты для реализации программ, которые способны обучаться на основе данных. В настоящем практическом руководстве показано, что и как следует делать.

За счет применения конкретных примеров, минимума теории и двух фреймворков Python производственного уровня — Scikit-Learn и TensorFlow — автор книги Орельен Жерон поможет вам получить интуитивное представление о концепциях и инструментах, предназначенных для построения интеллектуальных систем. Вы узнаете о ряде приемов, начав с простой линейной регрессии и постепенно добравшись до глубоких нейронных сетей. Учитывая наличие в каждой главе упражнений, призванных закрепить то, чему вы научились, для начала работы нужен лишь опыт программирования.

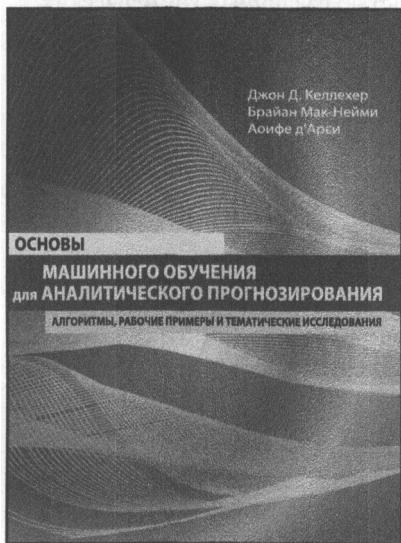
ISBN 978-5-9500296-2-2

в продаже

ОСНОВЫ МАШИННОГО ОБУЧЕНИЯ ДЛЯ АНАЛИТИЧЕСКОГО ПРОГНОЗИРОВАНИЯ

АЛГОРИТМЫ, РАБОЧИЕ ПРИМЕРЫ И ТЕМАТИЧЕСКИЕ ИССЛЕДОВАНИЯ

**Джон Д. Келлехер
Брайан Мак-Нейми
Аоифе д'Арси**



www.williamspublishing.com

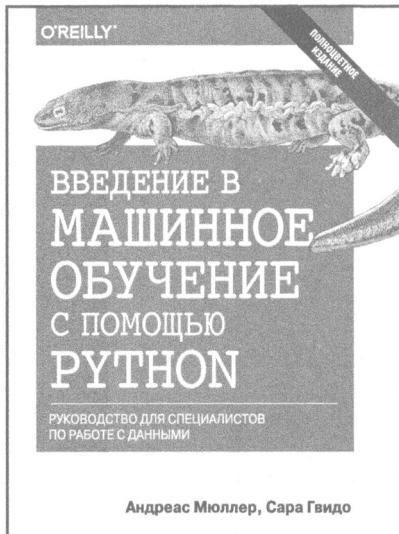
Книга представляет собой учебник по машинному обучению с акцентом на коммерческие приложения. Она предлагает подробное описание наиболее важных подходов к машинному обучению, используемых в интеллектуальном анализе данных, охватывающих как теоретические концепции, так и практические приложения. Формальный математический материал дополняется пояснительными примерами, а примеры исследований иллюстрируют применение этих моделей в более широком контексте бизнеса. В книге рассмотрены информационное обучение, обучение на основе сходства, вероятностное обучение и обучение на основе ошибок. Описанию каждого из этих подходов предшествует объяснение основополагающей концепции, за которой следуют математические модели и алгоритмы, иллюстрированные подробными рабочими примерами.

ISBN 978-5-6040044-9-4

в продаже

ВВЕДЕНИЕ В МАШИННОЕ ОБУЧЕНИЕ С ПОМОЩЬЮ PYTHON РУКОВОДСТВО ДЛЯ СПЕЦИАЛИСТОВ ПО РАБОТЕ С ДАННЫМИ

Андреас Мюллер
Сара Гвидо



www.williamspublishing.com

ISBN 978-5-9908910-8-1

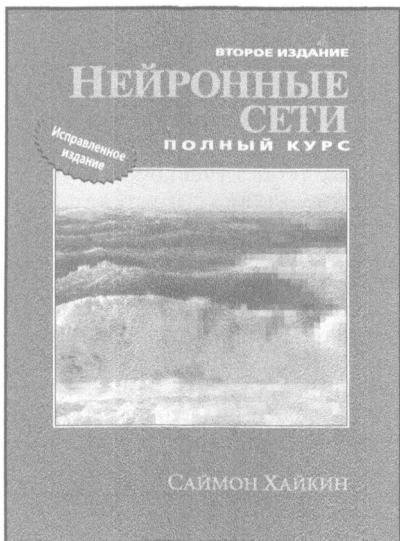
Эта полноцветная книга — отличный источник информации для каждого, кто собирается использовать машинное обучение на практике. Ныне машинное обучение стало неотъемлемой частью различных коммерческих и исследовательских проектов, и не следует думать, что эта область — прерогатива исключительно крупных компаний с мощными командами аналитиков. Эта книга научит вас практическим способам построения систем МО, даже если вы еще новичок в этой области. В ней подробно объясняются все этапы, необходимые для создания успешного проекта машинного обучения, с использованием языка Python и библиотек scikit-learn, NumPy и matplotlib. Авторы сосредоточили свое внимание исключительно на практических аспектах применения алгоритмов машинного обучения, оставив за рамками книги их математическое обоснование. Данная книга адресована специалистам, решающим реальные задачи, а поскольку область применения методов МО практически безгранична, прочитав эту книгу, вы сможете собственными силами построить действующую систему машинного обучения в любой научной или коммерческой сфере.

в продаже

НЕЙРОННЫЕ СЕТИ: ПОЛНЫЙ КУРС

2-Е ИЗДАНИЕ

Саймон Хайкин



www.dialektika.com

В книге рассматриваются основные парадигмы искусственных нейронных сетей. Представленный материал содержит строгое математическое обоснование всех нейросетевых парадигм, иллюстрируется примерами, описанием компьютерных экспериментов, содержит множество практических задач, а также обширную библиографию.

В книге также анализируется роль нейронных сетей при решении задач распознавания образов, управления и обработки сигналов. Структура книги очень удобна для разработки курсов обучения нейронным сетям и интеллектуальным вычислениям.

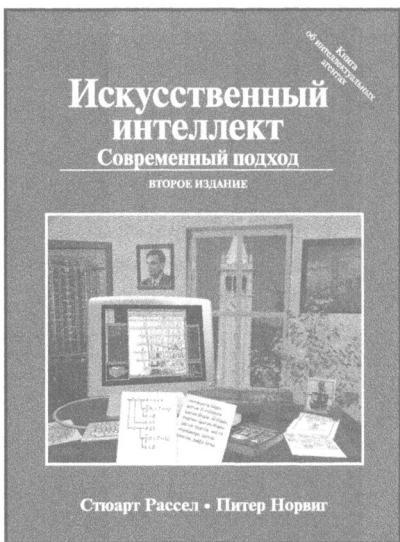
Книга будет полезна для инженеров, специалистов в области компьютерных наук, физиков и специалистов в других областях, а также для всех тех, кто интересуется искусственными нейронными сетями.

ISBN 978-5-907144-22-4

в продаже

ИСКУССТВЕННЫЙ ИНТЕЛЛЕКТ СОВРЕМЕННЫЙ ПОДХОД ВТОРОЕ ИЗДАНИЕ

**С. Рассел,
П. Норвиг**



В книге представлены все современные достижения и изложены идеи, которые были сформулированы в исследованиях, проводившихся в течение последних пятидесяти лет, а также собраны на протяжении двух тысячелетий в областях знаний, ставших стимулом к развитию искусственного интеллекта как науки проектирования рациональных агентов. Теоретическое описание иллюстрируется многочисленными алгоритмами, реализации которых в виде готовых программ на нескольких языках программирования находятся на сопровождающем веб-сайте. Книга предназначена для использования в базовом университете курсе или в последовательности курсов по специальности. Применима в качестве основного справочника для аспирантов, специализирующихся в области искусственного интеллекта, а также будет небезынтересна профессионалам, желающим выйти за пределы избранной ими специальности. Благодаря кристальной ясности и наглядности изложения вполне может быть отнесена к лучшим образцам научно-популярной литературы.

www.dialektika.com

ISBN 978-5-907114-65-4

в продаже

ВВЕДЕНИЕ В ГЛУБОКОЕ ОБУЧЕНИЕ

ЕВГЕНИЙ ЧЕРНЯК

Автор книги – давний исследователь искусственного интеллекта, специализирующийся на обработке естественного языка, революцию в котором сделало глубокое обучение. К сожалению, ему потребовалось много времени, чтобы это понять. Можно сказать, что нейронные сети угрожают революцией уже третий раз, а отнюдь не первый. Тем не менее автор внезапно оказался далеко позади и изо всех сил пытался наверстать упущенное. Именно поэтому он сделал то, что сделал бы на его месте любой уважающий себя профессор: запланировал преподавание курса и начал ускоренно изучать материал, просматривая веб-страницы.

Этим объясняется несколько выдающихся особенностей этой книги. Во-первых, краткость. Во-вторых, она сильно зависит от проекта. Автор считает, что материал по информатике лучше изучать при написании программ, поэтому книга во многом отражает его привычки в преподавании.

Эта книга, в первую очередь, задумана как учебник для курса по глубокому обучению. Курс, который автор преподает в Брауне, предназначен как для выпускников, так и для остальных студентов, и охватывает весь материал. Хотя фактическое количество материала по линейной алгебре не так уж велико, студенты сказали, что без него им было бы довольно сложно разобраться в многослойных сетях и необходимых им тензорах. И наконец, есть предпосылка для вероятности и статистики. Автор также предполагает элементарные знания читателей в программировании на языке Python. Хотя этот материал не включен в книгу, но у автора есть дополнительная “лаборатория” по основам языка Python.

Изображение на обложке:
©Depositphotos.com/5855636
Автор: agsandrew

Категория: искусственный интеллект/нейронные сети
Уровень: для пользователей средней и высокой квалификации



www.dialektika.com

The MIT Press
Cambridge, Massachusetts
London, England

ISBN 978-5-907203-10-5

19184

9 785907 203105