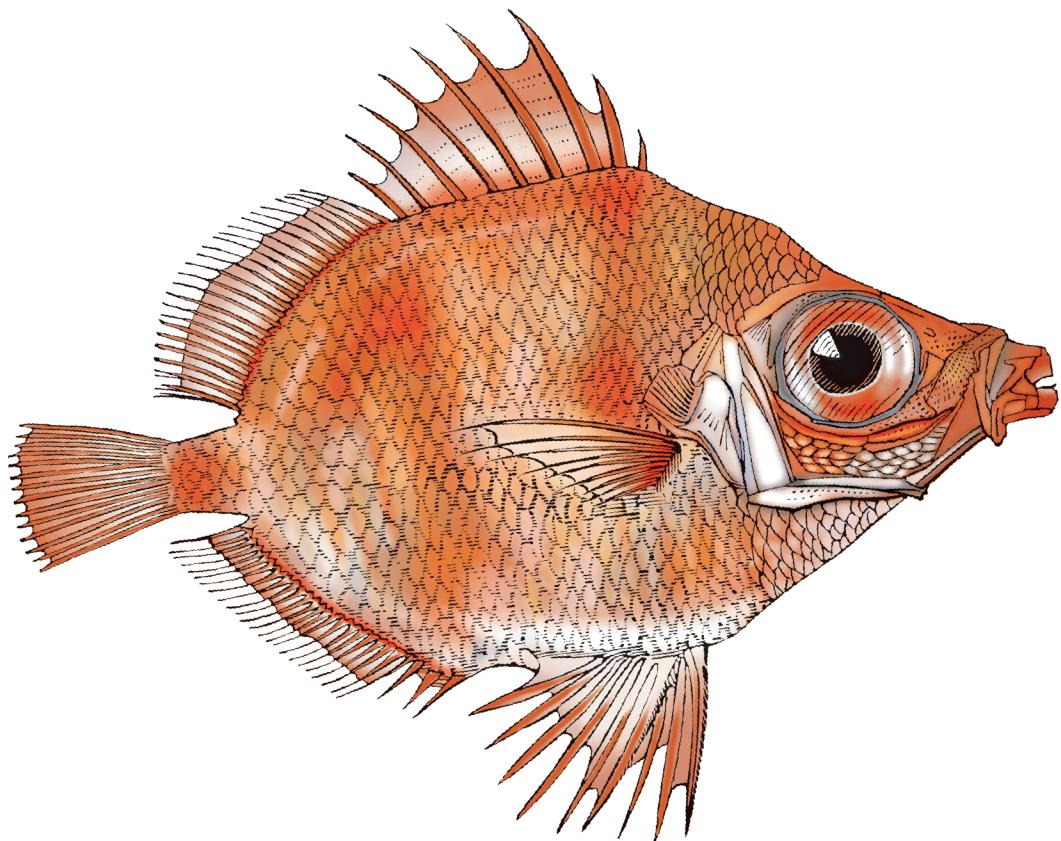


O'REILLY®

Глубокое обучение с fastai и PyTorch

Минимум формул, минимум кода,
максимум эффективности



Джереми Ховард и Сильвейн Гуггер

Deep Learning for Coders with fastai and PyTorch

AI Applications Without a PhD

Jeremy Howard and Sylvain Gugger

Beijing • Boston • Farnham • Sebastopol • Tokyo

O'REILLY®

Глубокое обучение с fastai и PyTorch

Минимум формул, минимум кода,
максимум эффективности

Джереми Ховард и Сильвейн Гуггер



Санкт-Петербург · Москва · Минск

2022

ББК 32.813
УДК 004.8
Х68

Ховард Джереми, Гуггер Сильвейн

Х68 Глубокое обучение с fastai и PyTorch: минимум формул, минимум кода, максимум эффективности. — СПб.: Питер, 2022. — 624 с.: ил. — (Серия «Бестселлеры O'Reilly»).
ISBN 978-5-4461-1475-7

Обычно на глубокое обучение смотрят с ужасом, считая, что только доктор математических наук или ботан, работающий в крутой айтишной корпорации, могут разобраться в этой теме. Отбросьте стереотипы: любой программист, знакомый с Python, может добиться впечатляющих результатов. Как? С помощью fastai — библиотеки, предоставляющей комфортный интерфейс для решения наиболее популярных задач.

Создатели fastai доказали, что самые модные и актуальные приложения можно делать быстро и не засыпать над скучными теоретическими выкладками и зубодробительными формулами.

16+ (В соответствии с Федеральным законом от 29 декабря 2010 г. № 436-ФЗ.)

ББК 32.813
УДК 004.8

Права на издание получены по соглашению с O'Reilly. Все права защищены. Никакая часть данной книги не может быть воспроизведена в какой бы то ни было форме без письменного разрешения владельцев авторских прав.

Информация, содержащаяся в данной книге, получена из источников, рассматриваемых издательством как надежные. Тем не менее, имея в виду возможные человеческие или технические ошибки, издательство не может гарантировать абсолютную точность и полноту приводимых сведений и не несет ответственности за возможные ошибки, связанные с использованием книги. Издательство не несет ответственности за доступность материалов, ссылки на которые вы можете найти в этой книге. На момент подготовки книги к изданию все ссылки на интернет-ресурсы были действующими.

ISBN 978-1492045526 англ.

© Authorized Russian translation of the English edition of Deep Learning for Coders with fastai and PyTorch

ISBN 9781492045526 © 2020 Sylvain Gugger and Jeremy Howard.

This translation is published and sold by permission of O'Reilly Media, Inc., which owns or controls all rights to publish and sell the same.

ISBN 978-5-4461-1475-7

© Перевод на русский язык ООО «Прогресс книга», 2022

© Издание на русском языке, оформление ООО «Прогресс книга», 2022

© Серия «Бестселлеры O'Reilly», 2022

Краткое содержание

Отзывы о книге.....	18
Введение	22
Предисловие.....	25

ЧАСТЬ I ГЛУБОКОЕ ОБУЧЕНИЕ НА ПРАКТИКЕ

Глава 1. Путешествие в мир глубокого обучения.....	28
Глава 2. От модели к продакшну.....	86
Глава 3. Этика данных.....	125

ЧАСТЬ II ПОНИМАНИЕ ПРИЛОЖЕНИЙ НА БАЗЕ FASTAI

Глава 4. Обучение классификатора цифр: взгляд изнутри	168
Глава 5. Классификация изображений.....	222
Глава 6. Другие задачи компьютерного зрения.....	256
Глава 7. Обучение современной модели	278
Глава 8. Коллаборативная фильтрация.....	294
Глава 9. Табличное моделирование	320
Глава 10. Погружение в NLP: рекуррентные нейронные сети	376
Глава 11. Преобразование данных с помощью Mid-Level API	405

ЧАСТЬ III ОСНОВЫ ГЛУБОКОГО ОБУЧЕНИЯ

Глава 12. Языковая модель с нуля	422
Глава 13. Сверточные нейронные сети.....	453
Глава 14. ResNet	492
Глава 15. Архитектуры приложений	511
Глава 16. Процесс обучения.....	524

ЧАСТЬ IV ГЛУБОКОЕ ОБУЧЕНИЕ С ЧИСТОГО ЛИСТА

Глава 17. Продвинутые основы нейронной сети	544
Глава 18. Интерпретация CNN с помощью CAM	571
Глава 19. Класс Learner с нуля.....	579
Глава 20. Подведем итог	598

ПРИЛОЖЕНИЯ

Приложение А. Создание блога.....	602
Приложение Б. Схема подготовки проекта по аналитике данных.....	611
Об авторах	619
Благодарности	620
Об обложке	622

Оглавление

Отзывы о книге	18
Введение.....	22
Для кого эта книга	22
Что нужно знать	23
Чему вы научитесь	23
Предисловие	25
От издательства.....	26

ЧАСТЬ I ГЛУБОКОЕ ОБУЧЕНИЕ НА ПРАКТИКЕ

Глава 1. Путешествие в мир глубокого обучения.....	28
Глубокое обучение для всех	28
Нейронные сети: краткая история	30
Кто мы	33
Как изучать глубокое обучение	35
Ваши проекты и мышление	37
ПО: PyTorch, fastai и Jupyter (почему это важно).....	38
Ваша первая модель	40
Настройка сервера глубокого обучения на GPU	40
Запуск первого блокнота	42
Что такое машинное обучение.....	46
Что такое нейронная сеть.....	50
Немного терминологии глубокого обучения.....	51
Характерные для ML ограничения.....	52
Как работает наш распознаватель изображений	54
Чему научился наш распознаватель изображений	61

Распознаватели изображений для других задач	64
Обобщение терминов	68
Глубокое обучение подходит не только для классификации изображений	70
Контрольные и тестовые выборки.....	77
Создавайте тестовую выборку обдуманно	79
Момент выбора собственного приключения	83
Вопросник.....	83
Дополнительные задания	85
Глава 2. От модели к продакшену	86
Практика глубокого обучения.....	86
Начало проекта	87
Текущий уровень глубокого обучения	89
Подход Drivetrain	93
Сбор данных	95
От данных к DataLoaders	100
Аугментация данных	105
Обучение модели и ее использование для чистки данных.....	105
Превращение модели в онлайн-приложение.....	109
Использование модели для вывода	109
Создание в блокноте приложения на основе модели	110
Превращение блокнота в реальное приложение.....	113
Развертывание приложения	114
Как избежать катастрофы.....	117
Непредвиденные последствия и петли обратной связи.....	120
Записывайте!.....	121
Вопросник.....	122
Дополнительные задания	124
Глава 3. Этика данных	125
Ключевые примеры этики данных	126
Баги и оказание помощи: неисправный алгоритм, распределявший медицинские льготы	127
Петли обратной связи: рекомендательная система YouTube.....	127
Предвзятость: «Арест» профессора Латаньи Суини.....	128
Почему это важно?	129

Тесное взаимодействие процессов ML и дизайна продукта.....	132
Темы этики данных.....	134
Защита прав и ответственность.....	134
Петли обратной связи	135
Необъективность.....	138
Дезинформация	150
Выявление этических проблем и их решение	152
Анализ проекта.....	152
Какие процессы нужно реализовать.....	153
Сила разнообразия.....	155
Справедливость, ответственность и прозрачность	157
Роль политики.....	159
Эффективность регулирования.....	159
Права и политика.....	161
Автомобили: исторический прецедент.....	161
Резюме.....	162
Вопросник.....	163
Дополнительные задания	164
Глубокое обучение на практике: итог!	165

ЧАСТЬ II ПОНИМАНИЕ ПРИЛОЖЕНИЙ НА БАЗЕ FASTAI

Глава 4. Обучение классификатора цифр: взгляд изнутри	168
Пиксели: основа компьютерного зрения.....	168
Первая попытка: сходство пикселей	172
Массивы NumPy и тензоры PyTorch	178
Вычисление метрик с помощью бродкастинга (Broadcasting)	180
Стохастический градиентный спуск	184
Вычисление градиентов.....	189
Определение шагов скорости обучения.....	191
Сквозной пример SGD.....	193
Подведение итогов темы градиентного спуска.....	198
Функция потерь MNIST	199
Сигмоида	205

SGD и мини-пакеты	206
Собрать все вместе	208
Создание оптимизатора	211
Добавление нелинейности	213
Углубляемся	217
Сводка терминов	218
Вопросник.....	220
Дополнительные задания	221
Глава 5. Классификация изображений	222
От собак и кошек к породам домашних животных.....	222
Подготовка размера.....	226
Проверка и отладка DataBlock	229
Перекрестная энтропия	231
Активации и метки	232
Softmax	233
Логарифмическая функция правдоподобия	236
Применение логарифма.....	238
Интерпретация модели	241
Улучшение модели.....	242
Поиск скорости обучения	242
Разморозка и перенос обучения.....	245
Дискриминативные скорости обучения.....	248
Выбор количества эпох	250
Углубленные архитектуры.....	251
Резюме	253
Вопросник.....	254
Дополнительные задания	255
Глава 6. Другие задачи компьютерного зрения	256
Классификация по нескольким меткам	256
Данные.....	257
Построение DataBlock	259
Бинарная перекрестная энтропия	264
Регрессия	269

Сборка данных.....	270
Обучение модели	273
Резюме.....	275
Вопросник.....	276
Дополнительные задания	277
Глава 7. Обучение современной модели	278
Imagenette	278
Нормализация	280
Прогрессивное изменение размера	282
Аугментация во время тестирования	284
Mixup	286
Сглаживание меток.....	289
Резюме.....	292
Вопросник.....	292
Дополнительные задания	293
Глава 8. Коллаборативная фильтрация	294
Первый взгляд на данные.....	295
Обучение скрытых факторов	297
Создание DataLoaders	299
Коллаборативная фильтрация с нуля	302
Сокращение весов	306
Создание собственного модуля вложений	307
Интерпретация вложений и смещений.....	309
Использование fastai.collab	311
Расстояние между вложениями.....	312
Бутстрэппинг модели коллаборативной фильтрации.....	312
Глубокое обучение для коллаборативной фильтрации	314
Резюме	317
Вопросник.....	317
Дополнительные задания	319
Глава 9. Табличное моделирование	320
Категориальные вложения	320
За гранью глубокого обучения	326

Датасет.....	328
Соревнования Kaggle	328
Знакомство с данными.....	330
Деревья решений.....	331
Обработка дат	333
Использование TabularPandas и TabularProc.....	334
Создание дерева решений.....	337
Категориальные переменные.....	342
Случайные леса	343
Создание случайного леса	344
Ошибка Out-of-Bag.....	346
Интерпретация модели	347
Дисперсия деревьев для уверенного прогнозирования	348
Важность признаков	349
Удаление переменных с низкой важностью	350
Удаление лишних признаков	351
Частичная зависимость	354
Утечка данных	357
Интерпретатор деревьев.....	358
Экстраполяция и нейронные сети.....	360
Проблема экстраполяции.....	361
Поиск несоответствующих области данных.....	362
Использование нейронной сети	365
Ансамблирование	369
Бустинг.....	370
Совмещение вложений с другими методами.....	371
Резюме.....	372
Вопросник.....	373
Дополнительные задания	375
Глава 10. Погружение в NLP: рекуррентные нейронные сети	376
Предварительная обработка текста.....	378
Токенизация	379
Токенизация слов с помощью fastai	380
Токенизация подслов	384

Нумеризация с помощью fastai	385
Разделение текстов на пакеты.....	387
Обучение классификатора текста.....	390
Создание языковой модели с помощью DataBlock.....	391
Тонкая настройка языковой модели	392
Сохранение и загрузка моделей	393
Генерация текста.....	395
Создание DataLoaders классификатора	395
Тонкая настройка классификатора	398
Дезинформация и языковые модели.....	399
Резюме.....	402
Вопросник.....	403
Дополнительные задания	404
 Глава 11. Преобразование данных с помощью Mid-Level API	405
Знакомство с многослойным API	405
Преобразования	406
Написание собственного преобразования.....	408
Pipeline.....	409
TfmdLists и датасеты: преобразованные коллекции.....	410
TfmdLists	410
Datasets.....	412
Использование промежуточного API: SiamesePair	414
Резюме.....	419
Вопросник.....	419
Дополнительные задания	420
Сфера применения fastai: обобщение	420
 ЧАСТЬ III	
ОСНОВЫ ГЛУБОКОГО ОБУЧЕНИЯ	
 Глава 12. Языковая модель с нуля	422
Данные	422
Первая языковая модель с нуля.....	424
Языковая модель в PyTorch.....	425
Первая рекуррентная нейронная сеть.....	428

Улучшение RNN	430
Управление состоянием RNN	430
Создание дополнительного сигнала	433
Многослойные RNN.....	436
Модель.....	437
Взрывающиеся или исчезающие активации	438
LSTM.....	439
Создание LSTM с нуля.....	440
Обучение языковой модели с помощью LSTM	442
Регуляризация LSTM	444
Dropout.....	444
Регуляризация активаций и регуляризация временных активаций.....	447
Обучение регуляризованной LSTM со связанными весами.....	447
Резюме.....	449
Вопросник.....	450
Дополнительные задания	452
Глава 13. Сверточные нейронные сети	453
Магия сверток	453
Отображение ядра свертки	457
Свертки в PyTorch	459
Штрихи и заполнение.....	461
Понимание сверточных уравнений	462
Первая сверточная нейронная сеть	465
Создание CNN.....	465
Разъяснение арифметики сверток	468
Рецептивные поля.....	469
О Twitter.....	471
Цветные изображения	473
Повышение стабильности обучения.....	476
Базовая модель	477
Увеличение размера пакета	480
Обучение 1cycle	480
Пакетная нормализация	485
Резюме.....	489

Вопросник.....	489
Дополнительные задания	491
Глава 14. ResNet.....	492
Возвращение к Imagenette	492
Построение современной CNN: ResNet	496
Пропускающие соединения	496
Актуальная ResNet	503
Зауженные слои	506
Резюме.....	508
Вопросник.....	509
Дополнительные задания	510
Глава 15. Архитектуры приложений	511
Компьютерное зрение	511
cnn_learner.....	511
unet_learner.....	513
Сиамская сеть.....	516
Обработка естественного языка	518
Табличные модели	519
Резюме.....	520
Вопросник.....	522
Дополнительные задания	523
Глава 16. Процесс обучения	524
Создание базовой модели	524
Универсальный оптимизатор	526
Импульс.....	527
RMSProp	530
Adam.....	531
Раздельное сокращение весов	532
Обратные вызовы.....	533
Создание обратного вызова.....	536
Упорядочивание обратных вызовов и исключения.....	539
Резюме.....	540
Вопросник.....	541

Дополнительные задания	542
Основы глубокого обучения: итог	542

ЧАСТЬ IV **ГЛУБОКОЕ ОБУЧЕНИЕ С ЧИСТОГО ЛИСТА**

Глава 17. Продвинутые основы нейронной сети	544
Создание слоя нейронной сети с нуля.....	544
Моделирование нейрона	544
Матричное умножение	546
Поэлементная арифметика.....	547
Уширение (broadcasting)	549
Соглашение Эйнштейна.....	553
Прямой и обратный проход	555
Определение и инициализация слоя.....	555
Градиенты и обратный проход.....	560
Рефакторинг модели	563
Переходим в PyTorch.....	564
Резюме.....	567
Вопросник.....	568
Дополнительные задания	570
Глава 18. Интерпретация CNN с помощью CAM.....	571
CAM и хуки.....	571
CAM градиентов.....	575
Резюме.....	577
Вопросник.....	577
Дополнительные задания	578
Глава 19. Класс Learner с нуля	579
Данные	579
Dataset	581
Module и Parameter	584
Простая CNN	587
Функция потерь	588

Learner	590
Обратные вызовы.....	591
Планирование скорости обучения	593
Резюме.....	595
Вопросник.....	595
Дополнительные задания	597
Глава 20. Подведем итог	598

ПРИЛОЖЕНИЯ

Приложение А. Создание блога.....	602
Блогинг на GitHub Pages.....	602
Создание репозитория	603
Настройка домашней страницы	604
Создание публикаций	606
Синхронизация GitHub и компьютера.....	608
Блогинг из Jupyter	610
Приложение Б. Схема подготовки проекта по аналитике данных.....	611
Специалисты по данным.....	612
Стратегия.....	613
Данные	615
Аналитика	616
Реализация	616
Обслуживание	617
Ограничения.....	618
Об авторах	619
Благодарности.....	620
Об обложке	622

Отзывы о книге

Если вам нужен гайд, который проведет вас от первых исследований до передовых рубежей в сфере глубокого обучения, эта книга для вас. Не только доктора наук могут наслаждаться преимуществами глубокого обучения — вы тоже можете использовать его для решения рабочих задач.

Хэл Вариан (Hal Varian), заслуженный профессор Калифорнийского университета в Беркли, главный экономист Google

Сегодня искусственный интеллект активно развивается в области глубокого обучения, поэтому настало время разобраться в принципах его работы. Данная книга позволяет сделать это даже непосвященным. Авторы смогли в упрощенном виде передать то, что иные сочли бы очень сложным.

Эрик Тополь (Eric Topol), автор книги Deep Medicine, профессор Scripps Research

Джереми и Сильвейн приглашают вас в интерактивное путешествие, где каждую строчку кода вы сможете выполнить на своем ноутбуке. Это будет путешествие по долинам потерь и пикам производительности. Приправленная остроумными шутками, эта книга поражает балансом изложения глубоких технических принципов и непринужденной манеры разговора. Она даст вам эталонные практические инструменты и реальные примеры для их использования. Независимо от того, новичок вы или опытный разработчик, книга позволит ускорить ваше обучение и постичь не только вершины, но и глубины машинного обучения.

Себастьян Рудер (Sebastian Ruder), ученый-исследователь Deepmind

Джереми Ховард и Сильвейн Гуггер написали удивительную книгу — она проектирует мост между областью ИИ и всем миром. Эта работа — исключительно содержательное и проницательное, но одновременно и понятное пособие по глубокому обучению для всех, кто интересуется этой областью.

Энтони Чанг (Anthony Chang), директор по информационным технологиям, детская больница округа Оранж

Как освоить глубокое обучение и не завязнуть в его деталях? Как быстро разобраться в его принципах, тонкостях и технических приемах с помощью при-

меров и кода? Все ответы здесь. Не пропустите новый источник практического глубокого обучения.

Орен Эциони (Oren Etzioni), профессор Вашингтонского университета, исполнительный директор, Институт ИИ Аллена

Эта книга — редкая жемчужина, тщательно проработанный и высокоэффективный учебный продукт, отточенный несколькими годами повторений не на одной тысяче студентов, одним из которых являюсь и я. Курс Fast.ai чудесным образом изменил мою жизнь и, без сомнений, может сделать то же самое для вас.

Джейсон Антик (Jason Antic), создатель DeOldify

Данная книга уже в первых главах учит эффективно использовать глубокое обучение. После этого в ней тщательно, но на доступном уровне рассматриваются внутренние процессы моделей машинного обучения и фреймворков. Хотел бы я иметь под рукой такую книгу, когда только начинал осваивать машинное обучение!

Эммануэль Амейсен (Emmanuel Ameisen), автор приложений на основе машинного обучения

«Глубокое обучение — для всех» — утверждается в первой главе. Подобное заявление можно встретить и в других книгах, но здесь оно раскрывается в полной мере. Авторы обладают обширными знаниями и при этом способны доступно донести их людям, знакомым с программированием, но не имеющим опыта в машинном обучении. Сначала даются примеры, и только затем рассматривается теория. Для большинства людей это наилучший способ обучения. Здесь вы найдете приложения глубокого обучения для области компьютерного зрения, обработки естественного языка и табличных данных, а также ознакомитесь с такой важной темой, как этика данных, которой так не хватает многим книгам. Материал представляет собой один из лучших ресурсов для программиста, изучающего глубокое обучение.

Питер Норвиг (Peter Norvig), директор по исследованиям Google

Гуггер и Ховард создали идеальный источник знаний для тех, кто хоть немного знаком с программированием. Благодаря практическому подходу и предварительно написанному коду они отлично раскрывают тему глубокого обучения. Больше не нужно мучиться с теоремами и доказательствами абстрактных понятий. Уже в первой главе вы создадите свою модель глубокого обучения, а к концу книги сможете читать и понимать раздел «Методы» любой научной работы по глубокому обучению.

Кертис Ланглотц (Curtis Langlotz), директор центра использования искусственного интеллекта для медицины и обработки изображений, Стэнфордский университет

Эта книга проливает свет на самый черный из всех черных ящиков — глубокое обучение. Она позволяет начать активно экспериментировать с кодом, используя возможности Python. Глубоко затрагиваются этические последствия применения ИИ. Авторы показывают, как не допустить появления антиутопии.

Гийом Шасло (Guillaume Chaslot), сотрудник Mozilla

Меня часто спрашивают, как начать осваивать глубокое обучение. Я же всегда указываю на fastai. Эта книга решает, казалось бы, невыполнимую задачу, выступая понятным руководством для сложного предмета. При этом она полна передовых открытий, которые, без сомнения, оценят даже опытные разработчики.

Кристин Пейн (Christine Payne), специалист OpenAI

Чрезвычайно практичная и доступная книга, которая поможет быстро начать создавать собственные проекты. Это очень понятное и легко читаемое руководство по практическому глубокому обучению, которое окажется полезным для всех: от начинающих программистов до руководителей и менеджеров проектов. Такую книгу я мечтала иметь много лет назад!

*Кэрол Райли (Carol Reiley),
президент-учредитель и председатель Drive.ai*

Компетентность Джереми и Сильвейна в области глубокого обучения, их практический подход, а также участие в открытых проектах сделали их ключевыми фигурами в сообществе PyTorch. Книга является продолжением работы, которую ее авторы совместно с сообществом fast.ai делают для повышения доступности машинного обучения, что в дальнейшем окажет серьезное влияние на всю сферу ИИ.

*Джером Песенти (Jerome Pesenti),
вице-президент подразделения ИИ в Facebook*

Сегодня глубокое обучение — одна из наиболее важных технологий. С ее помощью были достигнуты многие прорывы в сфере ИИ. Глубокое обучение больше не является областью, доступной только для ученых. Эта книга, основанная на популярном курсе fast.ai, делает глубокое обучение ближе всем тем, кто имеет опыт программирования. Обучение строится по целостному принципу, с помощью практических примеров и интерактивного сайта. Кстати, для кандидатов наук здесь тоже найдется много интересного.

*Григорий Пятицкий-Шапиро (Gregory Piatetsky-Shapiro),
президент KDnuggets*

Эта книга — продолжение курса fast.ai, который я неустанно рекомендовал всем на протяжении нескольких лет. Она превратит вас из новичка в опытного практика буквально за считанные месяцы. Наконец-то появилось что-то стоящее!

*Луи Монье (Louis Monier),
учредитель Altavista; бывший глава Airbnb AI Lab*

Рекомендуем эту книгу! В ней описаны продвинутые фреймворки, позволяющие без лишней суеты проработать реальные задачи ИИ и автоматизации. В итоге у вас остается время для часто игнорируемых тем, таких как безопасный перенос моделей в продакшен и этика данных.

*Джон Маунт (John Mount) и Нина Зумель (Nina Zumel),
авторы книги Practical Data Science with R*

Эта книга предназначена для программистов и подойдет даже тем, у кого нет ученой степени. Хоть я и обладаю такой степенью, но я не программист. Так почему же меня попросили оценить ее? Думаю, чтобы я рассказал вам, что она чертовски крута!

Прочитав первые две страницы первой главы, вы узнаете, как с помощью всего четырех строк кода и менее одной минуты вычислений получить эталонную сеть, способную отличать кошек от собак. Во второй главе вы перейдете от модели в продакшен, узнав, как быстро разместить приложение в вебе, не прибегая к HTML или JavaScript и не имея собственного сервера.

Эта книга как луковица. Вы получаете полноценный пакет с наилучшими возможными настройками. Если потребуется что-то переделать, то вы просто снимаете внешний слой. Нужны дополнительные настройки? Можно продолжить снимать оболочки. Еще? Углубляйтесь до уровня использования только PyTorch. На всем 600-страничном путешествии вас будут направлять три голоса, предлагая свою точку зрения.

*Альфредо Канциани (Alfredo Canziani),
профессор Нью-Йоркского университета NYU*

Это доступная и построенная на принципе диалога книга, которая рассказывает об основах и понятиях глубокого обучения. Авторы с самого начала приводят практические примеры, а теорию — по мере необходимости. Практик может начать знакомство с миром глубокого обучения на примерах в первой половине книги, а глубокие принципы открываются уже во второй части. И наконец-то все мифы будут окончательно развеяны.

Джош Паттерсон (Josh Patterson), Patterson Consulting

Введение

Глубокое обучение (deep learning, DL) — новая мощная технология, и мы считаем, что ее следует использовать в разных дисциплинах. Эксперты, скорее всего, найдут новые возможности ее применения, и нам хотелось бы, чтобы специалисты различного профиля также включились в процесс обучения.

По этой причине и не только Джереми принял участие в создании fast.ai, организации, которая стремится сделать глубокое обучение доступным с помощью бесплатных онлайн-курсов и ПО. Сильвейн — инженер-исследователь в Hugging Face. До этого он занимал должность ученого-исследователя в fast.ai, а также преподавал математику и computer science по программе подготовки студентов к поступлению в элитные университеты Франции. Мы вместе написали эту книгу, чтобы помочь как можно большему числу людей начать использовать глубокое обучение.

Для кого эта книга

Если вы новичок в области глубокого и машинного обучения, эта книга для вас. Но желательно уметь писать код на Python.



НЕТ ОПЫТА? НЕ ПРОБЛЕМА!

Если же код вы раньше не писали, то и это не проблема. Первые три главы написаны так, чтобы руководители, менеджеры по продукту и другие специалисты без труда поняли наиболее важные аспекты глубокого обучения. Встречая в тексте отрывки кода, просматривайте их, чтобы сформировать интуитивное понимание его работы. Мы же будем объяснять код строка за строкой. В этом случае не так важен синтаксис, как высокоуровневое понимание.

Если вы уверенно себя чувствуете в глубоком обучении, то все равно найдете в книге много полезного. Мы покажем, как добиваться высоких результатов, и познакомим вас с новейшими передовыми техниками. Вы увидите, что для этого не требуется глубокое знание математики. Нужны лишь здравый смысл и упорство.

Что нужно знать

Единственное (желательное) требование — умение писать код (года опыта будет достаточно), лучше всего на Python, а также знание математики из курса старших классов. Даже если вы ее подзабыли, мы поможем освежить знания. Отличным вспомогательным ресурсом послужит Khan Academy (<https://www.khanacademy.org/>).

Мы не говорим, что в глубоком обучении нет математики сложнее той, что изучают в старшей школе, но мы дадим вам основы, которые понадобятся для понимания рассматриваемых тем (или направим на нужные ресурсы).

Начнем с общей картины и постепенно будем углубляться, поэтому время от времени вам может потребоваться изучить дополнительные темы (программирование чего-либо или математику). Это нормально — именно так мы и будем строить работу. Читайте книгу и изучайте дополнительные ресурсы по мере необходимости.



Онлайн-ресурсы

Все примеры кода доступны онлайн в виде блокнотов Jupyter (в главе 1 мы расскажем, что это). Это интерактивная версия книги, где можно выполнять код и экспериментировать с ним. Дополнительно об этом можно прочитать на сайте книги (<https://book.fast.ai/>). Там же содержится актуальная информация о настройке различных инструментов и дополнительные главы.

Чему вы научитесь

Прочитав эту книгу, вы будете знать следующее.

- Как обучать модели, достигающие эталонных результатов:
 - в компьютерном зрении, включая распознавание изображений (например, классификацию фотографий домашних животных по породам), а также локализацию изображений и обнаружение (например, поиск животных на изображении);
 - в обработке естественного языка (NLP), включая классификацию документов (например, анализ тональности) и языковое моделирование;
 - в обработке табличных данных (например, прогнозировании продаж) с категориальными, непрерывными и смешанными данными, включая временной ряд;
 - в совместной фильтрации (например, рекомендации фильмов).
- Как преобразовывать модели в веб-приложения.

- Почему и как работают модели глубокого обучения, а также как на основе этих знаний повысить точность, скорость и надежность таких моделей.
- Как работают новейшие техники глубокого обучения, наиболее актуальные в практических задачах.
- Как читать исследовательские работы по глубокому обучению.
- Как реализовывать алгоритмы глубокого обучения с нуля.
- Как анализировать этические последствия вашей работы, чтобы в итоге она позволила сделать мир лучше и не могла быть использована во вред.

Полный перечень в содержании, но для краткого представления мы приведем несколько техник, которые рассмотрим (не переживайте, если слова кажутся непонятными, — скоро вы все узнаете):

- аффинные функции и нелинейности;
- параметры и активации;
- случайная инициализация и трансферное обучение (transfer learning);
- SGD, Momentum, Adam и другие оптимизаторы;
- свертки;
- пакетная нормализация;
- дропаут;
- увеличение (аугментация) данных;
- уменьшение весов;
- архитектуры ResNet и DenseNet;
- классификация и регрессия изображений;
- вложения;
- рекуррентные нейронные сети (RNN);
- сегментация;
- U-Net
- и многое другое!



ВОПРОСНИКИ

В конце каждой главы вы найдете тест. Он поможет закрепить полученные знания — если (мы надеемся!) вы сможете ответить на все вопросы. Кстати, один из рецензентов книги (спасибо, Фред!) сказал, что предпочитает начинать главу именно с вопросника, чтобы сразу понимать, на что обращать внимание.

Предисловие

Глубокое обучение весьма быстро стало популярной методикой, которая решает и автоматизирует задачи в области компьютерного зрения, робототехники, здравоохранения, физики, биологии и т. д. К приятным особенностям глубокого обучения можно отнести его сравнительную простоту. Есть много софта, благодаря которому буквально за несколько недель можно понять основы и освоить необходимые техники.

Здесь открывается целый мир творчества. Вы используете полученные знания для решения задач с данными, наблюдая, как машина справляется с этими задачами за вас. Затем вы приближаетесь к новому этапу, когда, создав модель глубокого обучения, понимаете, что работает она не так, как вы рассчитывали. И дальше вы начинаете искать и изучать новейшие исследования в глубоком обучении.

Эта область охватывает невероятный объем знаний — теорию и великое множество техник и инструментов, стоящих за ней. Надо признать, что люди склонны объяснять простое сложным языком. Ученые используют в своих работах непонятные слова и математические выражения, и кажется, что ни посты, ни туториалы не освещают нужную информацию доступным образом. Инженеры и программисты предполагают, что вы уже знаете, как работают GPU, и разбираетесь в малоизвестных инструментах.

Именно здесь закрадывается мысль, что неплохо бы иметь наставника или друга, которому можно задать вопросы. Кого-то, кто уже бывал в похожей ситуации и знает и инструменты, и математику, кто сможет рассказать о наилучших исследованиях, эталонных техниках и продвинутых инженерных решениях, до смешного упрощая задачу. Я был в похожей ситуации около десяти лет назад, когда только начинал знакомиться с машинным обучением (machine learning, ML). Несколько лет я с трудом понимал работы, где была математика. Меня окружали хорошие наставники, которые здорово помогали, но чтобы начать уверенно использовать машинное и глубокое обучение, потребовалось годы. Это послужило мотивом принять участие в разработке PyTorch, фреймворка, делающего глубокое обучение доступным.

Джереми Говард и Сильвейн тоже были на вашем месте. Они начали осваивать и применять глубокое обучение, не являясь учеными или инженерами в области ML. Как и я, Джереми и Сильвейн учились несколько лет и в итоге стали

не только экспертами, но и лидерами. Но в отличие от меня, они делали все возможное, чтобы в будущем другим людям не пришлось идти по такому же сложному пути. Они разработали отличный курс под названием fast.ai, который делает передовые техники глубокого обучения доступными для тех, кто умеет программировать. Сотни тысяч желающих окончили этот курс, став отличными практиками.

В книге Джереми и Сильвейн используют простые слова и поясняют каждое понятие, доступно описывая передовые техники глубокого обучения и эталонные исследования.

Вас познакомят с последними достижениями в области компьютерного зрения, научат обработке естественного языка и дадут основы математики. Но на этом веселье не заканчивается — ведь затем вам покажут, как реализовать ваши идеи в продакшене. Пускай сообщество fast.ai, объединяющее тысячи практиков онлайн, станет вашей семьей, где все люди могут обсуждать и находить решения любых проблем.

Я очень рад, что вы взялись за эту книгу, и надеюсь, она вдохновит вас на правильное использование глубокого обучения, независимо от области решаемой задачи.

Сумит Чинтала (Soumith Chintala), соавтор PyTorch

От издательства

Ваши замечания, предложения, вопросы отправляйте по адресу comp@piter.com (издательство «Питер», компьютерная редакция).

Мы будем рады узнать ваше мнение!

На веб-сайте издательства www.piter.com вы найдете подробную информацию о наших книгах.

ЧАСТЬ I

**Глубокое обучение
на практике**

ГЛАВА 1

Путешествие в мир глубокого обучения

Приветствуем вас и хотим поблагодарить за возможность попутешествовать с вами по миру глубокого обучения, как бы далеко вы в нем ни зашли. Поговорим сначала немного о том, что вас ждет на страницах книги, представим основные принципы, стоящие за глубоким обучением, и обучим нашу первую модель для различных задач. Неважно, есть ли у вас опыт работы в технической сфере (но если есть, то замечательно), — эту книгу мы написали понятным для широкой аудитории языком.

Глубокое обучение для всех

Многие предполагают, что для успеха в глубоком обучении нужны разные труднодоступные материалы. Но как вы увидите позже, это ошибочное мнение. В табл. 1.1 перечислено то, что вам *совершенно не понадобится* для достижения успеха.

Таблица 1.1. Что НЕ нужно для глубокого обучения

Миф	Правда
Много математики	Программы старших классов школы достаточно
Много данных	Мы встречали рекордные результаты с использованием <50 элементов данных
Дорогое оборудование	Все необходимое для работы вы можете получить бесплатно

Глубокое обучение — это компьютерная методика для извлечения и преобразования данных. Область ее использования простирается от распознавания человеческой речи до классификации образов животных. Делается это через

несколько слоев нейронных сетей, каждый из которых получает вводные от предыдущего и постепенно их уточняет. Слои обучаются через алгоритмы, которые минимизируют ошибки и одновременно повышают точность. Таким образом сеть учится выполнять конкретные задачи. В следующем разделе мы подробно рассмотрим алгоритмы обучения.

В глубоком обучении кроется мощь, гибкость и простота. Именно поэтому мы верим, что его следует применять во многих дисциплинах, включая гуманистические и технические науки, искусство, медицину, финансы и многое другое. Вот вам личный пример: несмотря на отсутствие опыта в медицине, Джереми организовал Entilic – компанию, использующую алгоритмы глубокого обучения для постановки диагнозов. Спустя всего несколько месяцев после запуска компании было объявлено, что используемые алгоритмы способны обнаруживать злокачественные опухоли *точнее, чем рентгенологи* (<https://oreil.ly/aTwdE>).

Вот список лишь некоторых из множества задач, для которых глубокое обучение или использующие его методы являются лучшими в мире.

Обработка естественного языка (NLP)

Ответы на вопросы; распознавание речи; классификация документов; обнаружение имен, дат и т. п. в документах; поиск статей, содержащих определенное понятие.

Компьютерное зрение

Интерпретация снимков со спутников и дронов (например, для предупреждения стихийных бедствий), распознавание лиц, захват изображений, чтение дорожных знаков, обнаружение автономными транспортными средствами пешеходов и других транспортных средств.

Медицина

Обнаружение аномалий на радиологических снимках, включая КТ, МРТ и рентген; подсчет признаков патологий на слайдах; измерение признаков в ультразвуке; диагностирование диабетической ретинопатии.

Биология

Свертывание белков; классификация белков; ряд задач геномики, включая секвенирование «опухоль – норма» и классификацию клинически значимых генетических мутаций; классификация клеток; анализ белков и белковых взаимодействий.

Создание изображений

Раскрашивание изображений, увеличение их разрешения, удаление шума, преобразование изображений в произведения искусства, оформленные в стиле известных художников.

Системы рекомендаций

Веб-поиск, рекомендации, создание домашней страницы.

Игры

Шахматы, го, большинство видеоигр от Atari, разные стратегии в реальном времени.

Робототехника

Обработка сложных для обнаружения объектов (например, прозрачных, блестящих, с отсутствием текстур) или тех, которые сложно захватить.

Другие сферы

Финансовое и логистическое прогнозирование, перевод текста в речь и многое другое...

Хотя глубокое обучение и находит столь разнообразное применение, в основе практически всех его решений лежит один инновационный тип модели: нейронная сеть.

Но нейросети не настолько новы. И чтобы во всем разобраться, начнем с небольшого исторического экскурса.

Нейронные сети: краткая история

В 1943 году нейрофизиолог Уоррен Маккаллок (Warren McCulloch) и логик Уолтер Питтс (Walter Pitts) занялись разработкой математической модели искусственного нейрона. В работе «Логическое исчисление идей, относящихся к нервной активности» они заявили следующее:

Из-за взаимодействия нейронов, характеризуемого как «все или никто», нейронные события и отношения между ними можно трактовать через логику высказываний. Выяснилось, что поведение каждой сети можно описать в этих терминах.

Маккаллок и Питтс поняли, что упрощенную модель реального нейрона можно представить с помощью простого сложения и определения порога (рис. 1.1). Питтс был самоучкой и к 12 годам уже получил предложение учиться в Кэмбриджском университете у великого Бертранда Рассела. Он отказался, как и от всех остальных предложений о присвоении ученых степеней и руководящих должностей, которые получал в течение жизни. Большая часть его выдающихся работ была создана им во времена, когда у него даже не

было собственного жилья. Несмотря на отсутствие официально признанного положения и увеличивающуюся социальную изоляцию Питтса, его работа с Маккаллоком серьезно заинтересовала психолога Фрэнка Розенблатта (Frank Rosenblatt).

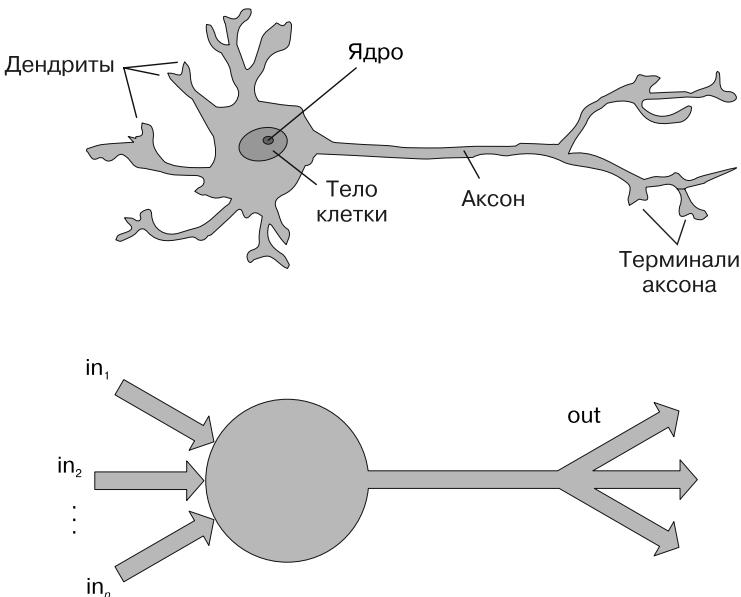


Рис. 1.1. Естественный и искусственный нейроны

Розенблatt усовершенствовал искусственный нейрон, чтобы тот смог обучаться. Более того, он работал над созданием первого устройства, использующего эти принципы: перцептрона «Марк-1». В своей работе *The Design of an Intelligent Automaton* Розенблatt писал об этом так: «Мы готовимся узреть рождение такой машины — машины, способной воспринимать, распознавать и идентифицировать свое окружение, не требуя какого-либо обучения или контроля со стороны человека». Перцептрон был создан и мог успешно распознавать простые формы.

Профессор Массачусетского технологического института (MIT) Марвин Мински (Marvin Minsky) (который в школе учился в одном классе с Розенблаттом!) совместно с Сеймуром Папертом (Seymour Papert) написал книгу *Perceptrons* (MIT Press), в которой рассказывалось об изобретении Розенблатта. В этой книге они показали, что один слой таких устройств не способен изучить простые, но критически важные функции, например XOR. Кроме этого, они показали, что устранить эти ограничения можно через использование нескольких слоев устройств. К сожалению, широко были приняты только пер-

вые из этих выводов, и в итоге глобальное научное сообщество практически полностью прекратило изучение нейронных сетей на последующие 20 лет.

Наиболее же значительным трудом в этой области за последние 50 лет стала многотомная работа *Parallel Distributed Processing (PDP)* («Параллельная распределенная обработка»), написанная Дэвидом Румельхартом (David Rumelhart), Джеймсом Маккландом (James McClelland) и исследовательской группой PDP. Она была издана MIT Press в 1986 году. В первой главе высказывается также идея, о которой когда-то говорил Розенблatt:

Люди умнее современных компьютеров, потому что человеческий мозг использует базовую вычислительную архитектуру, а она больше подходит для решения задач по обработке естественной информации, в которых люди так хороши... Мы создадим вычислительную структуру для моделирования когнитивного процесса, которая превзойдет другие структуры и окажется более приближенной к модели вычисления, характерной для мозга.

В качестве предпосылки PDP подразумевает, что принцип работы традиционных компьютерных программ отличается от принципа работы головного мозга, в связи с чем эти программы (на тот момент) плохо справлялись с задачами, которые мозг решал легко (например, распознавание объектов на картинках). Авторы заявляли, что подход PDP был «ближе других структур» к модели функционирования человеческого мозга, вследствие чего он сможет лучше справляться с такими задачами.

В действительности подход, заложенный в PDP, очень похож на используемый в современных нейронных сетях. Согласно той книге для параллельной распределенной обработки требовалось следующее:

- набор блоков обработки;
- состояние активации;
- функция вывода для каждого блока;
- паттерн связей между блоками;
- правило распространения для распространения шаблонов активности через сеть взаимосвязей;
- правило активации для объединения вводов, влияющих на блок, с текущим состоянием этого блока, чтобы можно было генерировать его вывод;
- правило обучения, согласно которому модели связей изменяются в процессе получения опыта;
- среда, в которой система должна функционировать.

В этой книге мы увидим, что современные нейронные сети отвечают всем перечисленным требованиям.

В 1980-х большинство моделей создавались со вторым слоем нейронов, что помогало избежать проблемы, выявленной Мински и Папертом (это была их «модель связей между нейронами» для использования предыдущей структуры). Нейронные сети действительно широко использовались на протяжении 80-х и 90-х годов прошлого века в реальных практических проектах. Тем не менее недопонимание теоретических сложностей в очередной раз привело к замедлению развития области. Считалось, что добавления всего одного дополнительного слоя нейронов должно хватать, чтобы такие сети могли аппроксимировать любую математическую функцию, но на практике эти сети зачастую оказывались слишком большими и медленными для эффективного использования.

И хотя исследователи еще 30 лет назад показали, что для получения хорошей практической производительности необходимо использовать больше нейронных слоев, этот принцип был широко принят и начал применяться только в течение последнего десятилетия. Нейронные сети наконец-то приближаются к реализации собственного потенциала. К этому привело развитие компьютерного оборудования, позволившего использовать большее число слоев, а также повышение доступности данных и доработка алгоритмических приемов, что в совокупности упростило и ускорило обучение нейронных сетей. Теперь мы получили то, что обещал нам Розенблattt: «машину, способную воспринимать, узнавать и идентифицировать окружение без какого-либо обучения или контроля со стороны человека».

Именно такие программы вы и будете учиться создавать по мере прочтения данной книги. Но сначала, раз уж мы собираемся провести вместе много времени, познакомимся поближе...

Кто мы

Мы – это Сильвейн и Джереми, ваши гиды в предстоящем путешествии. Надеемся, что вы сочтете нас подходящими на эту роль.

Джереми использует и преподает машинное обучение на протяжении почти 30 лет. Применять нейронные сети он начал 25 лет назад. За все это время он курировал многие компании и проекты, использующие в своей основе ML. К его достижениям можно отнести основание первой компании Enlitic, занимающейся глубоким обучением в сфере медицины, а также деятельность на посту президента и научного руководителя в Kaggle – крупнейшем мировом сообществе машинного обучения. Совместно с доктором Рейчел Томас (Rachel Thomas) он основал fast.ai, организацию, разработавшую учебный курс, на котором основывается данная книга.

Во врезках вы будете встречать наши комментарии:



СЛОВО ДЖЕРЕМИ

Всем привет! Меня зовут Джереми! Возможно, вам будет интересно узнать, что у меня нет формального технического образования. Я отучился на бакалавра философии и более серьезных ученых степеней не получал. Мне было гораздо интереснее заниматься практическими задачами, чем изучать теорию, поэтому в студенческие годы я на полную ставку работал в консалтинговой фирме McKinsey and Company. Если вы тоже предпочитаете заниматься делом, а не тратить годы, изучая абстрактные концепции, то поймете, о чем я говорю. Не пропускайте мои комментарии — в них вы найдете информацию, предназначенную для тех, кто не имеет глубокого математического или технического бэкграунда, то есть для таких же людей, как я сам.

Сильвейн же хорошо понимает суть технического образования. Он написал десять учебников по математике, охватывающих всю продвинутую программу изучения математики во Франции!



СЛОВО СИЛЬВЕЙНУ

В отличие от Джереми, я не тратил годы на программирование и применение алгоритмов машинного обучения. Напротив, я лишь недавно познакомился со сферой ML, просматривая видеокурсы Джереми по программе fast.ai. Так что если вы тоже не открывали терминал и не писали инструкции в командной строке, то мы найдем общий язык. В моих врезках представлена информация для тех, кто хорошо ориентируется в математике или имеет техническое образование, но не имеет реального опыта программирования.

Курс fast.ai прошли сотни тысяч студентов со всего мира. Сильвейн же, по мнению Джереми, выделялся среди всех встречавшихся ему за годы обучения студентов, в связи с чем он в итоге стал сначала сотрудником fast.ai, а затем в соавторстве с Джереми написал библиотеку ПО fastai.

В нашем дуэте представлено лучшее от обоих миров: эксперт в математике и эксперт в программировании и машинном обучении, которые при этом знают, каково это — не разбираться в математике или в написании кода.

Любой, кто смотрел спортивные передачи, знает, что если действия команды комментируют двое, то им нужен третий участник для «особых комментариев». Таким комментатором будет Алексис Галлахер (Alexis Gallagher). У Алексиса очень разнообразный опыт: он занимался исследованиями в области математической биологии, писал киносценарии, был исполнителем-импровизатором, консультантом в фирме McKinskey (как и Джереми!), программировал на Swift и даже занимал должность технического директора.



СЛОВО АЛЕКСИСУ

Я решил, что пришло время узнать об этих ваших искусственных интеллектах. Перепробовал я очень многое... Но на самом деле у меня нет опыта разработки моделей ML. И все же... Насколько это может быть сложно? По ходу книги я собираюсь учиться вместе с вами. Заглядывайте в мои комментарии — здесь вас ждут подсказки, которые пригодились мне самому на моем пути. Надеюсь, что и для вас они окажутся полезны.

Как изучать глубокое обучение

Профессор Гарварда Дэвид Перкинс, написавший книгу *Making Learning Whole*, имеет богатейший опыт преподавания. Основная его идея состоит в обучении с помощью целостного подхода. Это означает, что если вы учите людей играть в бейсбол, то сначала ведете их на матч или на площадку для игры. Вы не учите их, как наматывать шпагат, чтобы сделать бейсбольный мяч с нуля, и не рассказываете о физических свойствах или коэффициентах трения мяча о биту.

Пол Локхарт (Paul Lockhart), доктор математических наук Колумбийского университета, бывший профессор Университета Брауна, а также учитель математики полного курса, в своем известном эссе *A Mathematician's Lament* («Плач математика», <https://oreil.ly/yNimZ>) рисует мрачный мир, где музыке и искусству обучают так же, как математике. Детям не разрешается слушать музыку или играть, пока они не проведут более десяти лет за освоением нотной грамоты и теории, бесконечно перенося на уроках ноты в другую тональность. В художественном же классе студенты изучают краски и кисти, но рисовать им разрешается только после поступления в колледж. Звучит абсурдно, не так ли? А ведь именно так преподают математику: мы требуем, чтобы студенты годами занимались механическим запоминанием и изучением сухих, разрозненных основ, которые, как мы их заверяем, пригодятся в дальнейшем, когда многие из них уже просто бросят изучение этого предмета.

К сожалению, именно с этого и начинаются многие обучающие программы по глубокому обучению: в них учащихся просят выучить матрицы и теоремы Тейлора для аппроксимации функций потерь, не приводя никаких рабочих примеров кода. Мы не против теории аппроксимации. Мы ее любим — Сильвейн даже преподавал ее в колледже, но мы не считаем, что это лучшая отправная точка для знакомства с глубоким обучением.

В глубоком обучении очень важно регулярно корректировать модель, чтобы она с каждым разом все лучше выполняла поставленную задачу. Вот тогда вы и начинаете изучать соответствующую теорию. Но для начала у вас должна быть модель. Почти всему мы учим на реальных примерах. По мере создания этих примеров мы будем постепенно углубляться и покажем, как совершенствовать свои проекты. Шаг за шагом вы будете изучать необходимые

теоретические основы в соответствующем контексте и понимать, что и как работает.

Итак, мы обязуемся, что на протяжении книги будем следовать этим принципам.

Целостный подход обучения

Начнем с того, что покажем, как использовать завершенную рабочую эталонную сеть глубокого обучения для решения реальных задач с помощью простых, но наглядных инструментов. Затем постепенно мы будем все глубже и глубже вникать в то, как создаются эти инструменты, как создаются инструменты для создания инструментов и т. д.

Обучение только на примерах

Мы не будем жонглировать алгебраическими символами, а дадим вам контекст и цель, которые можно понять интуитивно.

Максимальное упрощение

Мы провели годы за разработкой инструментов и преподаванием методов, которые делают некогда сложные темы простыми.

Устранение барьеров

Глубокое обучение до этого времени было областью для избранных. Мы делаем ее открытой и даем возможность поработать в ней каждому.

Сложнейшая часть глубокого обучения при использовании на практике: откуда вам знать, что данных достаточно, что они находятся в нужном формате, правильно ли обучается ваша модель и что делать, если неправильно? Именно поэтому мы делаем ставку на обучение в процессе работы. Как и в случае с базовыми навыками в науке о данных, в глубоком обучении вы совершенствуетесь только через практический опыт. Попытки тратить чересчур много времени на теорию могут привести к обратным результатам. Главное — просто заниматься написанием кода и стараться решать задачи: теория придет позже, когда у вас уже будет контекст и мотивация.

Иногда вам будет тяжело, местами вы будете застревать, но не сдавайтесь! Возвращайтесь к той части книги, в которой полностью разобрались, а затем не спеша продолжайте читать с этого места в поиске первого непонятного момента. Затем постарайтесь самостоятельно поэкспериментировать с кодом, погуглите дополнительные уроки по возникшему вопросу — это поможет найти иной угол обзора и разобраться. Абсолютно нормально, что при первом прочтении некоторые моменты, особенно код, будут непонятны. Иногда бывает сложно понять материал последовательно, не забегая вперед. Бывает, что все встает на свои места, когда вы получаете больше контекста из последующего материа-

ла, показывающего более общую картину. Поэтому если вы вдруг застрянете в каком-то разделе, то попробуйте перейти дальше, отметив, что к этому разделу нужно вернуться позднее.

Чтобы достичь успеха в области глубокого обучения, не требуется особая академическая подготовка. Многие важнейшие прорывы в исследованиях и самой индустрии были сделаны людьми, не имеющими ученых степеней. Например, можно привести работу *Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks* (<https://oreil.ly/JV6rL>) — один из наиболее влиятельных трудов последнего десятилетия, который был процитирован более 5000 раз. А написал его Алек Рэдфорд (Alec Radford) еще в годы своего студенчества. Илон Маск — генеральный директор компании Tesla, решающей чрезвычайно сложные задачи по созданию беспилотных машин, говорит (<https://oreil.ly/nQCmO>):

Ученая степень точно не требуется. Важно лишь глубокое понимание ИИ и способность реализовывать нейронные сети так, чтобы они были действительно полезны (и последнее — очень трудно). При этом неважно, окончили вы вообще среднюю школу или нет.

Для достижения успеха вам нужно будет применять полученные в этой книге знания в собственных проектах и упорно продолжать делать свое дело.

Ваши проекты и мышление

Не так важно, хотите ли вы определять болезни растений по фотографиям их листьев, генерировать узоры для вязания, диагностировать туберкулез по рентгеновскому снимку, мы поможем начать использовать глубокое обучение для решения насущных задач максимально быстро (с помощью предварительно обученных другими моделей), а затем будем постепенно углубляться в подробности процесса. После прочтения следующей главы вы узнаете, как использовать глубокое обучение для решения своих задач с эталонной точностью. (Если вам не терпится поскорее приступить к работе с кодом, смело переходите к ней.) Есть миф, что для работы с глубоким обучением требуются огромные вычислительные ресурсы и наборы данных объемом как у Google, но это совершенно не так.

Какие же задачи подходят для хороших тестовых примеров? Можно обучить модель различать картины Пикассо и Моне или выбирать фотографии дочери, а не сына. Это помогает сосредоточиться на ваших увлечениях и пристрастиях — выбор четырех или пяти небольших проектов вместо стремления решить одну огромную суперзадачу упростит начало изучения данной области. Поскольку в этом деле очень легко застрять, излишнее рвение на ранних этапах может дать обратный эффект. Но когда вы уже освоите основы, можете нацеливаться на решение более серьезных задач.



СЛОВО ДЖЕРЕМИ

Глубокое обучение можно использовать для решения почти любой задачи. Например, моим первым стартапом была компания FastMail, запущенная в 1999 году. Она предоставляла расширенные сервисы электронной почты (кстати, она работает и по сей день). В 2002 году, стремясь улучшить качество сортировки писем и защиты пользователей от спама, я настроил ее на использование простейшей формы глубокого обучения, а именно однослойных нейронных сетей.

Как правило, люди, успешно справляющиеся с глубоким обучением, отличаются любопытством. Покойный физик Ричард Фейнман — пример того, кто, по нашему мнению, преуспел бы в глубоком обучении: в основу его работы над движением субатомных частиц легло увлеченное наблюдение за колеблющимися в воздухе тарелками для игры в фрисби.

А теперь сфокусируемся на том, что конкретно вы будете изучать, и начнем с софта.

ПО: PyTorch, fastai и Jupyter (почему это важно)

Мы завершили сотни проектов машинного обучения, используя десятки разных пакетов и множество языков программирования. В fast.ai мы подготовили онлайн-курсы с помощью большинства из основных пакетов ML и DL. После того как в 2017 году появился PyTorch, мы провели более тысячи часов за его тестированием, прежде чем решили, что будем использовать этот пакет в будущих курсах, разработке ПО и исследованиях. С тех пор PyTorch стал самой быстро растущей библиотекой глубокого обучения, которая уже применяется в большинстве исследовательских работ, представленных на престижных конференциях. По большому счету, это главный показатель используемости в индустрии, потому что эти работы в итоге применяются в коммерческих продуктах и сервисах. Мы выяснили, что PyTorch — это наиболее гибкая и выразительная библиотека для глубокого обучения. Она не жертвуя скоростью ради простоты, обеспечивая и то и другое.

PyTorch лучше всего работает как базовая библиотека низкого уровня, представляющая основные операции для высокуровневой функциональности. fastai же — это наиболее популярная библиотека для добавления этой высокуровневой функциональности поверх PyTorch. Кроме того, она отлично подходит для целей этой книги, потому что уникальным образом предоставляет глубокую многослойную архитектуру ПО (об этом многослойном API есть даже научная статья (<https://oreil.ly/Uo3GR>)). По мере освоения глубокого обучения мы будем также погружаться в слои fastai. Эта книга рассказывает о ее второй версии,

являющейся полностью переработанным вариантом со множеством уникальных возможностей.

Не имеет значения, какое ПО вы изучите, — для переключения с одной библиотеки на другую требуется всего несколько дней. По-настоящему же важным является изучение основ глубокого обучения и его техник. Мы сосредоточимся на коде, который максимально отчетливо выражает необходимые для освоения принципы. При объяснении высокоуровневых принципов мы будем использовать высокоуровневый код `fastai`. Там же, где будет рассказываться о низкоуровневых принципах, будем использовать низкоуровневый PyTorch или даже чистый код на Python.

Хотя может показаться, что новые библиотеки глубокого обучения появляются достаточно часто, вам нужно быть готовыми к существенному росту темпа изменений в ближайшие месяцы и годы. Поскольку все больше людей приходят в эту область, они будут привносить свои навыки и идеи, стремясь разработать больше новых продуктов. Следует понимать, что какие бы конкретные библиотеки и ПО вы ни изучали сегодня, они устареют буквально через год или два. Только подумайте о количестве изменений в библиотеках и наборах технологий, которые постоянно происходят в мире веб-программирования, — гораздо более зрелой и медленно растущей области, чем глубокое обучение. Мы твердо уверены, что при освоении этой сферы нужно уделять внимание как пониманию и применению ее основных техник, так и развитию способности быстро осваивать новые техники с инструментами по мере их появления.

К концу книги вы поймете практически весь код, лежащий внутри `fastai` (и большую часть PyTorch тоже), потому что в каждой главе мы будем уходить все глубже и глубже, показывая вам, что именно происходит при построении и обучении моделей. Это означает, что вы освоите наиболее важные практики, используемые в современном глубоком обучении. При этом вы узнаете не только о способах их использования, но и о принципах их работы. Если затем вы захотите использовать эти подходы в другом фреймворке, то для этого у вас уже будут необходимые знания.

Поскольку написание кода и эксперименты — это важнейшая часть освоения глубокого обучения, потребуется хорошая платформа для работы с кодом. Наиболее популярная — Jupyter (<https://jupyter.org/>). Именно ее мы и будем использовать в книге. Мы покажем, как с ее помощью обучать модели и экспериментировать с ними, рассмотрим каждый этап конвейера предварительной обработки данных и разработки модели. Jupyter неспроста считается самым популярным инструментом для научной работы с данными в Python. Он мощный, гибкий и легкий в использовании. Уверены, вы его полюбите!

Давайте же посмотрим его в деле, обучив нашу первую модель.

Ваша первая модель

Мы уже говорили, что сначала научим вас создавать что-либо, а потом расскажем, как это работает. Мы начнем с обучения классификатора изображений распознаванию кошек и собак. Для обучения этой модели и экспериментов с ней вам понадобится сделать кое-какую начальную настройку. Не волнуйтесь, это совсем не сложно.



СЛОВО СИЛЬВЕЙНУ

Не пропускайте этап настройки, даже если сперва он покажется пугающим, особенно если у вас мало или совсем нет опыта в работе с терминалом или командной строкой. Большая часть из всего этого необязательна, и вы увидите, что простейшие серверы можно настроить с помощью обычного браузера. Но для успешного обучения параллельно с книгой нужно пробовать экспериментировать.

Настройка сервера глубокого обучения на GPU

Для работы вам понадобится доступ к компьютеру с NVIDIA GPU (к сожалению, другие производители GPU не полностью поддерживаются основными библиотеками глубокого обучения). Но мы не призываем вас покупать именно такое устройство. Даже если оно у вас уже есть, мы пока не предлагаем его использовать. Для настройки компьютера потребуется время и силы, а сейчас вам нужно все их сосредоточить на освоении глубокого обучения. Поэтому мы предлагаем вам арендовать компьютер, где все необходимое уже установлено и готово к работе.



ТЕРМИН: ГРАФИЧЕСКИЙ ПРОЦЕССОР (GPU)

Иначе называется графическим адаптером, или видеокартой. Особый вид процессора, который может обрабатывать тысячи отдельных задач одновременно. Разработан, в частности, для отображения 3D-сред в видеоиграх. Выполняемые им базовые задачи очень похожи на те, с которыми работают нейросети, в связи с чем GPU может обучать нейросети в сотни раз быстрее обычного CPU. GPU установлен в каждом современном компьютере, но лишь в некоторых используется подходящий для глубокого обучения.

Предпочтительные GPU-серверы будут со временем меняться, так как компании появляются и исчезают, и цены тоже не стоят на месте. Мы поддерживаем актуальность списка рекомендуемых вариантов на сайте книги (<https://book.fast.ai>), так что можете перейти туда прямо сейчас и выполнить инструкции для подключения к предпочтительному GPU-серверу глубокого обучения. Не волнуйтесь: потребуется всего около двух минут, чтобы выполнить настройку

для большинства платформ. При этом в ряде случаев даже не придется платить или использовать кредитную карту.



СЛОВО АЛЕКСИСУ

Внесу свою лепту: прислушайтесь к этому совету! Если вы любите компьютеры, то наверняка захотите настроить собственный. Будьте осторожны! Это реально, но сильно затягивает и отвлекает. Неспроста мы не назвали эту книгу «Все, что вы хотели знать об администрировании системы Ubuntu, установке драйверов Nvidia, apt-get, conda, pip и конфигурации Jupyter Notebook». Собрав и развернув производственную инфраструктуру ML на работе, я могу гарантировать, что она справляется со своей задачей, но при этом не связана с моделированием, как обслуживание самолета не связано с управлением его полетом.

Каждый представленный на сайте вариант содержит руководство, проследовав которому вы увидите экран как на рис. 1.2.

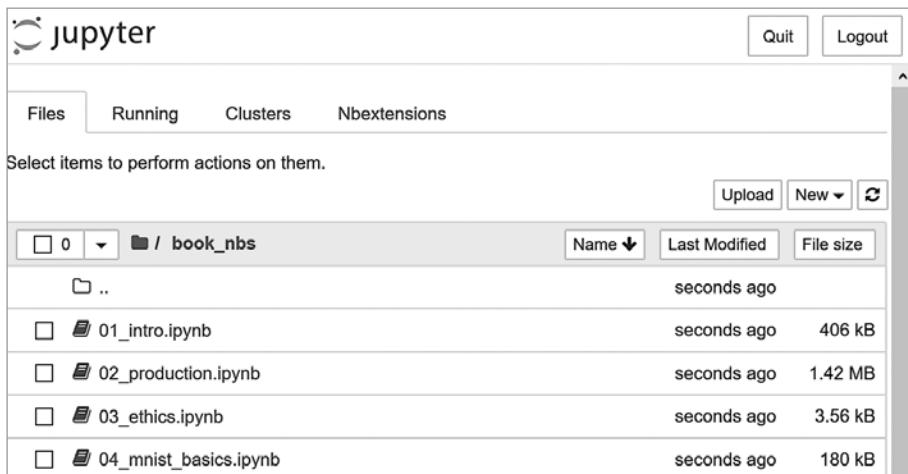


Рис. 1.2. Начальное представление Jupyter Notebook

Теперь вы готовы к запуску своего первого блокнота Jupyter!



ТЕРМИН: БЛОКНОТ JUPYTER

Элемент ПО, позволяющий включать форматированный текст, код, изображения, видео и многое другое в один интерактивный документ. Благодаря своему широкому использованию и невероятному влиянию на многие академические области и индустрию, Jupyter получил высочайшую награду ACM Software System. Блокноты Jupyter наиболее активно используются специалистами по работе с данными для разработки моделей глубокого обучения и взаимодействия с ними.

Запуск первого блокнота

Блокноты пронумерованы по главам в том же порядке, в каком представлены в книге. Итак, первым в списке вы увидите блокнот, который нужно использовать сейчас. С его помощью вы обучите модель, способную распознавать фотографии собак и кошек. Для этого вам потребуется загрузить датасет собак и кошек, с помощью которого вы и *обучите модель*.

Датасет (набор данных) — это просто куча данных, которыми могут быть изображения, электронные письма, финансовые показатели, звуки или что угодно другое. Существует множество бесплатных датасетов, подходящих для обучения моделей. Одни из них созданы учеными для использования в исследованиях, другие сделаны доступными для конкурсов (среди учёных, изучающих данные, проводятся турниры, в которых участники соревнуются, чья модель окажется наиболее точной), а третья являются побочными продуктами других процессов, например подготовки финансовой отчетности.



ЗАПОЛНЕННЫЕ (FULL) И ПУСТЫЕ (STRIPPED) БЛОКНОТЫ

Есть два каталога, в которых находятся разные версии блокнотов. Каталог full содержит блокноты, использованные для создания этой книги. В них прописано все содержание и выводы. Пустая версия содержит те же заголовки и ячейки кода, но все выводы и содержание из них удалены. После прочтения раздела книги мы рекомендуем закрывать ее и прорабатывать именно пустые блокноты, чтобы проверить, удастся ли вам выяснить, что покажет каждая ячейка, прежде чем вы ее выполните. Также старайтесь вспомнить, что тот или иной код демонстрирует.

Откройте блокнот, просто щелкнув на нем. В открытом виде он будет выглядеть примерно так, как показано на рис. 1.3 (обратите внимание, что могут быть небольшие отличия в зависимости от платформы, но можно не обращать на них внимания).

Блокнот состоит из *ячеек*, которые делятся на два основных вида:

- содержащие форматированный текст, изображения и т. д. Они созданы с помощью языка разметки Markdown, о котором вы можете прочесть в приложении A;
- включающие код, который можно выполнить, вследствие чего вывод тут же отобразится внизу (им может быть простой текст, таблицы, изображения, анимация, звуки или даже интерактивные приложения).

Блокноты Jupyter могут находиться в одном из двух режимов: режим редактирования или режим команд. В первом режиме при наборе с клавиатуры знаки вводятся в ячейки обычным способом. А вот в режиме команд вы уже не увидите мигающего курсора, и у каждой клавиши на клавиатуре будет своя функция.

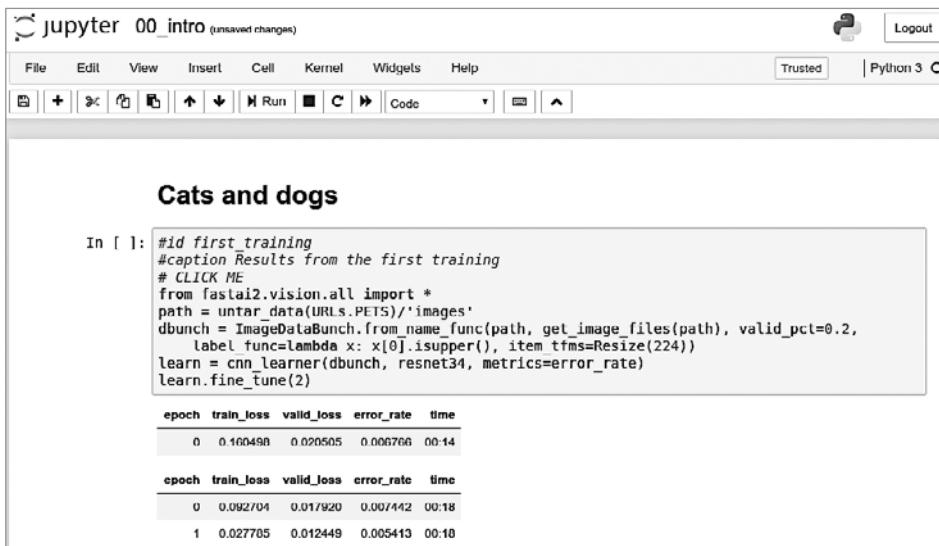


Рис. 1.3. Блокнот Jupyter

Прежде чем продолжить, нажмите клавишу Esc для перехода в режим команд (если вы уже находитесь в этом режиме, ничего не произойдет, поэтому нажмите на всякий случай). Вызвать полный список доступных функций можно нажатием клавиши H. Для выхода из этого вспомогательного экрана нажмите Esc. Обратите внимание, что в отличие от большинства программ в этом режиме выполнение команд не требует удержания Ctrl, Alt и т. п. — вы просто нажимаете нужную буквенную клавишу.

Скопировать ячейку можно нажатием клавиши C (сначала потребуется щелкнуть на этой ячейке, после чего она будет выделена рамкой; если она еще не выбрана, выберите ее). Для вставки скопированной ячейки нажмите клавишу V.

Щелкните на ячейке, начинающейся со строки `# CLICK ME`, чтобы выбрать ее. Первый знак в этой строке указывает, что далее следует комментарий Python, поэтому при выполнении ячейки он игнорируется. Остальная часть ячейки — это, как ни странно, полноценная система для создания и обучения эталонной модели, которая будет распознавать кошек и собак. Итак, приступим к ее обучению! Для этого просто нажмите Shift+Enter или кнопку Play на панели инструментов. Затем подождите несколько минут, в течение которых произойдет следующее.

1. Из коллекции датасетов fast.ai на используемый вами GPU-сервер будет загружен, а затем извлечен датасет под названием Oxford-IIIT Pet Dataset (https://oreil.ly/c_4Bv), содержащий 7349 изображений кошек и собак 37 пород.

2. Из интернета будет загружена *предварительно обученная* на 1,3 млн изображений модель.
3. Эта модель будет *скорректирована* с помощью последних достижений в технологии переноса обучения (transfer learning) для получения модели, распознающей собак и кошек.

Первые два шага выполняются только один раз на GPU-сервере. Если вы выполните ячейку повторно, она использует уже загруженные датасет и модель. Заглянем в содержимое ячейки и ее результаты (табл. 1.2):

```
# CLICK ME
from fastai.vision.all import *
path = untar_data(URLs.PETS) / 'images'

def is_cat(x): return x[0].isupper()
dls = ImageDataLoaders.from_name_func(
    path, get_image_files(path), valid_pct=0.2, seed=42,
    label_func=is_cat, item_tfms=Resize(224))
learn = cnn_learner(dls, resnet34, metrics=error_rate)
learn.fine_tune(1)
```

Таблица 1.2. Результаты первого обучения

epoch	train_loss	valid_loss	error_rate	time
0	0.169390	0.021388	0.005413	00:14

epoch	train_loss	valid_loss	error_rate	time
0	0.058748	0.009240	0.002706	00:19

Возможно, ваши результаты будут другими. На обучение влияет множество небольших произвольных вариаций. Однако в этом примере частота ошибок (error rate) обычно не превышает 0.02.



ВРЕМЯ ОБУЧЕНИЯ

В зависимости от скорости вашей нейросети для загрузки предварительно обученной модели и датасета может потребоваться несколько минут. Само же выполнение может занять около минуты. Нередко для обучения моделей из этой книги нужно несколько минут, столько же потребуется и для ваших, поэтому есть смысл найти этому времени хорошее применение, чтобы не тратить его впустую. Например, перейти к чтению следующего раздела, пока модель обучается, или открыть еще один блокнот и поэкспериментировать с кодом.

ЭТА КНИГА БЫЛА НАПИСАНА В БЛОКНОТАХ JUPYTER

Мы написали эту книгу, используя блокноты Jupyter, поэтому почти для каждого графика, таблицы и вычислений мы будем приводить точный код, с помощью которого вы сможете их повторить. По этой же причине вы очень часто будете видеть код сразу за таблицей, картинкой или просто текстом. Весь этот код можно найти на сайте книги (<https://book.fast.ai/>) и самостоятельно поэкспериментировать с выполнением и изменением каждого примера.

Вы только что видели, как в книге выглядит ячейка, выводящая таблицу. А вот пример ячейки, выводящей текст:

```
1+1  
2
```

Jupyter всегда выводит или показывает результат последней строки (если таковая присутствует). Например, вот образец ячейки, выводящей изображение:

```
img = PILImage.create('images/chapter1_cat_example.jpg')  
img.to_thumb(192)
```



Как же нам узнать, что модель работает хорошо? В последнем столбце таблицы указана *частота ошибок*, то есть доля неверно распознанных изображений. Частота ошибок служит в качестве метрики — полноценной и интуитивно понятной меры качества модели. В данном случае видно, что модель близка к идеалу, даже несмотря на то, что время обучения составило всего несколько секунд (без учета времени на загрузку датасета и предварительно обученной модели). В действительности еще десять лет назад такой точности не удавалось достичь никому.

В завершение убедимся, что эта модель на самом деле работает. Сделайте фото собаки или кошки. Если рядом с вами нет этих животных, то просто скачайте изображение из Google. Теперь выполните ячейку с определенным `uploader`. Это выведет кнопку для выбора изображения, которое требуется классифицировать.

```
uploader = widgets.FileUpload()  
uploader
```

 Upload (0)

Теперь можете передать загруженный файл в модель. Убедитесь, что это отчетливый снимок собаки или кошки, а не рисунок, мультфильм или нечто подобное. Блокнот сообщит вам, узнает ли он на снимке собаку или кошку и насколько он уверен в своем выборе. Надеемся, что ваша модель отлично справится с задачей:

```
img = PILImage.create(uploader.data[0])
is_cat,_,probs = learn.predict(img)
print(f"Is this a cat?: {is_cat}.")
print(f"Probability it's a cat: {probs[1].item():.6f}")
```

```
Is this a cat?: True.
Probability it's a cat: 0.999986
```

Поздравляем вас с созданием первого классификатора!

Но что это значит? Что вы на самом деле сделали? А теперь несколько отвлечемся и взглянем на картину обобщенно.

Что такое машинное обучение

Ваш классификатор — это модель глубокого обучения. Как уже говорилось, такие модели используют нейронные сети, которые ведут свою историю от 50-х годов и только недавно раскрыли свой мощный потенциал благодаря последним прорывам в их исследовании.

Не менее важной деталью контекста является и то, что глубокое обучение — это всего лишь новая область, относящаяся к более общей дисциплине: машинному обучению. Для понимания сути происходящего в процессе обучения вашей модели классификации вам не нужно понимать само глубокое обучение. Достаточно увидеть, что ваша модель и процесс ее обучения выражают принципы, характерные для машинного обучения в целом.

Так что в этом разделе мы расскажем именно о машинном обучении. Мы познакомим вас с его ключевыми понятиями и проследим их до момента первого появления в научном эссе.

Машинное обучение, как и обычное программирование, — это способ заставить компьютеры выполнять поставленную задачу. Но как же нам использовать стандартное программирование для выполнения только что решенной нами задачи: распознавания собак и кошек по фотографиям? Нам понадобилось бы прописать для компьютера точные шаги, необходимые для ее решения.

Обычно при написании программы несложно указать такие шаги для выполнения нужных действий. Мы просто представляем себе, какие шаги нужно проделать, если решать задачу самим, а затем переводим их в код. Например, можно написать функцию, сортирующую список. Как правило, для такой задачи мы пишем функцию, структура которой изображена на рис. 1.4, где *входные*

данные могут быть представлены неупорядоченным списком, а результаты — упорядоченным.

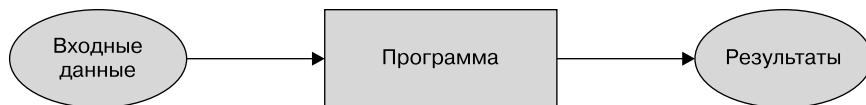


Рис. 1.4. Стандартная программа

Но для распознавания объектов на фото сделать это уже не так просто. Какие мы в таком случае предпринимаем шаги для распознавания? Мы не знаем этого, потому что весь процесс происходит у нас в голове неосознанно!

Еще на заре вычислительных технологий в 1949 году исследователь из IBM по имени Артур Сэмюэл (Arthur Samuel) начал разрабатывать иной подход к постановке компьютерам задач, который он назвал *машинным обучением*. В своем классическом эссе 1962 года *Artificial Intelligence: A Frontier of Automation* («Искусственный интеллект: передовая автоматизация») он писал:

Программирование компьютера для подобных вычислений — тяжелая задача. И не по причине сложности самого компьютера, а из-за необходимости прописывать ежеминутные шаги процесса в исчерпывающих подробностях. Любой программист вам скажет, что компьютеры — это не гигантские мозги, а гигантские тузики.

Его основной идеей было, что вместо того, чтобы сообщать компьютеру конкретные шаги решения задачи, ему нужно показывать примеры ее решения, позволяя определить, как ее решить, самому. Такой способ оказался очень эффективным: к 1961 году его программа игры в шашки обучилась до такой степени, что победила чемпиона Коннектикута! Вот как он описывал свою идею (из того же эссе):

Предположим, мы организуем некоторые автоматические средства проверки эффективности любого текущего назначения веса с позиции фактической производительности и предоставим механизм изменения назначения веса для максимального увеличения этой производительности. Нам не нужно углубляться в детали такой процедуры, чтобы увидеть, что она полностью автоматизирована, и понять, что машина, запрограммированная таким образом, будет «учиться» на собственном опыте.

В этом коротком заявлении кроется ряд важнейших принципов:

- идея «назначения веса»;
- то, что каждое назначение веса имеет определенную «фактическую производительность»;

- необходимость наличия «автоматических средств» тестирования этой производительности;
- потребность в «механизме» (то есть еще одном автоматическом процессе) для повышения производительности через изменение назначений веса.

Рассмотрим все эти принципы по очереди, чтобы увидеть, как они соотносятся друг с другом на практике. Для начала нам нужно понять, что Сэмюэл подразумевал под *назначением веса*.

Веса — это просто переменные, и их назначение — присвоение им конкретных значений. Входные данные программы — значения, которые она обрабатывает с целью получения их результатов: например, получая на входе пиксели изображения и возвращая в виде результата классификацию *dog*. Назначения весов программы — это уже другие значения, которые определяют, как программа будет выполняться.

Поскольку они влияют на выполнение программы, в некотором смысле их можно назвать другим видом входных данных. Учитывая это, мы изменим рис. 1.4 и получим рис. 1.5.

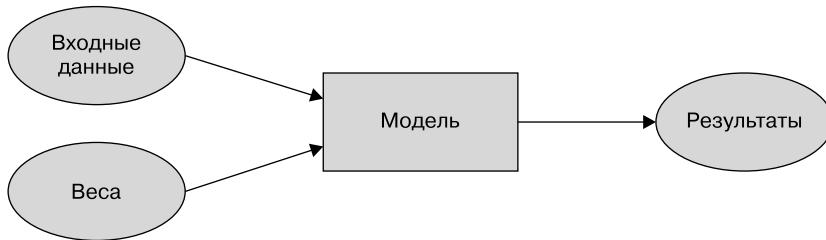


Рис. 1.5. Программа, использующая назначение весов

Мы изменили имя блока с *программа* на *модель*, чтобы следовать современной терминологии и отразить *модель* как особый вид программы: тот, что может решать много разных задач, количество и разнообразие которых зависит от весов. Ее можно реализовать многими способами. Например, в программе шашек Сэмюэла разные значения весов приводили к использованию различных стратегий игры.

(Кстати, то, что Сэмюэл называл «весами», теперь принято называть *параметрами* модели, а термин «веса» выделен для конкретного типа этих параметров.)

Затем Сэмюэл сказал, что *нам нужны автоматические средства тестирования эффективности любого текущего назначения веса с позиции фактической производительности*. В случае программы, играющей в шашки, «фактическая производительность» модели отражалась бы тем, как хорошо она играет. При этом

вы могли бы автоматически тестировать производительность двух моделей, настроив их играть друг против друга и наблюдая, какая из них обычно выигрывает.

Наконец, он говорит, что нам нужен *механизм изменения распределения веса для максимального повышения производительности*. Например, мы можем посмотреть на разницу в весах между побеждающей и проигрывающей моделью, а затем немного скорректировать эти веса для выигрыша.

Теперь мы можем понять, почему он сказал, что такая процедура *может быть сделана полностью автоматической и... машина, запрограммированная таким образом, будет «учиться» на собственном опыте*. Как только коррекция весов станет автоматической, обучение также станет полностью автоматическим. В этом случае уже не нам придется обучать модель, подстраивая ее веса вручную, а автоматизированному механизму, который произведет ее тонкую настройку на основе производительности.

На рис. 1.6 показана полная картина идеи Сэмюэла по обучению модели ML.

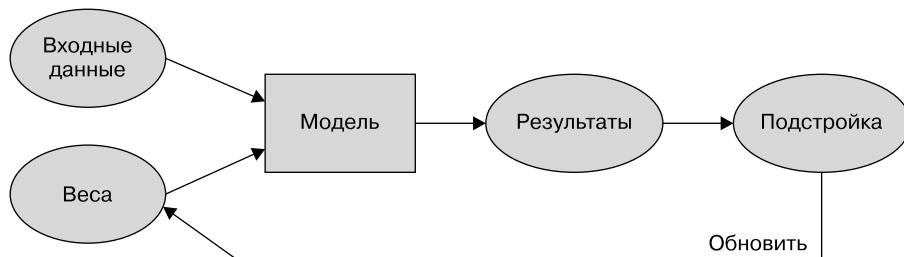


Рис. 1.6. Обучение модели ML

Обратите внимание на различие между *результатами* модели (например, ходами в шашках) и ее *производительностью* (например, выигрывает ли она игру или насколько быстро она ее выигрывает).

Также заметьте, что когда обучение модели завершено (то есть мы утвердили одно наилучшее итоговое назначение веса), тогда уже можно рассматривать веса как *часть модели*, поскольку дальше мы их не изменяем.

Следовательно, схема фактического использования модели после завершения ее обучения выглядит, как показано на рис. 1.7.

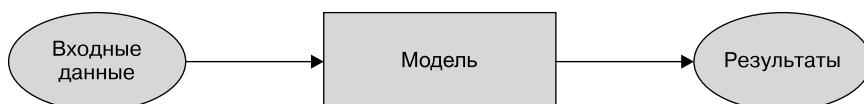


Рис. 1.7. Использование обученной модели как программы

Теперь она идентична начальной схеме на рис. 1.4, только слово *программа* заменено на *модель*. Это важное заключение: *обученную модель можно рассматривать как стандартную компьютерную программу*.



ТЕРМИН: МАШИННОЕ ОБУЧЕНИЕ

Обучение программ, подразумевающее самообучение компьютера на собственном опыте, а не ручное программирование всех отдельных шагов.

Что такое нейронная сеть

Не так уж сложно представить, как может выглядеть модель для программы по обучению игре в шашки. В ней может быть прописано несколько стратегий игры и определенный механизм поиска, после чего с помощью весов может выбираться определенная стратегия, область доски для выполнения поиска и т. д. Однако уже гораздо сложнее представить модель для распознавания изображений, понимания текста или многих других интересных задач.

Для этого нам понадобится определенная функция, которая будет достаточно гибкой, чтобы использоваться для решения любой поставленной задачи простым изменением весов. Удивительно, но такая функция есть! Это нейронная сеть, которую мы уже обсуждали. Таким образом, если рассматривать нейросеть как математическую функцию, то она оказывается чрезвычайно гибкой функцией, регулируемой весами. Математическое доказательство, известное как теорема об универсальной аппроксимации, показывает, что эта функция теоретически может решать любую задачу, достигая любого уровня точности. Факт такой гибкости нейросетей показывает, что на практике они часто оказываются подходящим видом модели, и вы можете сосредоточить свои усилия на их обучении, то есть на поиске удачных назначений веса.

Как же будет выглядеть этот процесс? Можно представить, что потребуется найти новый «механизм» для автоматического обновления веса для каждой задачи. Но это было бы очень утомительно. Нам же здесь нужен абсолютно общий способ обновления весов нейронной сети, чтобы она совершенствовалась при выполнении любой поставленной задачи. Удобно, что такой способ тоже есть!

Он называется *стохастическим градиентным спуском (SGD)*. Подробности работы нейросети и SGD мы увидим в главе 4, где также рассмотрим теорему об универсальной аппроксимации. Сейчас же мы обратимся к словам самого Сэмюэла: «Не нужно углубляться в детали этой процедуры, чтобы увидеть, что ее можно сделать полностью автоматической, и понять, что машина, запрограммированная таким образом, будет “учиться” на собственном опыте».



СЛОВО ДЖЕРЕМИ

Не волнуйтесь: ни SGD, ни нейронные сети не представляют математической сложности. Обе эти технологии практически полностью опираются на использование только сложения и умножения (но эти операции они производят в изобилии!). Обычно в ответ на разъяснение этих подробностей наши студенты с удивлением спрашивают: «И это все?»

Короче говоря, нейросеть — это конкретный вид модели ML, которая как раз соответствует изначальной концепции Сэмюэла. Особенность таких сетей в их повышенной гибкости, которая позволяет им решать необычайно широкий спектр задач простым нахождением подходящих весов. Это мощный инструмент, потому что стохастический градиентный спуск предоставляет нам способ находить значения этих весов автоматически.

Теперь уменьшим масштаб и вернемся к нашей задаче по классификации изображений с помощью схемы Сэмюэла.

Входными данными у нас являются изображения, а весами — веса в нейронной сети. Наша модель — это нейронная сеть, а результаты — значения, которые она вычисляет, например *dog* или *cat*.

А что насчет следующего элемента: *автоматического средства тестирования эффективности любого текущего назначения веса с позиции производительности?* Определить «фактическую производительность» достаточно легко: можно просто определить производительность нашей модели как точность ее прогнозирования правильных ответов.

Если все это объединить и предположить, что SGD является механизмом для обновления назначений весов, то становится очевидно, что классификатор изображений — это модель машинного обучения, как и предполагал Сэмюэл.

Немного терминологии глубокого обучения

Сэмюэл работал в 1960-е годы, и терминология с тех пор изменилась. Далее приведены современные термины глубокого обучения для всех рассмотренных нами элементов.

- Функциональная форма *модели* называется *архитектурой* (но будьте внимательны: иногда люди используют термин *модель* как синоним *архитектуры*, поэтому может возникнуть путаница).
- *Веса* теперь называются *параметрами*.
- *Прогнозы* вычисляются на основе *независимой переменной*, которой являются *данные*, не включающие *метки*.
- *Результаты* модели называются *прогнозами*.

- *Потери* — мера производительности.
- Потери зависят не только от прогнозов, но и от верности *меток* (также называемых *целями* или *зависимыми переменными*), например *dog* или *cat*.

На рис. 1.8 изображена обновленная схема с рис. 1.6 после внесения изменений.

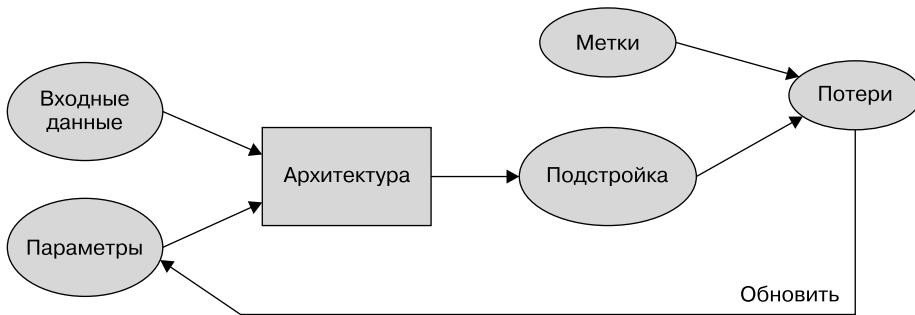


Рис. 1.8. Подробный цикл обучения

Характерные для ML ограничения

На основе этой схемы можно сделать фундаментальные выводы.

- Модель нельзя создать без данных.
- Модель может научиться работать только с шаблонами, найденными в используемых для ее обучения входных данных.
- Такой подход к обучению создает только *прогнозы*, но не рекомендуемые *действия*.
- Недостаточно иметь только образцы входных данных. Для этих данных также необходимы метки (например, картинок собак и кошек недостаточно для обучения модели; для каждой из них нужна метка, поясняющая, на какой из них собака, а на какой — кошка).

Исходя из наблюдений, мы заметили, что большинство организаций, заявляющих о недостатке данных, имеют в виду, что недостает именно *маркированных* данных. Если какая-либо организация заинтересована в выполнении какой-то задачи с помощью модели, то скорее всего, у них уже есть некоторые входные данные, в отношении которых они собираются применять эту модель. Кроме того, они наверняка уже выполняли эту задачу каким-то другим способом (например, вручную или с помощью эвристической программы), следовательно, от этих процессов у них тоже остались данные. Например, в практической радиологии наверняка будет архив медицинских снимков (поскольку они сохраняются для отслеживания изменений состояния пациентов), но эти снимки

могут не иметь структурированных меток, содержащих список диагнозов или вмешательств (так как радиологи обычно создают отчеты в произвольной форме, а не структурированными данными). Мы будем много говорить о подходах к маркировке данных, потому что для глубокого обучения этот аспект очень важен.

Поскольку такие виды моделей ML могут только делать *прогнозы* (то есть стараться воспроизвести метки), это может привести к существенному разрыву между целями организации и возможностями модели. Например, вы узнаете, как создавать *рекомендательную систему*, способную прогнозировать, какие товары может приобрести пользователь. Эта система часто используется в коммерции, например, для настройки показа товаров на главной странице исходя из их популярности. Однако такая модель обычно создается на основе изучения пользователя и его истории покупок (*входные данные*), а также того, что он собирается купить или что просматривал (*метки*). Это означает, что модель с большей вероятностью сообщит вам о товарах, которые у пользователя уже есть или о которых он уже знает, а не о новых, которые могли бы его реально заинтересовать. Это противоположно тому, что делает, например, эксперт в книжном магазине, который задает вопросы для выяснения вашего вкуса и уже затем рассказывает об авторах или сериях, о которых вы никогда не слышали.

Еще один важный нюанс можно заметить при рассмотрении того, как модель взаимодействует со средой. В результате этого взаимодействия могут создаваться *петли обратной связи*.

1. *Предиктивная полицейская* модель создается на основе данных о местах прошлых арестов. На практике это скорее предсказывает не преступление, а арест, что отчасти отражается в смещениях, наблюдаемых в текущей работе полиции.
2. Полицейские могут использовать эту модель для выбора приоритетных мест наблюдения, что, естественно, приводит к повышению числа арестов в этих местах.
3. Данные об этих дополнительных арестах в дальнейшем будут использоваться для повторного обучения будущих версий модели.

Это пример *положительной петли обратной связи*: чем больше используется модель, тем более необъективными становятся данные, что приводит ко все большему искажению модели, и т. д.

Петли обратной связи могут также создавать проблемы в коммерческой сфере. К примеру, рекомендательная система видео может быть смешена в сторону рекомендаций контента, популярного среди наиболее активных зрителей (например, любителей теорий заговоров и экстремистов), что приводит к повышению просмотра этими пользователями аналогичных видеоматериалов, в связи с чем такие материалы рекомендуются все больше и больше. Мы рассмотрим эту тему подробнее в главе 3.

Теперь, когда вы ознакомились с теорией, пора вернуться к нашему примеру кода и детальнее проанализировать элементы, соответствующие только что описанным процессам.

Как работает наш распознаватель изображений

Давайте проанализируем, как с описанными идеями соотносится код нашего распознавателя изображений. Мы поместим строчки в отдельные ячейки и посмотрим, что каждая из них делает (мы пока не будем объяснять каждую деталь каждого параметра, но опишем самые важные части, а подробности оставим на потом). Первая строка импортирует библиотеку fastai.vision целиком:

```
from fastai.vision.all import *
```

В итоге мы получаем все функции и классы, которые потребуются для создания моделей компьютерного зрения.



СЛОВО ДЖЕРЕМИ

Многие Python-программисты советуют избегать импорта всей библиотеки, подобной этой (используя синтаксис `import *`), потому что в крупных проектах это может вызвать проблемы. Однако для интерактивной работы, как, например, с блокнотом Jupyter, это подходит отлично. Библиотека fastai специально спроектирована для поддержки подобного вида интерактивного использования и импортирует в вашу среду только необходимые элементы.

Вторая строка загружает на ваш сервер стандартный датасет из коллекции fast.ai (<https://course.fast.ai/datasets>) (если он не загружался ранее), извлекает его (если он еще не извлекался) и возвращает объект `Path` с путем к месту извлечения:

```
path = untar_data(URLs.PETS) / 'images'
```



СЛОВО СИЛЬВЕЙНУ

Во времена обучения в fast.ai да и по сей день я многое узнаю о продуктивных практиках написания кода. Библиотека fastai и блокноты fast.ai содержат множество маленьких приемов, которые помогли мне повысить свои навыки программирования. В качестве примера можно даже обратить внимание на то, что библиотека fastai возвращает не просто строку, содержащую путь к датасету, но целый объект `Path`. Это очень полезный класс из стандартной библиотеки Python 3, который существенно упрощает доступ к файлам и каталогам. Если вы с ним еще не сталкивались, то обязательно ознакомьтесь с соответствующей документацией или руководством и попробуйте его в деле. Имейте в виду, что сайт книги предоставляет ссылки на рекомендемые материалы для каждой главы. Я продолжу обращать ваше внимание на советы по программированию, которые сам нашел полезными, по мере их появления в книге.

В третьей строке мы определяем функцию `is_cat`, которая ставит метки кошек на основе правила имени файла, предоставленного создателями датасета:

```
def is_cat(x): return x[0].isupper()
```

Мы используем эту функцию в четвертой строке, которая сообщает `fastai` о нашем датасете и его структуре:

```
dls = ImageDataLoaders.from_name_func(
    path, get_image_files(path), valid_pct=0.2, seed=42,
    label_func=is_cat, item_tfms=Resize(224))
```

Существуют различные классы, предназначенные для разных видов датасетов и задач, но здесь мы используем `ImageDataLoaders`. Первая часть имени класса обычно указывает на тип ваших данных, например `image` (изображение) или `text` (текст).

Еще один важный элемент информации, который нам нужно сообщить `fastai`, — это способ получения меток из датасета. Датасеты компьютерного зрения обычно структурированы так, что метка для изображения является частью имени файла или пути: чаще всего имени родительского каталога. `fastai` предлагает ряд стандартизованных методов разметки и возможности написания собственных. Здесь мы даем `fastai` команду использовать функцию `is_cat`, которую только что написали.

В завершение мы определяем нужные элементы `Transforms`. `Transform` содержит код, который автоматически применяется в процессе обучения. `fastai` включает множество предопределенных `Transforms`, а добавить новые так же просто, как создать функцию Python. Существует два их вида: `item_tfms` применяются к каждому элементу (в данном случае размер каждого элемента изменяется до квадрата размером 224 пикселя), а `batch_tfms` применяются сразу к набору элементов с помощью GPU, что делает их особенно быстрыми (в книге мы увидим много таких примеров).

Почему именно 224 пикселя? Так сложилось исторически — старые предварительно обученные модели требуют именно такой размер, но вы можете передать любой. Если этот размер увеличить, то зачастую вы получите модель с лучшими результатами (поскольку она сможет фокусироваться на большем числе деталей), но в итоге понизится скорость и увеличится потребление памяти. При понижении размера произойдет обратное.

Датасет домашних животных содержит 7390 картинок собак и кошек, относящихся к 37 породам. Каждое изображение размечено на основе имени его файла: например, файл `great_pyrenees_173.jpg` — это 173-й образец изображения пиренейской горной собаки в датасете. Имена файлов начинаются с верхнего регистра, если на их изображении кошка, и с нижнего регистра в любом другом случае. Нам нужно указать `fastai`, как получать метки из имен файлов, для чего мы используем вызов `from_name_func` (это означает, что метки можно извлекать,

применяя к ним функцию) и передаем `is_cat`, которая принимает значение `True`, если первая буква написана в верхнем регистре (то есть на изображении кошка).



ТЕРМИН: КЛАССИФИКАЦИЯ И РЕГРЕССИЯ

Классификация и *регрессия* в машинном обучении имеют очень специфичные значения. Они описывают два основных типа модели, которые мы будем изучать на протяжении книги. *Классификационная модель* — это та, которая старается прогнозировать класс или категорию, то есть делает прогноз из числа дискретных вероятностей, таких как `dog` или `cat`. *Регрессионная модель* — это та, которая старается прогнозировать одну или несколько численных величин, таких как температура или локация. Некоторые используют слово *регрессия* в отношении конкретного вида модели, называемой *моделью линейной регрессии*. Это плохая практика, и мы такую терминологию использовать в книге не будем.

Самый важный параметр здесь — это `valid_pct=0.2`. Он говорит о том, что `fastai` должен удержать 20 % данных и вообще не использовать их для обучения модели. Эти 20 % данных называются *контрольной выборкой*. Оставшиеся же 80 % называются *обучающей выборкой*. Контрольная выборка используется для измерения точности модели. По умолчанию эти 20 % выбираются произвольно. Параметр `seed=42` при каждом выполнении кода устанавливает для случайного начального числа одно и то же значение, и каждый раз мы получаем одинаковую контрольную выборку. Таким образом, если мы изменим нашу модель и повторно ее обучим, то мы будем знать, что любые отличия вызваны изменениями в модели, а не использованием другой контрольной выборки.

`fastai` будет *всегда* показывать точность вашей модели, используя *только* контрольную выборку, но не обучающую. Это критически важно, потому что если вы обучаете достаточно большую модель достаточно длительное время, то в итоге она запомнит метку каждого элемента вашего датасета! В результате получится бесполезная модель — нам важно, чтобы она работала на *ранее не встречавшихся изображениях*. Именно эту цель мы будем преследовать при каждом обучении модели: она должна быть эффективной только для будущих данных, которые встретит после своего обучения.

Даже когда ваша модель еще не полностью запомнила все данные, на ранних этапах обучения она может запомнить их конкретные части. Следовательно, чем дольше процесс обучения, тем выше точность в отношении обучающей выборки. Точность в отношении контрольной выборки тоже будет какое-то время увеличиваться, но в итоге уменьшится, так как модель будет запоминать обучающую выборку, вместо того чтобы находить базовые обобщаемые шаблоны в данных. Когда такое происходит, мы говорим, что модель *переобучилась*.

На рис. 1.9 показано, что происходит, если переборщить. Для демонстрации мы использовали упрощенный пример, где применяется только один параметр и не-

много произвольно сгенерированных на основе функции $x^{**}2$ данных. Несмотря на то что прогнозы в переобученной модели точны в отношении данных возле наблюдаемых точек, эта точность утрачивается при выходе за их диапазон.

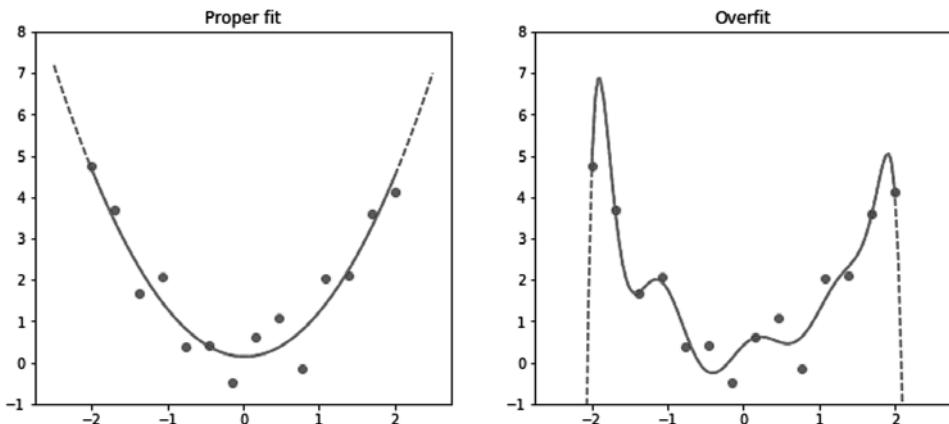


Рис. 1.9. Пример переобучения

Переобучение – это наиболее важная и сложная задача для всех практиков ML и алгоритмов. Как вы позднее увидите, легко создать модель, которая отлично справляется с составлением прогнозов только для тех данных, на которых она обучалась, но при этом куда сложнее делать точные прогнозы для данных, с которыми модель ранее не встречалась. Конечно же, именно такие данные и будут иметь значение на практике. Например, если вы создаете классификатор рукописных цифр (а мы скоро его создадим!) и используете для распознавания написанных на чеках чисел, то вы никогда не увидите таких чисел, на которых модель обучалась, ведь на каждом чеке будут как минимум небольшие вариации написания.

В этой книге мы рассмотрим много методов, позволяющих избежать переобучения. Однако использовать эти методы следует только после того, как вы убедитесь, что переобучение действительно имеет место (то есть если в процессе обучения вы наблюдаете ухудшение контрольной точности). Мы часто видим, что люди используют техники избегания переобучения, даже когда у них достаточно данных, чтобы этого не делать. В итоге их модель получается менее точной, чем могла бы быть.

Пятая строка кода, обучающего наш распознаватель изображений, дает fastai команду создать *сверточную нейронную сеть* (CNN) и указывает, какую архитектуру применять (то есть какой вид модели создать), на каких данных мы собираемся проводить обучение и какую *метрику* использовать:

```
learn = cnn_learner(dls, resnet34, metrics=error_rate)
```



КОНТРОЛЬНАЯ ВЫБОРКА

При обучении модели у вас всегда должна быть и обучающая выборка, и контрольная. Точность нужно измерять только с помощью контрольной. Если вы будете производить обучение слишком долго при недостаточном количестве данных, то обнаружите, что точность модели начнет ухудшаться, что и называется переобучением. `fastai` по умолчанию устанавливает `valid_pct` на 0.2, поэтому, даже если вы забудете это сделать, библиотека создаст контрольную выборку за вас!

Почему CNN? Это эталонный подход к созданию моделей компьютерного обучения. Позднее мы с вами подробно познакомимся с этим видом нейронных сетей, структура которых основана на принципе работы системы зрения человека.

В книге мы познакомим вас со множеством доступных в `fastai` архитектур, а также покажем, как создавать собственную. Но чаще всего выбор архитектуры не будет являться важной частью глубокого обучения. Эту тему больше любят обсуждать ученые умы, но на практике тратить на них много времени не придется. Есть стандартные архитектуры, которые подойдут для большинства ситуаций, и в конкретном случае мы используем *ResNet*, о которой будем говорить очень много. Для многих датасетов и задач она оказывается не только быстрой, но и точной. В `resnet34` число 34 указывает на число слоев в этой версии архитектуры (другие возможные варианты включают 18, 50, 101 и 152). Модели, использующие архитектуры с большим числом слоев, дольше обучаются и более подвержены переобучению (то есть вы не можете выполнять обучение определенное количество эпох до того, как точность в отношении контрольной выборки начнет падать). С другой стороны, при использовании большего объема данных они могут быть немного точнее.

Что такое метрика? *Метрика* — это функция, измеряющая качество прогнозов модели при использовании контрольной выборки, результаты которой выводятся в конце каждой эпохи обучения. В нашем случае мы используем `error_rate` — функцию, предоставляемую `fastai`, которая делает именно то, о чем говорит ее название: сообщает вам процент неверно классифицированных изображений контрольной выборки. Еще одной распространенной метрикой для классификации является `accuracy` (которая вычисляется как $1.0 - \text{error_rate}$). `fastai` предлагает гораздо больше вариантов метрик, которые мы также еще будем рассматривать.

Метрика может вызвать ассоциацию с потерями, но между ними есть важное отличие. Потери — это способ определения «показателя производительности», который обучающая система сможет использовать для автоматического обновления весов. Другими словами, для потерь хорошим вариантом будет тот, который удобно использовать стохастическому градиентному спуску. Метрика же определяется для оценки человеком, поэтому хорошая метрика — это та, которую легко понять вам и которая максимально близка к тому, что ваша модель

должна делать. Иногда вы можете решить, что функция потерь — это подходящая метрика, но это не обязательно так.

`cnn_learner` помимо прочего содержит параметр `pretrained`, который по умолчанию равен `True` (поэтому используется в данном случае, несмотря на то что мы его не указывали), что устанавливает значения весов вашей модели на те, которые уже были обучены экспертами для распознавания тысячи разных категорий среди 1,3 миллиона фотоснимков (при использовании известного датасета ImageNet (<http://www.image-net.org/>)). Модель, которая использует веса, обученные на другом датасете, называется *предварительно обученной моделью*. Вам следует использовать такой вид модели, так как это означает, что еще до знакомства с вашими данными она уже будет обладать рядом возможностей. И как вы увидите, в модели глубокого обучения многие из этих возможностей окажутся нужны независимо от деталей вашего проекта. Например, части предварительно обученных моделей будут обрабатывать границы, градиент и определение цвета, что необходимо для многих задач.

При использовании предварительно обученной модели `cnn_learner` удаляет последний слой, поскольку он всегда специально настроен для исходной задачи обучения (то есть классификации датасета ImageNet), и заменяет его одним или несколькими новыми слоями подходящего размера с произвольными весами для датасета, с которым работаете уже вы. Последняя часть модели называется *головой*.

Использование предварительно обученной модели — важнейший метод из имеющихся, который позволяет нам обучать более точные модели быстрее, используя при этом меньше данных, с меньшими затратами времени и денег. Может показаться, что использование предварительно обученных моделей является наиболее исследуемой областью в академическом глубоком обучении... Но это совершенно не так! Значимость таких моделей практически не признается и не учитывается в большинстве курсов и книг, а также редко рассматривается в академических трудах. На момент написания этой книги, а именно в первой половине 2020 года, ситуация только начинает меняться, но для существенных изменений потребуется какое-то время. Поэтому будьте внимательны: большинство людей, с которыми вы можете заговорить на эту тему, вероятно, будут недооценивать ваши достижения в глубоком обучении с помощью небольшого объема ресурсов, так как сами будут плохо разбираться в применении предварительно обученных моделей.

Использование таких моделей для задач, отличных от тех, для которых модель тренировалась изначально, называется *transfer learning*, то есть *переносом обучения*. К сожалению, из-за недостаточной изученности этой техники предварительно обученные модели доступны лишь для ограниченного числа областей. Например, в медицине таких моделей доступно очень мало, что усложняет использование переноса обучения в этой сфере. Кроме того, пока недостаточно ясно, как использовать эту технику для таких задач, как анализ временных рядов.



ТЕРМИН: ПЕРЕНОС ОБУЧЕНИЯ

Использование предварительно обученной модели для задачи, отличающейся от той, для которой модель обучалась изначально.

Шестая строка нашего кода сообщает fastai, как *тонко настраивать* модель:

```
learn.fine_tune(1)
```

Как мы уже говорили, архитектура описывает только *шаблон* для математической функции, но по факту ничего не делает, пока мы не предоставим значения для миллионов содержащихся в ней параметров.

Это ключевой момент в глубоком обучении — определение того, как настроить параметры модели так, чтобы она решала поставленную задачу. Для тонкой настройки модели нам нужно предоставить хотя бы один элемент информации: сколько раз просматривать каждое изображение (указать так называемое число *эпох*). Выбор числа эпох будет сильно зависеть от того, сколько у вас есть времени и сколько его фактически потребуется для тонкой настройки модели. Если вы выберете слишком маленькое число, то всегда сможете обуздать модель позднее.

Но почему метод называется *fine_tune*, а не *fit*? В fastai есть метод под названием *fit*, который действительно тонко настраивает модель (то есть просматривает изображения обучающей выборки несколько раз, каждый раз обновляя параметры, чтобы максимально приблизить прогнозы к целевым меткам). Но в данном случае мы начали с предварительно обученной модели и не хотим отказываться от всех тех возможностей, которые у нее уже есть. Как вы узнаете из этой книги, есть ряд полезных приемов для адаптации предварительно обученной модели к новому датасету, и этот процесс называется *тонкой настройкой* (*fine-tuning*).



ТЕРМИН: ТОНКАЯ НАСТРОЙКА

Техника переноса обучения, которая обновляет параметры предварительно обученной модели, обучая ее дополнительное количество эпох с помощью задачи, отличающейся от оригинальной.

Когда вы применяете метод *fine_tune*, fastai будет использовать эти приемы за вас. Существует ряд параметров, которые вы можете установить (о чем мы поговорим позже), но в приведенной здесь предустановленной форме библиотека выполняет два шага.

1. Использует одну эпоху для тонкой настройки только тех частей, которые необходимы для корректной работы новой произвольной головы с вашим датасетом.

2. Использует несколько эпох, запрашиваемых при вызове этого метода, чтобы тонко настроить всю модель, обновляя веса более поздних слоев (особенно головы) быстрее, чем ранних (которые, как мы увидим, обычно не требуют больших изменений).

Голова модели — это часть, которая заменяется, чтобы соответствовать новому датасету. *Эпоха* — это один полный проход через датасет. При вызове `fit` после каждой эпохи выводятся результаты, где показан ее номер, потери обучающей и контрольной выборки («показатель производительности», используемый для обучения модели) и любые другие запрашиваемые метрики (в этом случае `err` или `rate`).

Итак, используя весь этот код, наша модель научилась распознавать кошек и собак только по размеченным примерам. Но как мы этого добились?

Чему научился наш распознаватель изображений

На данной стадии у нас есть отлично работающий распознаватель изображений, но мы абсолютно не понимаем, что именно он делает. Хоть многие люди и жалуются, что глубокое обучение порождает непроницаемые модели «черных ящиков» (что-то, делающее прогнозы, но недоступное для понимания), это мнение очень далеко от истины. Существует огромное количество исследований, показывающих, как углубленно анализировать модели глубокого обучения и получать о них полную информацию. Тем не менее все виды моделей машинного обучения (включая модели глубокого обучения и традиционные статистические модели) могут оказываться весьма трудными для понимания, особенно когда мы рассматриваем их потенциальное поведение при получении данных, сильно отличающихся от тех, что использовались в обучении. Эту проблему мы будем обсуждать на протяжении всей книги.

В 2013 году аспирант Мэтт Зейлер (Matt Zeiler) и его научный руководитель Роб Фергус (Rob Fergus) опубликовали работу *Visualizing and Understanding Convolutional Networks* (<https://oreil.ly/iP8cr>) («Визуализация и понимание сверточных сетей»), которая показала, как визуализировать веса нейронной сети, обученные в каждом слое модели. Они внимательно проанализировали модель, победившую на соревновании ImageNet 2012 года, и использовали полученные данные для ее усовершенствования, что позволило им победить на соревновании 2013 года. На рис. 1.10 изображена опубликованная ими картинка весов первого слоя.

Эта картинка требует некоторого пояснения. Для каждого слоя часть изображения со светло-серым фоном показывает воспроизведенные веса, а более крупный набор картинок внизу показывает части обучающих изображений, которые больше других соответствовали каждому набору весов. Для слоя 1 мы видим, что модель нашла веса, представляющие диагональные, горизонтальные

и вертикальные границы, а также различные градиенты. (Обратите внимание, что для каждого слоя показано только подмножество признаков; на практике среди всех слоев их присутствуют тысячи.)

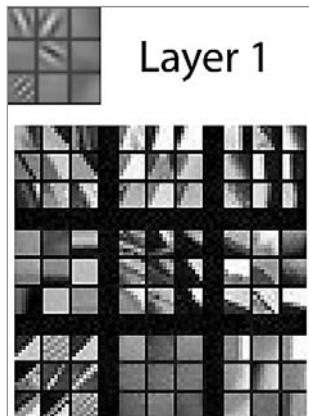


Рис. 1.10. Активации первого слоя CNN (изображение любезно предоставлено Мэттом Д. Зейлером и Робом Фергусом)

Это основные структурные элементы, которые усвоила модель для компьютерного зрения. Их широко анализировали нейробиологи и исследователи компьютерного зрения, выяснив, что эти усвоенные структурные элементы очень похожи на базовые зрительные механизмы человеческого глаза, а также на построенные вручную признаки компьютерного зрения, которые были разработаны еще до появления глубокого обучения. На рис. 1.11 показан следующий слой.

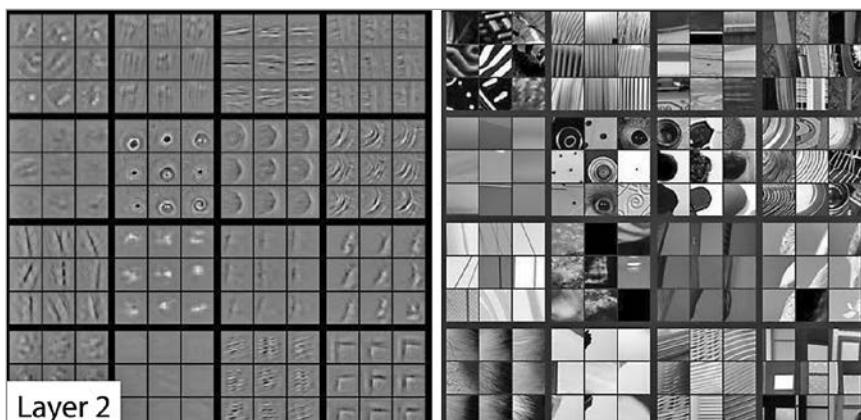


Рис. 1.11. Активации второго слоя CNN (изображение любезно предоставлено Мэттом Д. Зейлером и Робом Фергусом)

Для второго слоя приводится девять образцов воссоздания весов для каждого из найденных в модели признаков. Здесь мы видим, что модель научилась создавать детекторы признаков, которые ищут углы, повторяющиеся линии, круги и другие аналогичные паттерны. Эти признаки создаются из базовых конструктивных элементов, разработанных в первом слое. Для каждого из них правая сторона картинки показывает небольшие участки фактических изображений, с которыми эти признаки больше всего совпадают. Например, конкретный паттерн в первом столбце второго ряда совпадает с градиентами и текстурами, ассоциированными с закатами.

На рис. 1.12 показаны результаты воссоздания признаков третьего слоя.

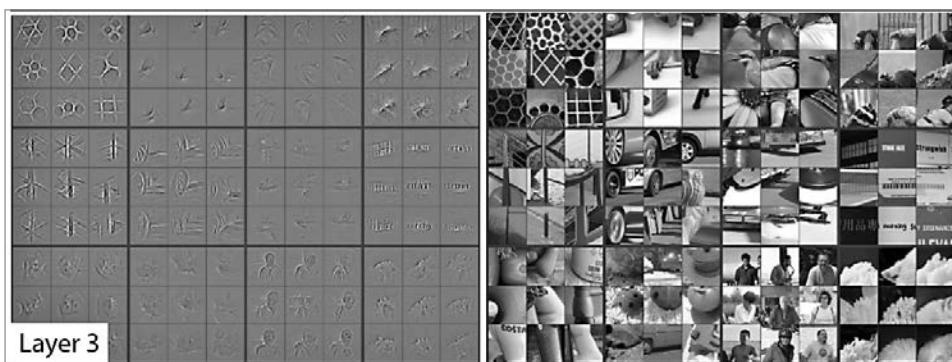


Рис. 1.12. Активации третьего слоя CNN (изображение любезно предоставлено Мэттом Д. Зейлером и Робом Фергусоном)

Как вы видите в правой части картинки, теперь признаки могут идентифицироваться и сопоставляться с высокоуровневыми семантическими компонентами, такими как колеса машины, текст и лепестки цветка. Как показано на рис. 1.13, на основе этих компонентов слои 4 и 5 могут определять еще более высокоуровневые понятия.

В этой работе изучалась более старая модель под названием *AlexNet*, которая содержала только пять слоев. У сетей, разработанных с тех пор, могут быть уже сотни слоев, так что несложно представить, насколько богатыми могут быть разрабатываемые такими моделями признаки.

Когда мы ранее настраивали предварительно обученную модель, то адаптировали то, на чем фокусируются эти последние слои (цветы, людей, животных), специализируя ее для задачи различения собак и кошек. Если говорить более обобщенно, то мы могли бы специализировать предварительно обученную модель для работы со многими разными задачами. Обратимся к нескольким примерам.

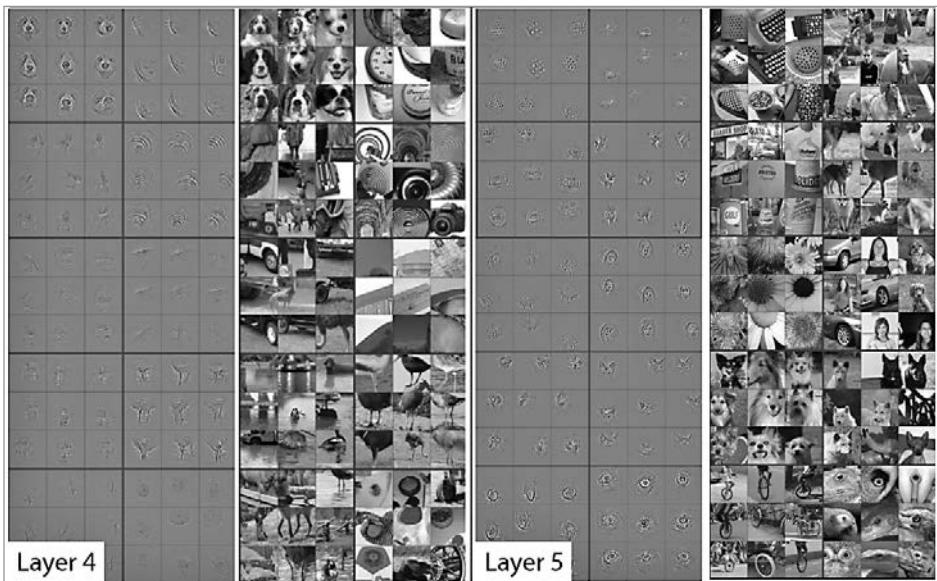


Рис. 1.13. Активации четвертого и пятого слоев CNN (изображение любезно предоставлено Мэттом Д. Зейлером и Робом Фергусом)

Распознаватели изображений для других задач

Как следует из названия, распознаватель изображений способен распознавать изображения. Но с помощью изображений можно представить много разных вещей, а это значит, что такой распознаватель может научиться выполнять множество других задач.

К примеру, звук из аудиофайла можно преобразовать в спектрограмму — график, отражающий насыщенность каждой частоты в каждый момент времени. Студент fast.ai Этан Сутин (Ethan Sutin), используя этот подход, с легкостью превзошел зарегистрированную точность эталонной модели обнаружения звуков окружающей среды (<https://oreli.ly/747uv>), использующей датасет из 8732 городских звуков. Метод `show_batch`, использующийся в fast.ai, отчетливо показывает, что каждый звук имеет весьма характерную спектрограмму (рис. 1.14).

Временной ряд можно легко преобразовать в изображение, просто отразив его на графике. Однако зачастую следует искать способ представить данные наиболее доступным для извлечения самых важных компонентов способом. В случае с временными рядами такие явления, как сезонность и аномалии, представляют наибольший интерес.

Для данных временных рядов доступны различные преобразования. Например, студент fast.ai Игнасио Огуиза (Ignacio Oguiza) создал изображения на основе

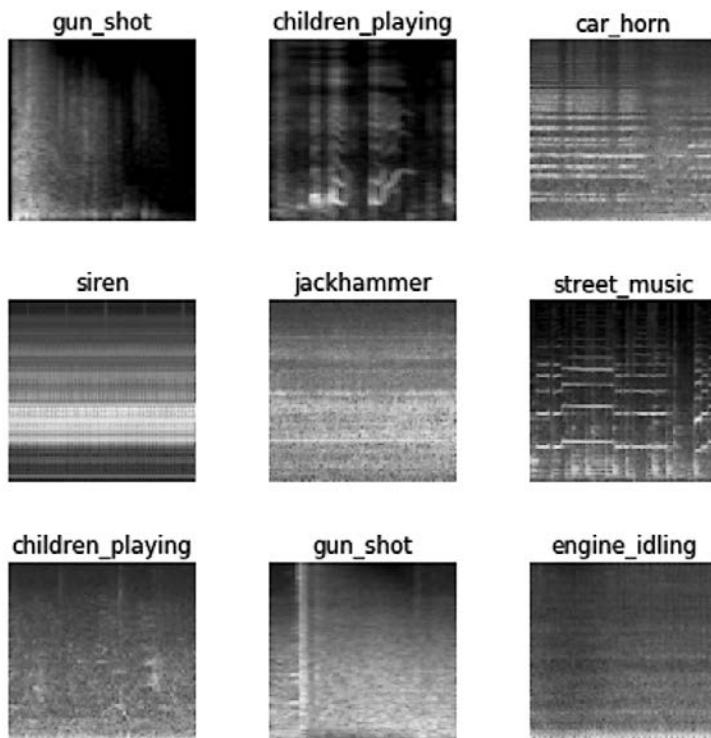


Рис. 1.14. `show_batch` показывает спектрограммы звуков

датасета временных рядов для классификации оливкового масла, используя технику под названием Gramian Angular Difference Field (поле углового смещения матрицы Грэма) (GADF). Результат вы можете посмотреть на рис. 1.15. После этого он отправил полученные изображения в модель классификации, аналогичную той, с которой вы познакомились в данной главе. Итоговые результаты, несмотря на использование в обучающей выборке всего 30 изображений, показали точность выше 90 %, приблизившись к эталонному уровню.

Еще один примечательный проект fast.ai был разработан студентом Глебом Эсманом (Gleb Esman). Он работал над системой обнаружения мошеннических действий в Splunk, используя датасет с движениями и кликами мыши, совершаемыми пользователем. Он преобразовывал их в картинки путем рисования цветными линиями изображения, показывающего позицию, скорость и ускорение указателя мыши, а щелчки при этом отображал мелкими цветными кругами (https://oreil.ly/6-I_X). Все это показано на рис. 1.16. Полученные рисунки он передавал модели распознавания изображений так же, как мы делали с нашими, и сработало это настолько хорошо, что в итоге он даже получил патент на этот подход в аналитике мошенничества.

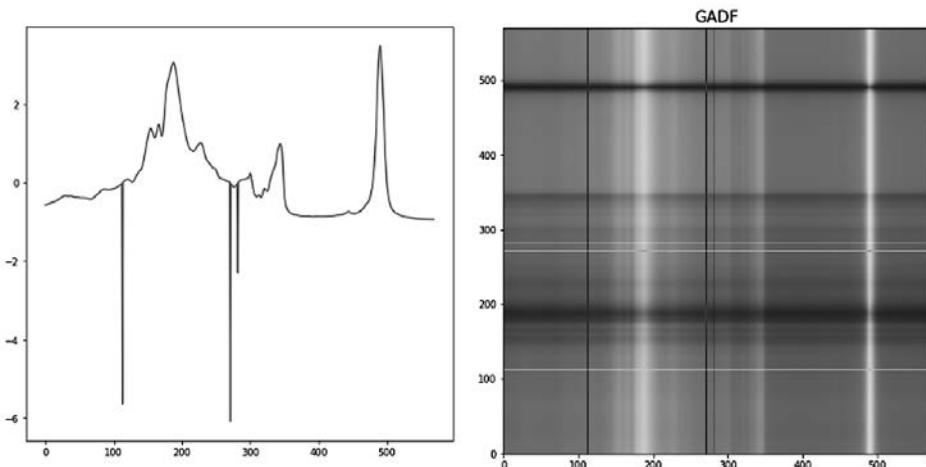


Рис. 1.15. Преобразование временного ряда в изображение

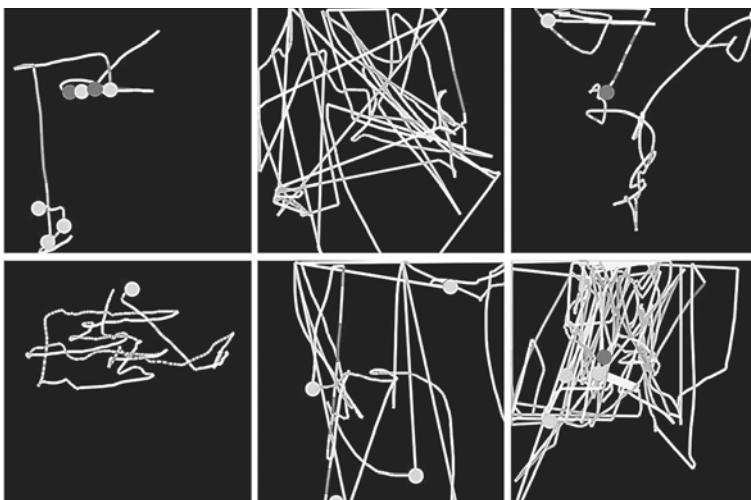


Рис. 1.16. Преобразование поведения компьютерной мыши в изображение

Еще один пример можно привести из работы *Malware Classification with Deep Convolutional Neural Networks* (https://oreil.ly/l_kna) (классификация вредоносного ПО с помощью сверточных нейронных сетей), написанной Махмудом Калашом (Mahmoud Kalash) и др., где объясняется, что «вредоносный бинарный файл разделяется на 8-битные последовательности, которые затем преобразуются в эквивалентные десятичные значения. Этот десятичный вектор изменяется, и создается полутоновое изображение, представляющее образец вредоносного ПО». Схематично этот процесс показан на рис. 1.17.

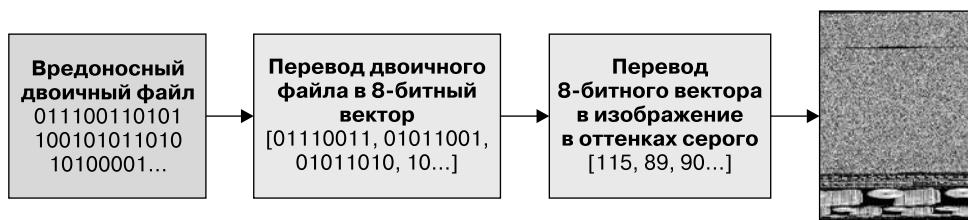


Рис. 1.17. Процесс классификации вредоносного ПО

На рис. 1.18 показаны «картинки» сгенерированных этим процессом образцов вредоносного ПО в различных категориях.

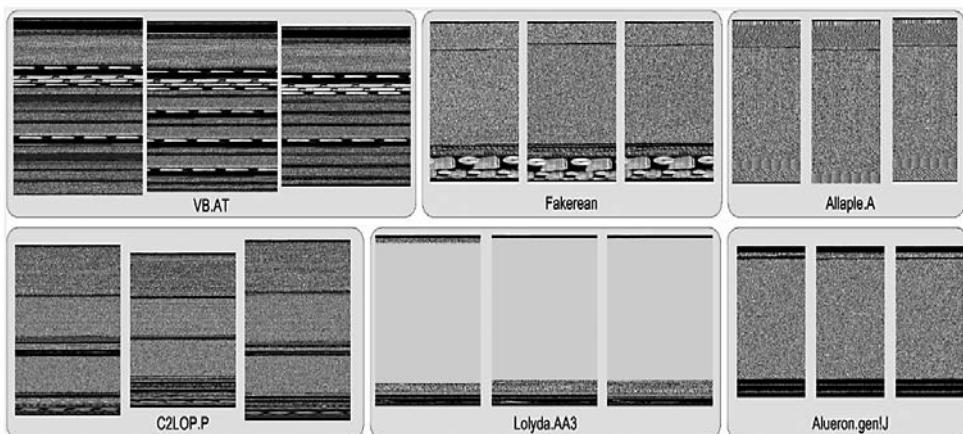


Рис. 1.18. Образцы вредоносного ПО

Как видите, разные типы рассматриваемого ПО для человеческого глаза выглядят совсем по-разному. Модель, обученная исследователями на основе этих представлений изображений, оказалась более точной, чем любые предшествовавшие ей подходы из академической литературы. Из этого следует хорошее эмпирическое правило для преобразования датасета в изображение: если человеческий глаз может распознать категории по изображениям, то и модель глубокого обучения должна это сделать.

Вы поймете, что, как правило, даже с помощью небольшого числа общих подходов в глубоком обучении можно добиться многоного, если вы хоть немного креативны в представлении своих данных. Не стоит воспринимать похожие на описанные здесь подходы как «хитрые обходные пути», потому что они нередко превосходят существующие эталонные результаты и действительно оказываются очень уместными для решения этих проблемных задач.

Обобщение терминов

Мы уже рассмотрели очень много информации, поэтому подытожим. В табл. 1.3 приводится перечень терминов.

Таблица 1.3. Словарь глубокого обучения

Термин	Значение
Метка	Данные, которые мы пытаемся предсказать, такие как «dog» или «cat»
Архитектура	Шаблон модели, которую мы стараемся настроить, то есть реальная математическая функция, в которую мы передаем входные данные и параметры
Модель	Архитектура с конкретным набором параметров
Параметры	Значения модели, изменяющие выполняемое ею задание и обновляющиеся при ее обучении
Fit (Настройка)	Обновление параметров модели, чтобы ее прогнозы на основе входных данных совпадали с целевыми метками
Обучение	Синоним настройки
Предварительно обученная модель	Модель, обученная ранее, как правило, на большом датасете и предполагающая точную настройку
Точная настройка	Обновление предварительно обученной модели для выполнения другой задачи
Эпоха	Один полный проход через входные данные
Потери	Показатель качества модели, выбираемый для обучения через SGD
Метрика	Показатель качества модели при использовании контрольной выборки, выбираемый для оценки человеком
Контрольная выборка	Набор данных, исключаемый из обучения и используемый только для измерения качества модели
Обучающая выборка	Данные, используемые для настройки модели; не включают данные из контрольной выборки
Переобучение	Такое обучение модели, после которого она запоминает конкретные признаки входных данных вместо формирования обобщенного представления для оценки данных, не встречавшихся во время обучения
CNN	Сверточная нейронная сеть — тип нейросети, которая особенно хорошо справляется с задачами компьютерного зрения

Вооружившись этим словарем, мы можем совместить все ключевые понятия, с которыми познакомились. Уделите время просмотру этих определений и прочтению следующего далее краткого изложения пройденного материала. Если вы сможете его понять, значит, к восприятию последующих тем мы вас подготовили достаточно.

Машинное обучение – это дисциплина, в которой мы определяем программу, не прописывая всю ее вручную, а с помощью обучения на основе данных. *Глубокое обучение* – это специализация в машинном обучении, использующая *нейронные сети* со множеством слоев. *Классификация изображений* – ее типичный пример (также известный как *распознавание изображений*). Мы начинаем с *маркированных данных* – набора изображений, в котором мы присвоили каждому изображению *метку*, указывающую, что оно представляет. Наша цель – создать программу, называемую *моделью*, которая при получении нового изображения сможет сделать точный *прогноз* того, что оно представляет.

Каждая модель начинается с выбора *архитектуры* – общего шаблона внутреннего устройства и функционирования этой модели. Процесс *обучения* (или *настройки*) модели – это процесс нахождения набора *значений параметров* (или *весов*), которые преобразуют обобщенную архитектуру в модель, хорошо работающую в отношении конкретного вида данных. Чтобы выяснить, насколько успешно модель справляется с одним прогнозом, нам нужно создать *функцию потерь*, определяющую, как мы оцениваем предсказание: как плохое или хорошее.

Чтобы ускорить процесс обучения, мы можем начать с *предварительно обученной модели* – модели, которая уже обучалась на других данных. После этого можно адаптировать ее к нашим данным, обучая ее на них в течение некоторого времени. Данный процесс называется *тонкой настройкой*.

При обучении модели в первую очередь нужно убедиться, что она успешно *генерализует*: усваивает из данных общие принципы, применимые к новым элементам, с которыми она будет встречаться в дальнейшем, чтобы делать в отношении этих элементов правильный прогноз. При этом есть риск, что если мы обучим нашу модель плохо, то вместо усвоивания общих принципов она эффективно запомнит то, что уже видела, и будет делать плохие прогнозы в отношении новых изображений. Такое отклонение называется *переобучением*.

Чтобы этого избежать, мы всегда разделяем данные на две части: *обучающую выборку* и *контрольную*. Мы обучаем модель, показывая ей только *обучающую выборку*, а затем оцениваем ее качество, наблюдая, как хорошо она прогнозирует элементы из *контрольной выборки*. Чтобы человек мог проанализировать, насколько хорошо модель справляется с *контрольной выборкой*, мы определяем *метрику*. В процессе обучения цикл, в течение которого модель просмотрела каждый элемент *обучающей выборки*, мы называем *эпохой*.

Все эти понятия применимы к машинному обучению в общем, то есть ко всем видам схем для разработки модели путем обучения ее на данных. Отличается же глубокое обучение использованием конкретного класса архитектур, которые основаны на *нейронных сетях*. В частности, задачи вроде классификации изображений в существенной степени опираются на *сверточные нейронные сети*, которые мы вскоре рассмотрим.

Глубокое обучение подходит не только для классификации изображений

Последние годы активно обсуждалась эффективность глубокого обучения для классификации изображений, включая его способность лучше человека выполнять такие сложные задачи, как распознавание злокачественных опухолей на снимках КТ. Однако эта технология способна на гораздо большее, как будет нами показано далее.

Поговорим о важнейшей программной составляющей автономных транспортных средств: локализации объектов на картинке. Если беспилотный автомобиль не знает, где находится пешеход, то, естественно, не сможет избежать столкновения с ним. Создание модели, способной распознавать содержимое каждого отдельного пикселя изображения, называется *сегментацией*. Вот как мы можем обучить модель сегментации с помощью fastai, используя подмножество датасета CamVid (<https://oreil.ly/rDy1i>) из работы *Semantic Object Classes in Video: A High-Definition Ground Truth Database* (<https://oreil.ly/Mqclf>) (Классы семантических объектов в видео: база достоверных данных высокого разрешения), написанной Гэбриэлом Дж. Бростоу (Gabriel J. Brostow) и др.:

```
path = untar_data(URLs.CAMVID_TINY)
dls = SegmentationDataLoaders.from_label_func(
    path, bs=8, fnames = get_image_files(path/"images"),
    label_func = lambda o: path/'labels'/'f'{o.stem}_P{o.suffix}',
    codes = np.loadtxt(path/'codes.txt', dtype=str)
)
learn = unet_learner(dls, resnet34)
learn.fine_tune(8)
```

epoch	train_loss	valid_loss	time
0	2.906601	2.347491	00:02

epoch	train_loss	valid_loss	time
0	1.988776	1.765969	00:02
1	1.703356	1.265247	00:02
2	1.591550	1.309860	00:02
3	1.459745	1.102660	00:02
4	1.324229	0.948472	00:02
5	1.205859	0.894631	00:02
6	1.102528	0.809563	00:02
7	1.020853	0.805135	00:02

Не будем разбирать этот код построчно, потому что он практически идентичен нашему предыдущему примеру. (Углубленно с моделями сегментации мы познакомимся в главе 15, где помимо этого изучим все остальные вкратце упомянутые здесь модели, а также познакомимся со многими другими.)

Можно визуализировать, насколько хорошо модель справилась со своей задачей, попросив ее определить цвет каждого пикселя изображения. Обратите внимание, что все машины выделены одним цветом, так же как и деревья, которые выделены другим (в каждой паре изображений левое представляет истинную метку, а правое — прогноз модели):

```
learn.show_results(max_n=6, figsize=(7,8))
```



Еще одна область, в которой глубокое обучение продвинулось за последние пару лет, — это обработка естественного языка (NLP). Теперь компьютеры могут генерировать текст, выполнять автоматический перевод из одного языка в другой, анализировать комментарии, помечать слова в предложениях и многое другое. Ниже приведен весь код для обучения модели, которая в итоге сможет классифицировать тональность отзывов лучше любых методов, существовавших всего пять лет назад:

```
from fastai.text.all import *

dls = TextDataLoaders.from_folder(unzip_data(URLs.IMDB), valid='test')
learn = text_classifier_learner(dls, AWD_LSTM, drop_mult=0.5, metrics=accuracy)
learn.fine_tune(4, 1e-2)
```

epoch	train_loss	valid_loss	accuracy	time
0	0.594912	0.407416	0.823640	01:35

epoch	train_loss	valid_loss	accuracy	time
0	0.268259	0.316242	0.876000	03:03
1	0.184861	0.246242	0.898080	03:10
2	0.136392	0.220086	0.918200	03:16
3	0.106423	0.191092	0.931360	03:15

Эта модель использует датасет IMDb Large Movie Review (<https://oreil.ly/tl-wp>) из работы *Learning Word Vectors for Sentiment Analysis* (<https://oreil.ly/L9vre>) (изучение векторов слов для анализа тональности текста), написанной Эндрю Маасом (Andrew Maas) и др. Она отлично справляется с отзывами, состоящими из многих тысяч слов, но давайте протестируем ее на короткой рецензии, чтобы увидеть, как она работает:

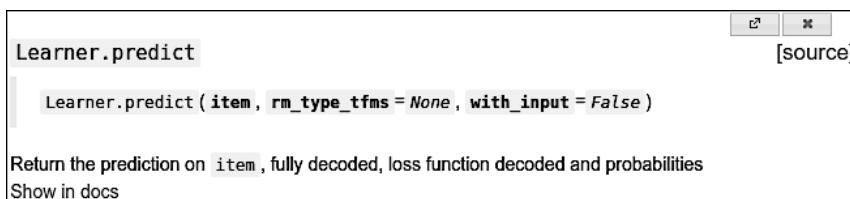
```
learn.predict("I really liked that movie!")
('pos', tensor(1), tensor([0.0041, 0.9959]))
```

Здесь мы видим, что модель оценила отзыв как положительный. Вторая часть результата является индексом *pos* в нашем словаре данных, а последняя — вероятностями, приписываемыми каждому классу (99,6 % для *pos* и 0,4 % для *neg*).

Теперь ваша очередь! Напишите свой собственный отзыв или скопируйте любой из интернета и узнайте, что эта модель о нем думает.

При возникновении вопросов об одном из методов `fastai` вы можете использовать функцию `doc`, передав ей имя интересующего метода:

```
doc(learn.predict)
```



В результате всплывает окно с кратким односторонним объяснением. Ссылка `Show in docs` выдаст документацию (<https://docs.fast.ai/>), где вы найдете все подробности со множеством примеров. Кроме того, большая часть методов `fastai` представлены всего несколькими строками, поэтому вы можете просто щелкнуть на ссылке `source`, чтобы увидеть, что конкретно происходит за кадром.

ПОРЯДОК ВАЖЕН

В блокноте Jupyter важна последовательность выполнения каждой ячейки. Этим он отличается от Excel, где при вводе данных в любом месте все обновляется одновременно. У блокнота же есть внутреннее состояние, которое обновляется при каждом выполнении ячейки. Например, выполняя первую ячейку (с комментарием CLICK ME), вы создаете объект `learn`, который содержит модель и данные для задачи классификации изображений.

Если бы мы выполнили только что приведенную ячейку, прогнозирующую положительность отзыва, сразу после первой, то получили бы ошибку, так как объект `learn` не содержит модель классификации текста. Эту ячейку нужно выполнять после той, что содержит следующее:

```
from fastai.text.all import *

dls = TextDataLoaders.from_folder(untar_data(URLs.IMDB), valid='test')
learn = text_classifier_learner(dls, AWD_LSTM, drop_mult=0.5,
                                 metrics=accuracy)
learn.fine_tune(4, 1e-2)
```

Сами по себе данные вывода могут сбивать с толку, так как включают результаты последнего выполнения ячейки. Если вы измените код внутри ячейки, не выполняя ее, то сохранятся старые (неверные) результаты.

За исключением случаев, когда мы указываем на это явно, блокноты на сайте книги (<https://book.fast.ai/>) выполняются в порядке сверху вниз. Обычно при экспериментировании вы будете выполнять ячейки в любом порядке, чтобы продвигаться как можно быстрее (что очень удобно в блокноте Jupyter). Но как только вы все перепробуете и придете к заключительной версии кода, убедитесь, что можете выполнять ячейки своего блокнота последовательно (ведь в дальнейшем вы можете и не вспомнить тот запутанный путь, которым выполняли их изначально).

В режиме команд двойной набор 0 перезапустит ядро (движок, на котором работает ваш блокнот). Эта процедура очистит состояние блокнота так, как будто вы только что его запустили. Для запуска всех ячеек выше той точки, в которой вы находитесь, выберите в меню Cell команду Run All Above. При разработке библиотеки `fastai` такой вариант выполнения показался нам очень удобным.

Перейдем к гораздо менее увлекательной, но при этом намного более популярной в коммерческом смысле теме: построение моделей на основе простых *табличных* данных.



ТЕРМИН: ТАБЛИЧНЫЙ

Данные, представленные в виде таблицы, например, из электронной таблицы, базы данных или файла значений, разделенных запятыми (CSV). Табличная модель — это модель, которая старается прогнозировать один столбец таблицы на основе информации из других ее столбцов.

Оказывается, и этот пример очень похож на наш изначальный. Ниже приведен код для обучения модели, которая на основе социально-экономического контекста человека будет прогнозировать, имеет ли он высокий доход:

```
from fastai.tabular.all import *
path = untar_data(URLs.ADULT_SAMPLE)

dls = TabularDataLoaders.from_csv(path/'adult.csv', path=path, y_names="salary",
    cat_names = ['workclass', 'education', 'marital-status', 'occupation',
        'relationship', 'race'],
    cont_names = ['age', 'fnlwgt', 'education-num'],
    procs = [Categorify, FillMissing, Normalize])

learn = tabular_learner(dls, metrics=accuracy)
```

Как видите, нам пришлось сообщить fastai, какие столбцы *категориальные* (содержат значения, являющиеся одним из дискретных наборов вариантов, например *occupation* (род занятых)), а какие — непрерывные (содержат число, представляющее количество, например *age* (возраст)).

Для этой задачи нет предварительно обученной модели (как правило, такие модели мало распространены для любых задач моделирования табличных данных, хотя некоторые организации создали их для чисто внутреннего использования), поэтому в данном случае мы не используем *fine_tune*. Вместо этого мы используем *fit_one_cycle*, наиболее распространенный метод для обучения моделей fastai с нуля (то есть без использования переноса обучения):

```
learn.fit_one_cycle(3)
```

epoch	train_loss	valid_loss	accuracy	time
0	0.359960	0.357917	0.831388	00:11
1	0.353458	0.349657	0.837991	00:10
2	0.338368	0.346997	0.843213	00:10

Эта модель использует датасет Adult (<https://oreil.ly/Gc0AR>) из работы *Scaling Up the Accuracy of Naive-Bayes Classifiers: a Decision-Tree Hybrid* (<https://oreil.ly/qFOSc>) (повышение точности наивных байесовских классификаторов: гибрид дерева решений), написанной Роном Кохави (Ron Kohavi), которая содержит некоторые демографические данные об отдельных людях (их образование, семейный статус, расу, пол, а также информацию о том, превышает ли их ежегодный доход 50 тысяч долларов). Точность этой модели выше 80 %, а для ее обучения потребовалось около 30 секунд.

Еще взглянем на рекомендательную систему, которая представляет высокую важность, особенно для электронной коммерции. Такие компании, как Amazon и Netflix, прилагают много усилий для рекомендации товаров или фильмов,

которые могут понравиться пользователям. Ниже показано, как, используя датасет MovieLens (<https://oreil.ly/LCfwH>), обучить модель прогнозировать фильмы, которые понравятся зрителям, на основе истории их просмотров:

```
from fastai.collab import *
path = untar_data(URLs.ML_SAMPLE)
dls = CollabDataLoaders.from_csv(path/'ratings.csv')
learn = collab_learner(dls, y_range=(0.5,5.5))
learn.fine_tune(10)
```

epoch	train_loss	valid_loss	time
0	1.554056	1.428071	00:01

epoch	train_loss	valid_loss	time
0	1.393103	1.361342	00:01
1	1.297930	1.159169	00:00
2	1.052705	0.827934	00:01
3	0.810124	0.668735	00:01
4	0.711552	0.627836	00:01
5	0.657402	0.611715	00:01
6	0.633079	0.605733	00:01
7	0.622399	0.602674	00:01
8	0.629075	0.601671	00:00
9	0.619955	0.601550	00:01

Модель прогнозирует рейтинги фильмов по шкале от 0,5 до 5,0 со средней ошибкой около 0,6. Поскольку мы прогнозируем непрерывное число, а не категорию, то необходимо сообщить fastai, какой диапазон имеет наша цель, используя параметр `y_range`.

Несмотря на то что по факту мы не используем предварительно обученную модель (по той же причине, что и для табличной модели), этот пример показывает, что fastai все равно позволяет нам здесь использовать `fine_tune` (о том, как и почему это работает, вы узнаете в главе 5). Иногда лучше поэкспериментировать с использованием `fine_tune` против `fit_one_cycle`, чтобы выяснить, какой лучше подойдет для вашего датасета.

Мы можем использовать тот же вызов `show_results`, который видели ранее, чтобы просмотреть несколько примеров ID пользователей и фильмов, фактический рейтинг и прогнозы:

```
learn.show_results()
```

	userId	movieId	rating	rating_pred
0	157	1200	4.0	3.558502
1	23	344	2.0	2.700709
2	19	1221	5.0	4.390801
3	430	592	3.5	3.944848
4	547	858	4.0	4.076881
5	292	39	4.5	3.753513
6	529	1265	4.0	3.349463
7	19	231	3.0	2.881087
8	475	4963	4.0	4.023387
9	130	260	4.5	3.979703

ДАТАСЕТЫ: ПИЩА ДЛЯ МОДЕЛЕЙ

В этом разделе вы уже увидели немало моделей, каждая из которых обучалась с помощью разных датасетов и для разных задач. В машинном и глубоком обучении мы не можем обойтись без данных. Поэтому люди, создающие для нас датасеты, могут считаться героями, хоть зачастую и недооцененными. Некоторыми из наиболее полезных и важных датасетов являются те, которые становятся *академическими основами*, — это датасеты, которые активно изучаются исследователями и используются для сравнения алгоритмических изменений. Некоторые из них становятся общеизвестными (по крайней мере среди тех, кто занимается обучением моделей), например MNIST, CIFAR-10 и ImageNet.

Датасеты, используемые в этой книге, были выбраны потому, что предоставляют отличные примеры тех видов данных, с которыми вы, скорее всего, столкнетесь. В академической литературе есть много примеров результатов моделей, использующих эти датасеты, с которыми вы можете сравнить итоги собственной работы.

Большинство использованных нами в книге датасетов потребовали от своих создателей больших усилий. Например, несколько позднее мы покажем вам, как создавать модель, способную делать перевод между французским и английским. В качестве входных данных будет использован свод параллельных текстов на французском/английском, подготовленный профессором Крисом Каллисон-Берч (Chris Callison-Burch) из Университета Пенсильвании в 2009 году. Этот датасет содержит более 20 миллионов пар предложений. Профессор создал его поистине умным способом: просканировал миллионы канадских веб-страниц (которые зачастую мультиязычны), а затем использовал набор простых эвристик для преобразования URL-адресов французского контента в URL-адреса, указывающие на тот же контент на английском.

По мере знакомства с датасетами в книге подумайте, откуда они могли взяться и как могли быть организованы. Затем подумайте о том, какие виды интересных датасетов вы могли бы создать для собственных проектов. (Вскоре мы даже пошагово проведем вас через процесс создания своего датасета изображений.)

fast.ai вложила немало времени в создание урезанных версий популярных датасетов, которые специально спроектированы для поддержки быстрого прототипирования и экспериментов, а также для упрощенного обучения с их помощью. В этой книге мы будем часто начинать с использования одной из урезанных версий, а после переходить к их полноразмерным вариантам (также, как мы делаем в этой главе). Именно так ведущие практики во всем мире выполняют свое моделирование. Большую часть экспериментов и прототипирования они производят с подмножествами имеющихся данных, а полные датасеты используют только тогда, когда у них уже сформировано четкое понимание необходимых действий.

Каждая из обученных нами моделей показала обучающие и контрольные потери. Хорошая контрольная выборка — это один из важнейших элементов процесса обучения. Давайте узнаем почему и научимся создавать такую.

Контрольные и тестовые выборки

Как мы уже говорили, цель модели — прогнозирование данных. Однако процесс обучения модели в корне примитивен. Если бы мы обучали модель на всех имеющихся данных, а затем оценивали ее с помощью тех же данных, то не сумели бы понять, насколько хорошо она справится с данными, которые не видела. Без использования этого ценного элемента информации в качестве ориентира при обучении модели велик шанс, что она будет хорошо делать прогнозы в отношении использованных данных, но при этом ошибаться, встречая новые.

Чтобы этого избежать, нашим первым шагом было разделение датасета на два набора: *обучающей выборки* (предоставляемой модели во время обучения) и *контрольной выборки*, также называемой *development set* (которая используется только для оценки). Это позволяет нам проверить, извлекает ли модель общие принципы из обучающих данных, используя их в отношении новых, то есть контрольных, данных.

Для лучшего представления можно сказать, что в определенном смысле мы не хотим, чтобы модель получала хорошие результаты, «жульничая». Если она делает точный прогноз для элемента данных, это должно происходить из-за того, что она усвоила характеристики этого вида элемента, а не потому, что была сформирована путем *фактического просмотра конкретно этого элемента*.

Отделение контрольных данных означает, что модель никогда не увидит их во время обучения и, следовательно, ничего о них не узнает. Жульничать у нее уже не получится, ведь так?

По правде говоря, не совсем. Ситуация несколько тоньше. Все потому, что в реальных сценариях мы редко строим модель, обучая ее параметры всего раз. Чаще мы пробуем много версий модели, используя различные варианты

моделирования в отношении архитектуры сети, скорости обучения, стратегий аугментации данных и других факторов, которые мы будем обсуждать в последующих главах. Многие из этих вариантов можно описать как варианты *гиперпараметров*. Это слово означает, что они представляют параметры параметров, поскольку являются более высокоуровневыми вариантами, управляющими значением параметров весов.

Проблема в том, что в обычном процессе обучения при изучении значений для параметров весов рассматриваются только прогнозы обучающих данных. Мы же, как создатели модели, решая попробовать новые значения гиперпараметров, оцениваем модель путем изучения прогнозов для контрольной выборки. Следовательно, последующие версии модели косвенно сформированы нами после просмотра контрольных данных. Подобно тому как автоматический процесс обучения подвержен риску переобучения, так и мы, пробуя различные варианты и совершая ошибки, подвержены риску переобучения контрольными данными.

Решением этой головоломки явилось введение дополнительного уровня еще более скрытых данных: *тестовой выборки*. Аналогично тому, как мы удерживаем контрольные данные от процесса обучения, нужно удерживать тестовые даже от самих себя. Их нельзя использовать для улучшения модели. Они применяются только для ее окончательной оценки в самом конце всех наших стараний. В результате мы определяем иерархию нарезок данных исходя из того, насколько полноценно мы хотим спрятать их от процессов обучения и моделирования: обучающие данные полностью раскрыты, контрольные данные раскрыты частично, а тестовые полностью скрыты. Эта иерархия параллельна разным способам моделирования и оценки — автоматическому процессу обучения с обратным распространением ошибок, более ручному способу проверки разных гиперпараметров в разных сессиях обучения и оценке окончательного результата.

Тестовые и контрольные выборки должны содержать достаточно данных, чтобы обеспечить вам хорошую оценку точности. Если вы создаете, к примеру, обнаружитель кошек, то обычно в контрольной выборке должно быть не менее 30 их экземпляров. Это означает, что если у вас датасет с тысячами элементов, то использование 20 % в качестве контрольной выборки может быть излишним. С другой стороны, если у вас избыток данных, то использование их части в качестве контрольных не должно навредить.

Наличие двух уровней «зарезервированных данных» — контрольной выборки и тестовой, где один из этих уровней вы как бы прячете от самих себя, может показаться несколько чрезмерным. Однако зачастую это необходимо, так как модели склонны к поиску простейшего пути успешного прогнозирования (запоминанию), а мы же, наивные люди, склонны к самообману относительно того, как наши модели хорошо работают. В этом случае проверка тестовой выборкой помогает нам оставаться интеллектуально честными. Это не значит, что нам *всегда* нужно отделять тестовую выборку — если у вас очень мало данных, то,

возможно, будет достаточно и контрольного набора, — однако по возможности лучше использовать также и тестовый.

Это же правило может оказаться критически важным в случае найма стороннего исполнителя для выполнения моделирования от вашего лица. Сторонний человек может не до конца понять ваши требования либо ввиду своих корыстных интересов в принципе им не последовать. Хорошая тестовая выборка может существенно снизить эти риски и позволить вам успешно оценить, решает ли проделанная наемной стороной работа вашу фактическую задачу.

Говоря прямо, если в компании вы занимаете роль лица, принимающего решения (либо являетесь его советником), то главное понять одно: если вы достаточно разберетесь в том, что такая контрольная и тестовая выборки, а также почему они важны, то избежите одного крупнейшего источника провалов, который мы встречали в организациях, стремящихся внедрить ИИ. Например, если вы рассматриваете привлечение внешнего исполнителя или службы, убедитесь, что удерживаете тестовые данные, которые исполнитель *никогда не увидит*. Затем вы сможете проверить созданную ими модель на этих тестовых данных, используя метрики, выбранные *вами* на основе реально значимых практических задач, и вы же решите, какой уровень производительности считать достаточным. (Также нeliшним будет попробовать создать базовую модель самостоятельно, чтобы понять, каких результатов она сможет достичь. Зачастую выясняется, что ваша простая модель работает не хуже той, что разработал сторонний «эксперт»!)

Создавайте тестовую выборку обдуманно

Для успешного определения контрольной (и, возможно, тестовой выборки) иногда потребуется не просто изъять случайную долю исходного датасета. Помните, что ключевая черта контрольной и тестовой выборки в том, что они должны представлять новые данные, которые вы увидите только в будущем. Это может показаться невыполнимым, ведь по определению вы эти данные еще не видели. Однако обычно вам все-таки кое-что известно.

Полезно взглянуть на несколько примеров случаев. Многие из этих примеров взяты с соревнований по предиктивному моделированию, проводимых на платформе Kaggle (<https://www.kaggle.com/>), и хорошо отражают задачи и методы, с которыми вы можете столкнуться на практике.

Одним из таких примеров может стать рассмотрение данных временных рядов. Здесь выбрать случайное подмножество данных будет очень просто (можно просмотреть данные как до, так и после дат, которые вы хотите прогнозировать), но вместе с тем такие данные не будут представлять большинство бизнес-сценариев использования (где вы используете исторические данные, чтобы построить модель для работы в будущем). Если ваши данные включают даты и вы создаете модель для использования в будущем, то в качестве контрольной выборки

понадобится выбрать непрерывный отрезок с последними датами (например, последние две недели последнего месяца).

Предположим, нужно разделить данные временного ряда (рис. 1.19) на обучающую и контрольную выборки.

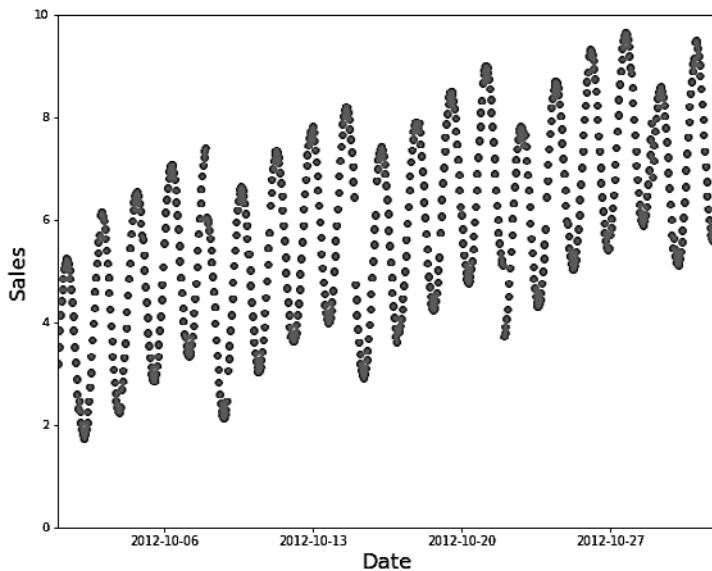


Рис. 1.19. Временной ряд

Как видно на рис. 1.20, случайное подмножество будет неудачным выбором (слишком легко заполнить в нем пробелы, и оно не является показательным).

Вместо этого используйте в качестве обучающей выборки более ранние данные (а поздние данные — для контрольной выборки) (рис. 1.21).

К примеру, на Kaggle проводилось соревнование по прогнозированию продаж в сети эквадорских продуктовых магазинов (<https://oreil.ly/UQoXe>). Обучающая выборка Kaggle включала даты от 1 января 2013 года до 15 августа 2017 года. Тестовые данные охватывали период с 16 по 31 августа 2017 года. Таким образом, организатор соревнований гарантировал, что с точки зрения модели участники делали прогнозы для временного отрезка в будущем. Аналогичным образом трейдеры количественных хедж-фондов проводят *ретроспективное тестирование*, чтобы выяснить, позволяют ли их модели прогнозировать будущие периоды на основе прошлых данных.

Второй случай — это когда вы легко предвидите, как данные, для которых будете делать прогнозы в работе, могут *качественно отличаться* от обучающих данных.

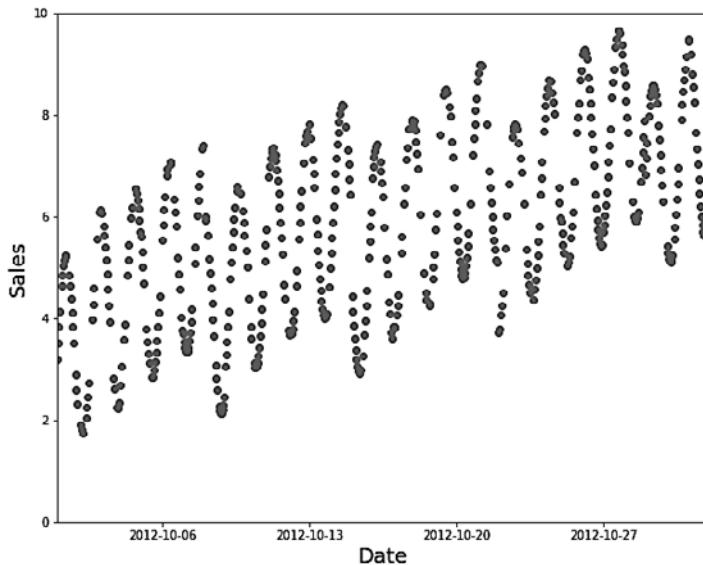


Рис. 1.20. Плохое обучающее подмножество

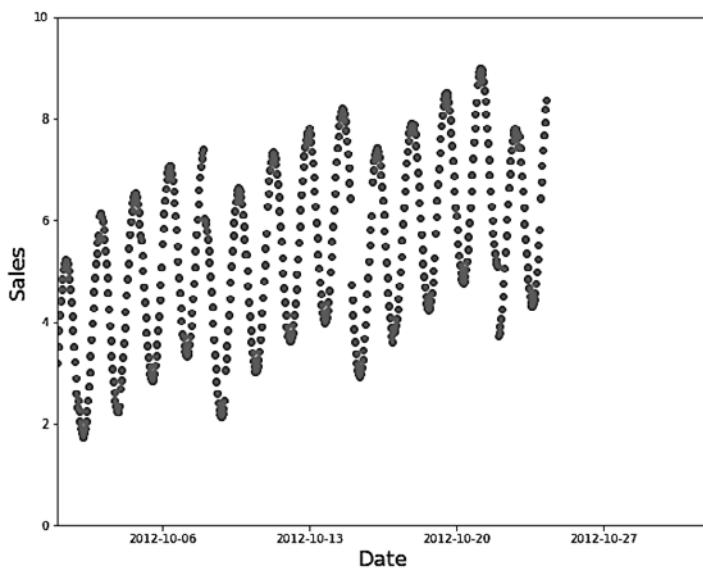


Рис. 1.21. Хорошее обучающее подмножество

В соревновании по распознаванию отвлекающихся водителей (https://oreil.ly/zT_tC), также проводимом Kaggle, независимыми переменными являлись картинки водителей за рулем автомобиля, а зависимыми — такие категории, как

отправка сообщений, принятие пищи или внимательный взгляд вперед. Многие изображения представляют тех же водителей в разных положениях (рис. 1.22). Если вы представите себя страховой компанией, разрабатывающей модель на основе этих данных, то тогда в первую очередь вас будет интересовать то, как модель справится с прогнозом в отношении водителей, которых она ранее не видела (поскольку у вас наверняка будут данные только для небольшого числа людей). С учетом этого условия тестовые данные для соревнования состояли из изображений людей, которые не встречаются в обучающей выборке.



Рис. 1.22. Два изображения из обучающих данных

Если вы поместите одно из приведенных на рис. 1.22 изображений в обучающую выборку, а второе — в контрольную, то ваша модель легко сделает прогноз для последней, то есть в итоге она покажет лучший результат, чем при работе с изображениями новых людей. Еще возможен вариант, когда вы используете все изображения людей только для обучающей выборки, и в итоге модель может переобучиться распознавать особенности именно этих людей, не изучив их общие состояния (набирает сообщение, ест и т. д.).

Похожая динамика наблюдалась и в рыбоохранном соревновании Kaggle (<https://oreil.ly/iJwFf>), где задачей было определить вид рыбы, пойманной рыбаками лодками, чтобы снизить нелегальный отлов исчезающих популяций. Тестовая выборка состояла из изображений с лодок, не встречающихся в обучающих данных, следовательно, и для контрольной выборки также понадобились бы данные с лодок, не участвующих в обучающей выборке.

Иногда может быть не до конца понятно, как должны отличаться контрольные данные. Например, в задаче, использующей данные спутниковой съемки, вам понадобится собрать больше информации о том, содержит ли обучающая выборка только конкретные географические локации или же была получена из географически разрозненных данных.

Теперь, когда вы уже почувствовали процесс создания модели, настало время решать, в каком направлении двигаться дальше.

Момент выбора собственного приключения

Если вы хотите побольше узнать о том, как использовать модели глубокого обучения на практике, включая обнаружение и исправление ошибок, создание реальных веб-приложений и избежание неожиданного причинения моделью вреда вашей организации или обществу, тогда продолжайте читать следующие две главы. Если же вы хотите приступить к изучению основ внутренних процессов глубокого обучения, переходите к главе 4. (Читали ли вы в детстве книги в стиле «*Выбери свой путь?*»? Мы предлагаем вам нечто похожее... Только здесь у нас гораздо больше глубокого обучения, чем в них.)

Тем не менее для успешного продвижения по ходу книги вам потребуется прочитать все эти главы, хотя порядок их прочтения можете выбирать сами, так как друг от друга они не зависят. Если же вы сразу перейдете к главе 4, то в ее конце мы напомним вам, что лучше вернуться и прочесть пропущенные главы, прежде чемходить далее.

Вопросник

Прочитав много страниц текста, может быть сложно выделить ключевые моменты, на которых важнее всего заострить внимание. Поэтому в конце каждой главы мы подготовили список вопросов и список дальнейших шагов. Все ответы — в тексте пройденной главы, поэтому если вы в чем-то не уверены, перечитайте соответствующую часть текста, убедившись, что поняли ее. Ответы на все приводимые вопросы также доступны на сайте книги (<https://book.fast.ai/>). В случае затруднений вы можете посетить форумы (<https://forums.fast.ai/>), где вам помогут разобраться люди, изучающие этот материал.

1. Что вам нужно для глубокого обучения?
 - Много математики Д/Н
 - Много данных Д/Н
 - Много дорогих компьютеров Д/Н
 - Ученая степень Д/Н
2. Назовите пять областей, в которых глубокое обучение является самым эффективным инструментом.
3. Как называлось первое устройство, которое основывалось на принципе искусственного нейрона?
4. Исходя из одноименной книги, перечислите требования к параллельной распределенной обработке (PDP).
5. Какие два теоретических заблуждения привели к задержке развития технологий нейронных сетей?

6. Что такое GPU?
7. Откройте блокнот и выполните ячейку, содержащую $1 + 1$. Что произойдет?
8. Пройдите по каждой ячейке пустой версии блокнота для этой главы. Перед выполнением каждой этой ячейки постараитесь угадать, что в итоге произойдет.
9. Выполните инструкции в онлайн-приложении Jupyter Notebook (<https://oreil.ly/9uPZe>).
10. Почему сложно использовать для распознавания образов на фото традиционную компьютерную программу?
11. Что Сэмюэл подразумевал под «назначением весов»?
12. Какой термин мы обычно используем в глубоком обучении для того, что Сэмюэл назвал «весами»?
13. Нарисуйте картинку, в целом отражающую модель машинного обучения, как ее видел Сэмюэл.
14. Почему сложно понять причину, по которой модель глубокого обучения делает конкретный прогноз?
15. Как называется теорема, показывающая, что нейронная сеть может решить любую математическую задачу с любым уровнем точности?
16. Что вам необходимо для обучения модели?
17. Как может петля обратной связи повлиять на внедрение модели прогнозирования полицейской деятельности?
18. Всегда ли мы должны использовать для модели распознавания кошек изображения размером 224×224 пикселя?
19. В чем разница между классификацией и регрессией?
20. Что такое контрольная выборка? Что такая тестовая выборка? Зачем они нужны?
21. Что сделает fastai, если вы не предоставите контрольную выборку?
22. Можем ли мы всегда использовать для контрольной выборки случайные образцы? Почему?
23. Что такое переобучение? Приведите пример.
24. Что такое метрика? Чем она отличается от потерь?
25. Как могут помочь предварительно обученные модели?
26. Что такое голова модели?
27. Какие виды признаков обнаруживают ранние слои CNN, а какие — поздние?

28. Модели изображений полезны только для работы с фотографиями?
29. Что такое архитектура?
30. Что такое сегментация?
31. Для чего используется метод `y_range`? Когда он нам нужен?
32. Что такое гиперпараметры?
33. Каков наилучший способ избежать провалов использования ИИ в организациях?

Дополнительные задания

В каждой главе также есть раздел «Дополнительные задания», который ставит некоторые не до конца раскрытые в тексте вопросы или дает более продвинутые задания. Ответы на эти вопросы отсутствуют на сайте книги, так что вам понадобится провести собственное исследование.

1. Почему GPU полезен для глубокого обучения? Чем отличается CPU и почему в этой области он менее эффективен?
2. Постарайтесь подумать о трех областях, где петли обратной связи могут повлиять на использование машинного обучения. Попробуйте также найти задокументированные примеры того, что в этих случаях происходит на практике.

ГЛАВА 2

От модели к продакшену

Шесть строк кода, которые мы видели в главе 1, описывают лишь малую часть процесса практического применения глубокого обучения. В этой главе мы будем использовать пример компьютерного зрения, чтобы рассмотреть сквозной процесс создания приложения глубокого обучения. Говоря конкретнее, мы собираемся построить классификатор медведей! В течение этого процесса мы будем обсуждать возможности и ограничения техники глубокого обучения, рассмотрим возможные подводные камни ее практического применения, узнаем, как создавать датасеты и не только. Многие ключевые пункты будут одинаково применимы и к другим задачам глубокого обучения наподобие рассмотренных нами в главе 1. Если вы будете решать похожую на один из наших примеров задачу, то с большой вероятностью быстро добьетесь блестящих результатов, используя минимум кода.

Сначала научимся правильно формулировать задачу.

Практика глубокого обучения

Мы видели, что глубокое обучение может решать множество сложных задач быстро и с помощью небольшого количества кода. Для новичков есть благоприятная область задач, достаточно схожих с приводимыми нами примерами, в которых вы можете очень быстро добиться чрезвычайно полезных результатов. Однако глубокое обучение — это не магия. Те же шесть строк кода не будут работать для каждой задачи, о которой можно подумать сегодня.

Недооценка ограничений и переоценка возможностей глубокого обучения может приводить к весьма разочаровывающим результатам, по меньшей мере до тех пор, пока вы не наберетесь опыта и сможете решать возникающие задачи. И наоборот, переоценка ограничений и недооценка возможностей этой техники может означать, что вы не беретесь за решаемые задачи, потому что сами себя отговариваете.

Мы часто общаемся с людьми, которые недооценивают и ограничения, и возможности глубокого обучения. И то и другое может быть проблемой: недооценка возможностей означает, что вы можете так и не попробовать то, что могло принести много пользы, а недооценка ограничений может означать, что вы не учтете существенные сложности и не сможете отреагировать на них.

В этом смысле лучше оставаться открытым возможностям. Если вы не станете отрицать перспективу, что глубокое обучение сможет решить часть вашей задачи с меньшим количеством данных или более простым путем, чем вы ожидали, то вам удастся разработать процесс, который поможет вам найти конкретные возможности и ограничения вашей задачи. Это не означает, что нужно предпринимать рискованные шаги, — мы покажем вам, как поэтапно развертывать модели, чтобы они не несли существенных рисков, и даже научим производить ретроспективное тестирование перед их запуском.

Начало проекта

Итак, с чего же начать свое путешествие в мир глубокого обучения? Самое важное — это убедиться, что у вас есть проект для работы: только работая над собственными проектами, вы сможете получить реальный опыт создания и использования моделей. При выборе проекта важнее всего учесть доступность данных.

Независимо от того, делаете вы проект только для собственного обучения или для практического применения в организации, вам захочется иметь возможность быстро приступить к работе. Мы видели многих студентов, исследователей и практиков этой отрасли, которые тратили месяцы и даже годы в поиске своего идеального датасета. Цель не в том, чтобы найти «идеальный» датасет или проект, а в том, чтобы просто начать и продолжать пробовать. Если вы используете такой подход, то окажетесь уже на третьем этапе обучения и совершенствования, в то время как perfectionists продолжат топтаться на стадиях планирования!

Мы также предлагаем реализовывать весь проект от и до, не тратя месяцы на точную настройку модели, оттачивание идеального GUI или разметку совершенного датасета.

Завершайте каждый шаг настолько хорошо, насколько можете, в течение разумного промежутка времени, продвигаясь до конца. Например, если вашей конечной целью является приложение, выполняемое на мобильном телефоне, то именно его вы и должны получать после каждого подхода. Возможно, на ранних этапах вы будете использовать короткие пути, например, выполняя всю обработку на удаленном сервере и используя простое отзывчивое веб-приложение. Выполнив проект от начала и до конца, вы увидите, какие его части оказались самыми витиеватыми, а какие максимально повлияли на итоговый результат.

По мере изучения книги наряду с постепенной разработкой вами собственных проектов мы будем предлагать много небольших экспериментов, которые вы сможете проводить, выполняя и настраивая предоставляемые нами блокноты. Таким образом вы будете получать опыт работы со всеми инструментами и техниками по ходу их объяснения.



СЛОВО СИЛЬВЕЙНУ

Для получения максимальной пользы от книги уделите время экспериментам, предлагаемым после каждой главы, реализуя их либо в собственном проекте, либо в предоставленных нами блокнотах. Затем попробуйте переписать эти блокноты с нуля для нового датасета. Выработать чутье к правильному обучению модели вы сможете только на практике и неизбежных ошибках.

Используя описанный сквозной подход с повторениями, вы лучше поймете, сколько данных вам реально необходимо. Например, вы можете обнаружить, что вам легко доступны только 200 элементов размеченных данных, и пока вы не попробуете их использовать, то никак не узнаете, достаточно ли этого количества для достижения нужной вашему приложению производительности модели.

Вы сможете показать коллегам, что ваша идея работает, продемонстрировав им рабочий прототип. Мы неоднократно убеждались, что такой подход — это секрет получения поддержки проекта со стороны организации.

Поскольку легче всего начать работать над проектом, для которого у вас уже есть доступные данные, это означает, что наверняка легче всего начать проект, связанный с тем, что вы уже делаете. Например, если вы работаете в музыкальном бизнесе, то можете иметь доступ ко многим аудиозаписям. Если вы работаете радиологом, то наверняка имеете доступ ко множеству медицинских снимков. Если вы интересуетесь сохранением дикой природы, то можете иметь доступ ко многим изображениям ее представителей.

Иногда нужно быть чуточку креативными. Возможно, вы сможете найти предыдущий проект машинного обучения, например соревнование Kaggle, связанное с вашей областью интересов. Иногда нужно идти на компромисс. К примеру, не удается найти точные данные, необходимые для конкретного проекта, который вы задумали. Но тогда можно попробовать найти что-то из схожей области или в ином формате и решить немного другую задачу. Работа над подобными похожими проектами по-прежнему будет давать вам хорошее понимание всего процесса и может помочь найти другие короткие пути, источники данных и т. д.

Имейте в виду, что особенно в начале обучения будет плохой идеей рассеивать свое внимание по нескольким совершенно разным областям, начиная работать там, где глубокое обучение еще не применялось. Причина в том, что если ваша модель не будет с первого раза работать отлично, то вы не сможете узнать, явилась ли причиной ваша ошибка или же просто сама задача не поддается

решению с помощью глубокого обучения. Вы также не будете знать, где искать в этом вопросе помощи. Поэтому лучше всего начать с поиска онлайн-примера, в котором уже были достигнуты хорошие результаты и который хоть в какой-то степени похож на то, что вы хотите реализовать. В этом случае вы сможете преобразовать свои данные в формат, аналогичный использованному ранее (например, создать изображения для ваших данных). Теперь взглянем на уровень глубокого обучения, чтобы вы понимали, в чем сейчас эта область хороша.

Текущий уровень глубокого обучения

Начнем с рассмотрения того, насколько глубокое обучение может подойти к задаче, с которой вы хотите работать. В этом разделе приводится краткий обзор уровня глубокого обучения на начало 2020 года. Однако все быстро меняется, и к моменту, когда вы будете это читать, некоторые из приведенных ограничений могут уже не существовать. Мы, в свою очередь, постараемся поддерживать актуальную информацию на сайте книги, а вы, помимо этого, можете также воспользоваться поиском в Google по запросу *What can AI do now* (Современные возможности ИИ), чтобы получить общее представление о текущей ситуации.

Компьютерное зрение

Существует много областей, в которых глубокое обучение еще не использовалось для анализа изображений, но те, в которых оно опробовано, практически во всех случаях показали, что компьютеры способны распознавать элементы на изображениях по крайней мере не хуже людей, даже таких специально обученных профессионалов, как радиологи. Эта их способность называется *распознаванием объектов*. Глубокое обучение также хорошо справляется с определением положения объектов на изображении и может выделять занимаемые ими области и называть каждый обнаруженный объект. Эта функция уже называется *обнаружением объектов* (в одном из ее вариантов, который мы видели в главе 1, каждый пиксель классифицировался на основе типа объекта, частью которого являлся, — это называется *сегментацией*).

Алгоритмы глубокого обучения обычно плохо справляются с распознаванием изображений, существенно отличающихся в структуре или стиле от использованных для обучения модели. Например, если в обучающих данных не было черно-белых изображений, то модель будет плохо справляться с их распознаванием. Аналогичным образом, если в обучающих данных отсутствовали изображения, нарисованные от руки, то модель наверняка будет плохо оценивать такие на практике. Нет никакого общего способа проверить, какого типа изображений недостает в вашей обучающей выборке, но в этой главе мы покажем некоторые способы распознать появление неожиданных типов изображений в данных при использовании модели в работе (эта техника известна как проверка на наличие данных за пределами области применимости модели —out-of-domain data).

Наша главная сложность в системах обнаружения объектов в том, что маркировка изображений может быть медленной и дорогостоящей. Сейчас активно разрабатываются методики и инструменты для ускорения и облегчения процесса маркировки, а также для уменьшения числа вручную проставленных меток, необходимого для обучения точных моделей обнаружения объектов. Один из подходов, который особенно полезен, — это синтетическая генерация вариаций входных изображений, например, их вращением или изменением яркости/контраста. Эта техника называется *аугментация данных* и также отлично работает для текста и других типов моделей. Чуть позже мы обсудим ее подробнее.

Помимо этого, стоит учитывать, что хоть ваша задача может и не выглядеть относящейся к компьютерному зрению, но, возможно, с помощью некоторой доли воображения ее можно в таковую преобразовать. Например, если вы хотите классифицировать звуки, то можете попробовать преобразовать эти звуки в изображения форм их звуковых волн, на которых затем и обучить свою модель.

Текст (обработка естественного языка)

Компьютеры хорошо справляются с классификацией как коротких, так и длинных документов на основе таких категорий, как спам/не спам, тональность (например, положительный отзыв или отрицательный), авторство, сайт-источник и т. д. Мы не знаем о каких-либо серьезных проделанных в этом направлении работах, чтобы сравнить компьютеры с людьми, но по некоторым наблюдениям можно сделать вывод, что производительность глубокого обучения в подобных задачах аналогична человеческой.

Глубокое обучение также хорошо зарекомендовало себя в генерации соответствующего контексту текста, например ответов на посты в социальных сетях и имитации стиля конкретного автора. Причем такие модели вполне успешно делают этот контент убедительным для людей: порой даже более убедительным, чем текст, написанный человеком. Но глубокое обучение слабо справляется с генерацией *верных* ответов. У нас нет надежного способа, чтобы, например, совместить базу медицинских знаний с моделью глубокого обучения для генерации верных с медицинской точки зрения ответов на естественном языке. Это опасно, потому что очень легко создать контент, который для неспециалиста покажется убедительным, а на деле окажется абсолютно неверен.

Озабоченность также вызывает то, что соответствующие контексту убедительные ответы в социальных сетях могут использоваться в огромных масштабах, в тысячи раз превышающих любую известную фабрику троллей по распространению дезинформации, созданию беспокойства и провоцированию конфликтов. В качестве одного из основных правил можно отметить, что модели для генерации текста всегда будут технологически на шаг опережать модели для распознавания автоматически сгенерированного текста. Например, можно использовать модель, которая способна распознавать искусственно сгенериро-

ванный контент, чтобы улучшать создающий его генератор до тех пор, пока эта модель классификации не перестанет справляться с задачей.

Несмотря на описанные сложности, глубокое обучение широко применяется в NLP: его можно использовать для перевода текста с одного языка на другой, обобщения объемных документов в форму для более быстрого восприятия, нахождения всех упоминаний интересующих понятий и др. К сожалению, в таких случаях перевод или обобщение могут вполне включать неверную информацию. Их текущего качества уже достаточно для того, чтобы люди эти системы использовали: например, онлайн-переводчик Google (и все другие известные нам аналогичные онлайн-сервисы) основан на глубоком обучении.

Совмещение текста и изображений

Способность глубокого обучения совмещать текст и изображения в единую модель оказывается намного лучше, чем многие люди могут изначально предположить. Например, модель можно обучить на входных изображениях с выходными подписями на английском языке, и она научится генерировать удивительно подходящие подписи автоматически для новых изображений. Но важно понимать, что такие подписи не обязательно будут всегда верны.

Поскольку это серьезная проблема, мы обычно советуем использовать глубокое обучение не в качестве полностью автоматической процедуры, а как часть процесса, в котором происходит близкое взаимодействие модели и человека. Такой подход может потенциально сделать людей на порядок более продуктивными, чем при использовании ими только ручных методов, что приведет к более точному выполнению задач.

Например, автоматическая система может использоваться для определения потенциальных жертв инсульта прямо по КТ-снимкам и отправки срочного сигнала для скорейшего просмотра этих снимков специалистом. В случае инсульта на принятие эффективных мер у врача есть всего три часа, следовательно, такой быстрый цикл обратной связи может спасти жизни. При этом в то же время все снимки могут по-прежнему отправляться радиологам обычным способом, а это значит, что снижения доли человеческого участия не произойдет. Другие модели глубокого обучения могут автоматически измерять элементы со снимков и вставлять полученные показатели в отчеты, предупреждая радиологов о результатах, которые они могли пропустить, а также сообщая о других случаях, которые могут иметь отношение к делу.

Табличные данные

Глубокое обучение в последнее время добилось больших успехов в анализе временных рядов и табличных данных. Но эта технология обычно используется как часть ансамбля из нескольких типов моделей. Если у вас уже есть система,

использующая случайные леса или градиентный бустинг (популярные инструменты табличного моделирования, о которых вы скоро узнаете подробнее), тогда переход на глубокое обучение или его добавление может не привести к существенному улучшению.

Глубокое обучение отлично расширяет вариативность столбцов, которые вы можете включать: например, столбцы, содержащие естественный язык (заголовки книг, отзывы и т. д.), и категориальные столбцы с высокой кардинальностью (то есть такие, которые содержат большое число дискретных вариантов, таких как zip-код или ID товара). С другой стороны, модели глубокого обучения, как правило, требуют больше времени на обучение, чем случайные леса или градиентный бустинг. Хотя ситуация в этом плане изменяется благодаря, например, библиотеке RAPIDS (<https://rapids.ai/>), которая обеспечивает ускорение GPU для всего конвейера моделирования. Плюсы и минусы всех названных методов мы подробно рассмотрим в главе 9.

Рекомендательные системы

Рекомендательные системы, по сути, являются просто особым видом табличных данных. В частности, они обычно имеют одну категориальную переменную высокой кардинальности, которая представляет пользователей, и еще одну, представляющую товары (или их аналог). Такие компании, как Amazon, представляют все когда-либо совершенные клиентами покупки как гигантскую разреженную матрицу, где клиенты отражены в строках, а покупки — в столбцах. Когда у компании появляются данные в подобном формате, специалисты применяют определенную форму совместной фильтрации, чтобы эту *матрицу заполнить*. Например, если клиент А покупает товары 1 и 10, а клиент Б покупает товары 1, 2, 4 и 10, то система порекомендует А купить 2 и 4.

Так как модели глубокого обучения хороши в обработке категориальных переменных высокой кардинальности, они также хороши в обработке рекомендательных систем. Поэтому, в частности, их потенциал раскрывается, как и в случае с табличными данными, при совмещении этих переменных с другими видами данных, такими как естественный язык или изображения. Они также могут успешно справляться с совмещением всех этих типов информации с дополнительными метаданными, представленными в виде таблиц, наподобие пользовательской информации, предыдущих транзакций и т. д.

Тем не менее практически все подходы ML имеют недостаток: они сообщают вам только о том, какие товары могут понравиться конкретному пользователю, но не дадут реально полезных рекомендаций. Многие варианты рекомендаций товаров, потенциально интересующих пользователя, могут на деле оказаться бесполезными: например, если пользователь уже знаком с этими товарами или если это уже приобретенные им товары, но иначе упакованные (как вариант, коробочный набор романов, которые клиент уже приобрел по отдельности).

Джереми любит читать книги Терри Пратчетта, в связи с чем Amazon какое-то время рекомендовал ему произведения только этого автора (рис. 2.1), что не было полезным, потому что об этих книгах Джереми уже знал.



Рис. 2.1. Не очень полезные рекомендации

Другие типы данных

Со временем вы обнаружите, что специфичные для конкретной области типы данных очень хорошо подходят под существующие категории. Например, белковые цепочки схожи с документами на естественном языке в том, что представляют длинную последовательность дискретных токенов со сложными связями и значением на протяжении этой последовательности. На деле так и выясняется, что использование NLP-методов глубокого обучения — эталонный подход для многих типов белкового анализа. Еще один пример: звуки можно представить в виде спектрограмм и использовать как изображения. При этом стандартные подходы глубокого обучения для обработки изображений отлично работают со спектрограммами.

Подход Drivetrain

Многие точные модели оказываются никому не нужными, в то время как многие неточные оказываются очень полезны. Чтобы обеспечить полезность своей работы, вам нужно проанализировать, как она будет в итоге использоваться. В 2012 году Джереми вместе с Маргит Звемер (Margit Zwemer) и Майком Лукидесом (Mike Loukides) для подобного анализа ввели метод под названием *Drivetrain Approach* (трансмиссионный подход).

Этот подход, изображенный на рис. 2.2, был детально описан в работе *Designing Great Data Products* (<https://oreil.ly/KJPa>) («Проектирование грамотных продуктов для обработки данных»). Основная идея в том, чтобы начать с рассмотрения поставленной задачи, а затем подумать, какие действия можно предпринять

для ее решения и какие вспомогательные данные для этого есть (или могут быть получены). Затем нужно построить модель, которую можно использовать, чтобы определить оптимальные действия для достижения наилучших результатов этой задачи.

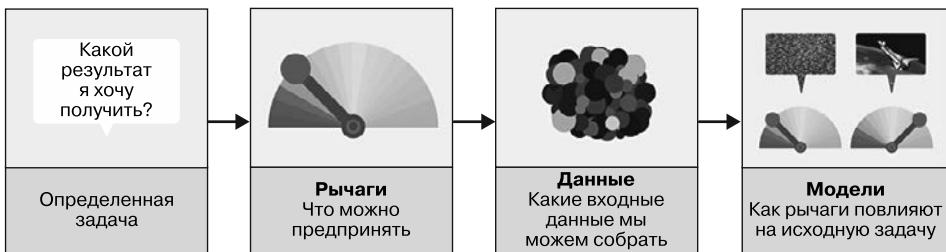


Рис. 2.2. Трансмиссионный подход

Рассмотрим модель в автономном транспортном средстве: вы хотите помочь автомобилю безопасно доехать из пункта А в пункт Б без человеческого вмешательства. Важной частью решения этой задачи будет грамотное предиктивное моделирование, но строится все не только на нем. По мере того как продукты становятся все сложнее, оно теряется в общей картине. Водители современных автономных машин даже не догадываются о том, что здесь стоят сотни (или даже тысячи) моделей и петабайты данных. Но так как ученые по работе с данными создают все более сложные продукты, возникает необходимость в систематическом подходе проектирования.

Мы используем данные не только для генерации еще большего количества данных (в виде прогнозов), но и для получения *действенных результатов*. В этом и состоит цель трансмиссионного подхода. Начните с определения четкой *задачи*. Например, Google при создании своего первого поискового движка рассматривал вопрос «Какую основную цель преследует пользователь при вводе поискового запроса?» Это позволило компании сформулировать четкую задачу: «Показывать максимально релевантные результаты поиска». Следующим шагом будет рассмотреть, какие *рычаги* вы можете применить (то есть какие действия предпринять) для достижения наилучшего результата. В случае с Google ответом было ранжирование результатов поиска. Третьим шагом компании было проанализировать, какие новые *данные* могут понадобиться, чтобы это ранжирование осуществить. Они поняли, что для этого можно использовать неявную информацию о том, какие страницы связаны с другими страницами.

Только совершив все эти три шага, мы начинаем думать о построении предиктивных *моделей*. Возможность построения тех или иных моделей определяется нашей задачей, доступными рычагами, а также имеющимися данными и теми, которые еще нужно собрать. Эти модели в качестве входных данных будут

получать рычаги и выбранные неконтролируемые переменные. А на основе полученных выводов можно будет прогнозировать итоговое состояние для нашей задачи.

Рассмотрим еще один пример: рекомендательные системы. *Задача* рекомендательного движка — обеспечить дополнительные продажи, удивляя клиента и угощая ему рекомендациями товаров, которые иначе он бы не купил. *Рычаг* в данном случае — это ранжирование рекомендаций. А чтобы генерировать рекомендации, *ведущие к повышению продаж*, необходимы новые *данные*. Для этого требуется проведение многих разнообразных экспериментов, которые позволяют собрать данные о широком спектре рекомендаций для большого числа клиентов. На деле этот шаг реализуют немногие организации, но без него у вас не будет информации, необходимой для оптимизации рекомендаций на основе вашей исходной задачи — увеличения продаж.

И наконец, вы можете создать две *модели* для оценки вероятностей покупок при условии просмотра рекомендации и без него. Разница между этими двумя вероятностями — это функция полезности для данной рекомендации клиенту. Ее показатель будет низок в случаях, когда алгоритм рекомендует знакомую книгу, от покупки которой клиент уже отказался (оба компонента будут малы), или книгу, которую он бы купил, даже без рекомендации (оба компонента будут велики и компенсируют друг друга).

Как вы видите, в реальности зачастую практическая реализация ваших моделей потребует гораздо больше, чем просто обучения модели. Вам нередко потребуется проводить эксперименты для сбора дополнительных данных и рассматривать, как вписать ваши модели в общую систему, которую вы разрабатываете. Что же касается данных, то сфокусируемся на вариантах их поиска для вашего проекта.

Сбор данных

Для многих типов проектов может получиться найти все необходимые данные онлайн. В этой главе мы будем создавать *детектор медведей*, который будет различать три их вида: гризли, черных и плюшевых. В интернете есть много подходящих изображений этих видов. Нужен только способ найти и загрузить их.

Мы предоставили инструмент, с помощью которого вы сможете это сделать, чтобы создать собственное приложение для распознавания изображений любого интересующего вас вида объектов. На курсе fast.ai тысячи студентов выставляли свои работы на форумах курса, показывая все что угодно, начиная с разновидностей колибри, обитающих в Тринидаде, и заканчивая типами автобусов в Панаме. Один студент даже создал приложение в помощь супруге в распознавании его 16 кузенов на семейных торжествах!

На момент написания Bing Image Search является наилучшим известным нам сервисом поиска и загрузки изображений. Он выполняет бесплатно до 1000 запросов в месяц, при этом каждый запрос может содержать до 150 изображений. Тем не менее за время, прошедшее с момента написания и до момента прочтения вами этой книги, могло появиться и что-то лучшее, так что не поленитесь посетить ее сайт (<https://book.fast.ai/>) для просмотра актуальных рекомендаций.



СЛЕДИТЕ ЗА НОВЕЙШИМИ СЕРВИСАМИ

Сервисы, которые можно использовать для создания датасетов, появляются и исчезают, и параллельно с этим постоянно меняются как их функции и интерфейсы, так и прайсы на услуги. В этом разделе мы продемонстрируем, как на момент написания книги можно было использовать Bing Image Search API (<https://oreil.ly/P8VtT>), доступный как часть Azure Cognitive Services (сервисы когнитивных вычислений Azure).

Для скачивания изображений с Bing Image Search зарегистрируйте бесплатный аккаунт в Microsoft. В результате вы получите ключ, который можно скопировать и ввести в ячейку следующим образом (заменив XXX на свой ключ и выполнив ее):

```
key = 'XXX'
```

Как вариант, если вы привыкли использовать командную строку, то можете установить ключ через терминал так:

```
export AZURE_SEARCH_KEY=ваsh_ключ
```

после чего перезапустить сервер Jupyter, ввести следующий код в ячейку и выполнить ее:

```
key = os.environ['AZURE_SEARCH_KEY']
```

После установки ключа вы можете использовать `search_images_bing`. Эта функция предоставляется небольшим классом `utils`, включенным в онлайн-блокноты (если вы не уверены, где именно эта функция определена, то ниже показано, как можно это выяснить, просто набрав ее в блокноте):

```
search_images_bing
<function utils.search_images_bing(key, term, min_sz=128)>
```

Проверим эту функцию:

```
results = search_images_bing(key, 'grizzly bear')
ims = results.attrgot('content_url')
len(ims)
150
```

Мы успешно загрузили URL-адреса 150 медведей-гризли (или, по крайней мере, изображения, найденные для этого запроса поиском Bing). Взглянем на одно из этих изображений:

```
dest = 'images/grizzly.jpg'  
download_url(ims[0], dest)  
  
im = Image.open(dest)  
im.to_thumb(128,128)
```



Вроде бы все сработало отлично, поэтому воспользуемся fastai-методом `download_images`, чтобы загрузить все URL-адреса для каждого из наших поисковых выражений, и разместим полученные изображения по соответствующим каталогам:

```
bear_types = 'grizzly','black','teddy'  
path = Path('bears')  
  
if not path.exists():  
    path.mkdir()  
for o in bear_types:  
    dest = (path/o)  
    dest.mkdir(exist_ok=True)  
    results = search_images_bing(key, f'{o} bear')  
    download_images(dest, urls=results.attrgot('content_url'))
```

Как мы и ожидали, теперь файлы изображений находятся в нашем каталоге:

```
fns = get_image_files(path)  
fns  
  
(#421) [Path('bears/black/00000095.jpg'), Path('bears/black/00000133.  
jpg'), Path('bears/black/00000062.jpg'), Path('bears/black/00000023.  
jpg'), Path('bears/black/00000029.jpg'), Path('bears/black/00000094.  
jpg'), Path('bears/black/00000124.jpg'), Path('bears/black/00000056.  
jpeg'), Path('bears/black/00000046.jpg'), Path('bears/black/00000045.jpg')...]
```

СЛОВО ДЖЕРЕМИ



Мне очень нравится работать в блокнотах Jupyter! Очень легко постепенно создавать то, что мне нужно, и проверять работу на каждом ее этапе. Я совершаю много ошибок, так что это мне очень помогает.

Зачастую при загрузке файлов из интернета некоторые из них оказываются поврежденными. Давайте это проверим:

```
failed = verify_images(fns)
failed

(#0) []
```

Для удаления всех некорректных изображений можно использовать `unlink`. Как и большинство функций fastai, возвращающих коллекцию, `verify_images` возвращает объект типа L, который включает метод `map`. Этот метод вызывает переданную функцию для каждого элемента коллекции:

```
failed.map(Path.unlink);
```

Здесь нужно иметь в виду, что модели могут отражать только те данные, на которых они обучались. При этом мир полон данных с искажениями, которые встречаются, например, в использованном нами поиске Bing. Предположим, вы хотите создать приложение, которое поможет пользователям определять, здоровая ли у них кожа, для чего вы обучили модель на результатах поиска по запросу, скажем, «здоровая кожа». На рис. 2.3 показан вариант полученных вами результатов.

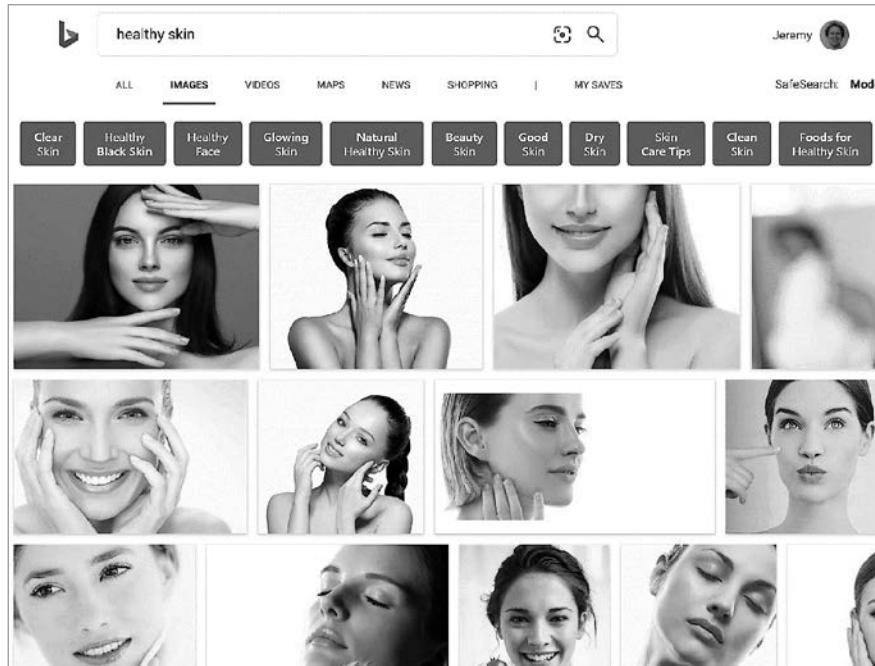


Рис. 2.3. Данные для детектора здоровой кожи

ПОМОЩЬ В БЛОКНОТАХ JUPYTER

Блокноты Jupyter отлично подходят для экспериментов и мгновенного просмотра результатов каждой функции, но в них также есть много сопутствующей функциональности, помогающей разобраться в тех или иных функциях. Можно даже непосредственно просматривать их исходный код. Допустим, вы вводите в ячейку следующее:

```
??verify_images
```

В результате будет выведено:

```
Signature: verify_images(fns)
Source:
def verify_images(fns):
    "Find images in `fns` that can't be opened"
    return L(fns[i] for i,o
             enumerate(parallel(verify_image, fns)) if not o)
File:      ~/git/fastai/fastai/vision/utils.py
Type:      function
```

Здесь говорится, какой аргумент данная функция принимает (`fns`), а затем приводится исходный код и имя содержащего его файла. Глядя на исходный код, мы можем видеть, что он параллельно применяет функцию `verify_image` и сохраняет только те файлы изображений, для которых результаты этой функции равны `False`, что согласуется со строкой документации: `Find the images in fns that can't be opened` (находит в `fns` изображения, которые не могут быть открыты).

Вот список некоторых других весьма полезных возможностей блокнотов Jupyter.

- Если вы не помните точное написание функции или имени аргумента, можете нажать `Tab` для автоподстановки.
- Из положения внутри скобок функции одновременное нажатие `Shift` и `Tab` отобразит окно с сигнатурой этой функции и кратким описанием. Двойное нажатие этих клавиш развернет документацию, а тройное откроет полное окно с той же информацией в нижней части экрана.
- Из положения внутри ячейки ввод и выполнение `?имя_функции` откроет окно с сигнатурой функции и ее кратким описанием.
- Из положения внутри ячейки ввод и выполнение `??имя_функции` откроет окно с сигнатурой функции, ее кратким описанием и исходным кодом.
- Если вы используете библиотеку `fastai`, то мы добавили в нее функцию `doc`: выполнение `doc(имя_функции)` в ячейке откроет окно с сигнатурой этой функции, ее кратким описанием, а также ссылками на исходный код на `Github` и полное описание этой функции в документации библиотеки (<https://docs.fast.ai/>).
- Не связанная с документацией, но очень полезная возможность: в любой момент для получения помощи при появлении ошибки введите `%debug` в следующую ячейку и выполните ее. Это откроет отладчик Python (<https://oreil.ly/RShnP>), который позволит вам просмотреть содержимое каждой переменной.

Если использовать это в качестве обучающих данных, то в итоге получится детектор не здоровой кожи, а *молодой белой женщины, дотрагивающейся до лица!* Постарайтесь хорошоенько продумать типы данных, которые следует ожидать в своем приложении на практике, а затем внимательно проверьте и убедитесь, что все эти типы отражаются в исходных данных для вашей модели. (Благодарим Деб Раджи (Deb Raji), которая придумала этот пример со здоровой кожей. Больше удивительных выводов на тему необъективности моделей вы можете найти в статье *Actionable Auditing: Investigating the Impact of Publicly Naming Biased Performance Results of Commercial AI Products* (https://oreil.ly/POS_C) («Действенный аудит: исследование влияния публичного объявления необъективных результатов работы коммерческих ИИ продуктов»).

Теперь, когда мы загрузили нужные данные, пора объединить их в формат, подходящий для обучения нашей модели. В fastai это означает создание объекта, называемого **DataLoaders**.

От данных к DataLoaders

DataLoaders — это тонкий класс, который только хранит любые передаваемые ему объекты **DataLoader** и делает их доступными в виде **train** и **valid**. Несмотря на свою простоту, это очень важный класс в fastai: он предоставляет данные для вашей модели. Ключевая функциональность **DataLoaders** обеспечивается всего четырьмя нижеприведенными строчками кода (у него есть еще некоторая менее значительная функция, которую мы пока опустим):

```
class DataLoaders(GetAttr):
    def __init__(self, *loaders): self.loaders = loaders
    def __getitem__(self, i): return self.loaders[i]
    train, valid = add_props(lambda i, self: self[i])
```



ТЕРМИН: DATALOADERS

Класс в fastai, хранящий несколько объектов **DataLoader**, которые вы в него передаете, — обычно ими являются **train** и **valid**, хотя допускается любое их количество. Первые два доступны в качестве свойств.

Позже вы узнаете о классах **Dataset** и **Datasets**, которые имеют такие же отношения. Для преобразования загруженных нами данных в объект **DataLoaders** нужно сообщить fastai по меньшей мере четыре следующих пункта:

- с какими видами данных мы работаем;
- как получить список элементов;
- как эти элементы разметить;
- как создать контрольную выборку.

До сих пор мы видели ряд *фабричных методов* для конкретных комбинаций этих пунктов, что удобно, когда у вас есть приложение и структура данных, подходящая к этим предопределенным методам. Для иных случаев в fastai есть очень гибкая система под названием *data block API*. С помощью этого API вы можете полностью настроить каждый этап создания **DataLoaders**. Чтобы создать **DataLoaders** для только что загруженного датасета, потребуется сделать следующее:

```
bears = DataBlock(
    blocks=(ImageBlock, CategoryBlock),
    get_items=get_image_files,
    splitter=RandomSplitter(valid_pct=0.2, seed=42),
    get_y=parent_label,
    item_tfms=Resize(128))
```

Давайте рассмотрим каждый из аргументов по очереди. Сначала мы предоставляем кортеж, определяющий типы, нужные нам для независимой и зависимой переменных:

```
blocks=(ImageBlock, CategoryBlock)
```

Независимая переменная — это та, на основе которой мы делаем прогнозы, а *зависимая* — это наша цель. В данном случае независимой переменной выступает набор изображений, а зависимой — категории (виды медведей) для каждого из этих изображений. На протяжении оставшейся части книги мы увидим также много других типов блоков данных.

Для этого **DataLoaders** нашими определяющими элементами будут пути файлов. Нужно сообщить fastai, как получить список этих файлов. Функция `get_image_files` получает путь и возвращает список всех соответствующих ему изображений (по умолчанию рекурсивно).

```
get_items=get_image_files
```

Зачастую в загружаемых датасетах контрольная выборка уже будет определена. Иногда это делается путем размещения изображений для обучающей и контрольной выборок в разные каталоги. Иногда для этого предоставляется CSV-файл, где рядом с каждым перечисленным именем файла указано, к какой выборке он относится. Это можно выполнить и многими другими способами, fastai же обеспечивает общий подход, позволяющий вам использовать для этого один из ее предопределенных классов или написать собственный.

В данном случае нужно разделить нашу обучающую и контрольную выборки произвольно. Но при каждом запуске этого блокнота мы хотели бы получать одинаковое разделение на обучающую/контрольную выборки, поэтому мы исправим случайное *начальное число* (на самом деле компьютеры не знают, как

создавать случайные числа, а просто создают списки чисел, которые выглядят случайными; если вы будете каждый раз предоставлять для списка одну и ту же стартовую точку, называемую начальным числом, то будете всякий раз получать одинаковый список).

```
splitter=RandomSplitter(valid_pct=0.2, seed=42)
```

Независимую переменную часто обозначают как *x*, а зависимую как *y*. Здесь мы сообщаем fastai, какую функцию вызвать для создания меток в нашем датасете:

```
get_y=parent_label
```

parent_label — это функция fastai, которая просто получает имя каталога, в котором находится файл. Так как мы помещаем каждое из изображений медведей в каталоги на основе их видов, то в итоге эта функция обеспечит их необходимыми метками.

Все наши изображения имеют разный размер, что в глубоком обучении является проблемой: мы передаем модели не по одному изображению, а сразу несколько (так называемый *мини-пакет*). Чтобы сгруппировать их в большой массив (обычно называемый *тензор*), который будет проходить через нашу модель, они все должны иметь одинаковый размер. Итак, нам нужно добавить преобразование, которое изменит размер изображений до одинакового. *Преобразователи элементов* — это части кода, выполняемые для каждого отдельного элемента, будь то изображение, категория или другое. fastai включает много предопределенных преобразователей. Мы же используем преобразователь *Resize* и указываем размер 128 пикселей:

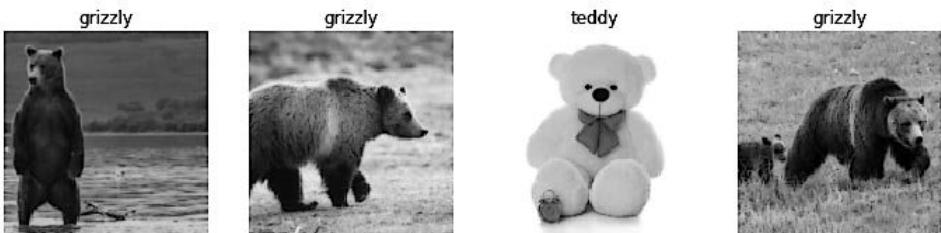
```
item_tfms=Resize(128)
```

Эта команда создала для нас объект *DataBlock*, который является своеобразным *шаблоном* для создания *DataLoaders*. Нам все еще нужно сообщить fastai фактический источник наших данных: в этом случае им будет путь, ведущий к изображениям:

```
dls = bears.dataloaders(path)
```

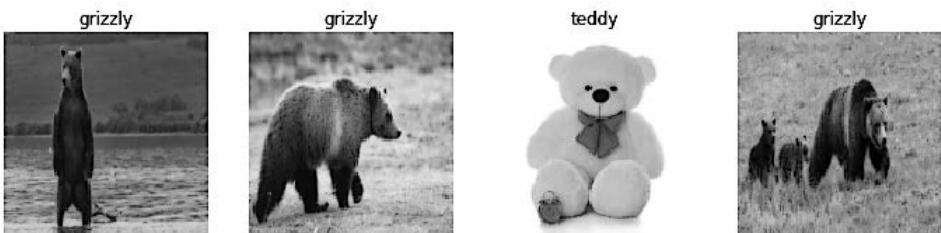
DataLoaders включает контрольный и обучающий *DataLoader*. *DataLoader* — это класс, поочередно предоставляющий GPU пакеты из нескольких элементов. Мы будем более подробно разбирать этот класс в следующей главе. При каждом прохождении через *DataLoader* fastai будет выдавать вам по 64 (по умолчанию) элемента, составленных в тензор. Мы можем взглянуть на некоторые из них, вызвав для *DataLoader* метод *show_batch*:

```
dls.valid.show_batch(max_n=4, nrows=1)
```

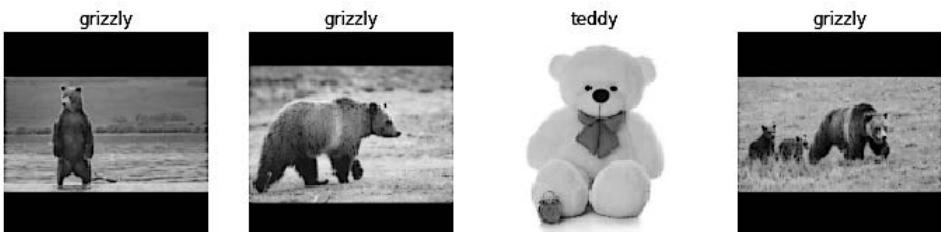


По умолчанию `Resize` *обрезает* изображения, чтобы вписать их в квадратную форму согласно установленному размеру, используя при этом полную ширину и высоту. В результате могут утратиться некоторые важные детали. В качестве альтернативы можно попросить fastai заполнить изображения нулями (черным) или сжать/растянуть их:

```
bears = bears.new(item_tfms=Resize(128, ResizeMethod.Squish))
dls = bears.dataloaders(path)
dls.valid.show_batch(max_n=4, nrows=1)
```



```
bears = bears.new(item_tfms=Resize(128, ResizeMethod.Pad, pad_mode='zeros'))
dls = bears.dataloaders(path)
dls.valid.show_batch(max_n=4, nrows=1)
```



Все эти подходы выглядят проблематичными и не очень эффективными. Если мы сожмем или растянем изображения, то в итоге они обретут нереалистичную форму, и модель обучится прогнозировать искаженные представления

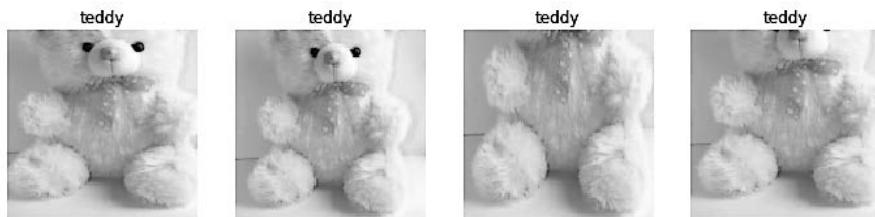
элементов, что снизит ее итоговую точность. Если же мы обрежем изображения, то удалим некоторые признаки, которые позволяют модели их распознавать. Например, если мы постараемся распознать породы собак или кошек, то рискуем обрезать ключевую часть тела или морды, необходимую для обнаружения различий между похожими породами. Если же заполнить изображения, то будет много пустого пространства, а это приведет к излишней трате вычислительной мощности и снизит разрешение фактически используемой части изображения.

Вместо этого обычно на практике мы произвольно выбираем часть изображения и затем обрезаем его до этой части. В каждой эпохе (которая представляет собой один полный просмотр моделью всех изображений датасета) мы произвольно выбираем разные части каждого изображения. Это означает, что наша модель может учиться фокусироваться на разных признаках изображений и распознавать их. Этот принцип также отражает то, как работают изображения в реальном мире: разные фотографии одного и того же предмета могут быть оформлены по-разному.

В действительности полностью необученная нейронная сеть абсолютно ничего не знает о поведении изображений. Она даже не распознает, что объект, повернутый всего на один градус, является все тем же объектом. Поэтому обучение нейросети примерами, на которых объекты находятся в немного других местах и имеют немного другой размер, помогает ей понять основную концепцию того, что такое объект и как он может быть представлен на изображении.

Вот еще один пример, в котором мы заменяем `Resize` на `RandomResizedCrop` — преобразователь, обеспечивающий только что описанное поведение. Самый важный параметр для передачи в него — это `min_scale`, который определяет, какую часть изображения нужно выбирать при минимальном масштабе:

```
bears = bears.new(item_tfms=RandomResizedCrop(128, min_scale=0.3))
dls = bears.dataloaders(path)
dls.train.show_batch(max_n=4, nrows=1, unique=True)
```



Здесь мы использовали `unique=True`, чтобы одно и то же изображение повторилось с разными вариантами преобразователя `RandomResizedCrop`.

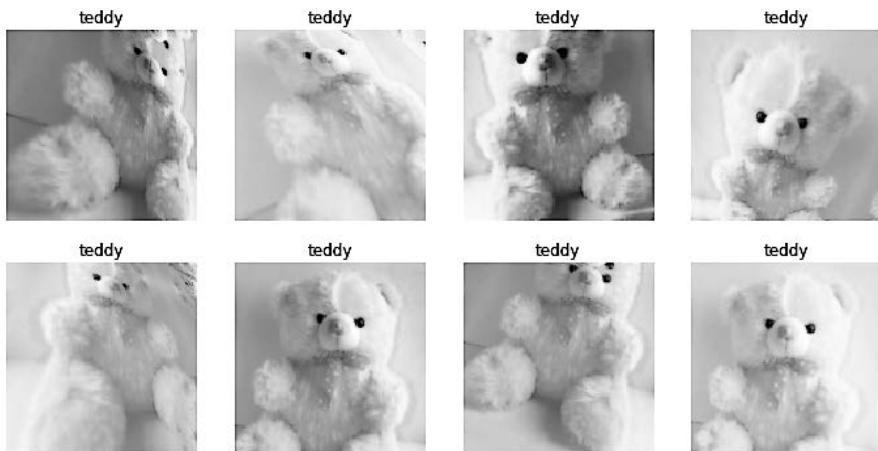
`RandomResizedCrop` — это конкретный пример более общей техники, называемой аугментацией данных.

Аугментация данных

Аугментация данных означает создание случайных вариаций входных данных так, чтобы они выглядели иначе, но сохраняли смысл. Примерами распространенных техник аугментации данных в отношении изображений являются вращение, переворачивание, искажение перспективы, изменение яркости и контраста. Для естественных фотоизображений, таких как мы используем здесь, стандартный, хорошо показавший себя набор аугментаций предоставляется функцией `aug_transforms`.

Поскольку теперь все наши изображения одного размера, мы можем применить эти аугментации ко всему их пакету, используя GPU, что сэкономит нам уйму времени. Чтобы сообщить `fastai`, что мы хотим использовать эти преобразователи для пакета, мы используем параметр `batch_tfms` (обратите внимание, что в этом примере мы уже не используем `RandomResizedCrop`, чтобы вы могли более отчетливо увидеть различия; кроме того, по той же причине мы задействуем вдвое больший объем аугментации, чем определен по умолчанию):

```
bears = bears.new(item_tfms=Resize(128), batch_tfms=aug_transforms(mult=2))
dls = bears.dataloaders(path)
dls.train.show_batch(max_n=8, nrows=2, unique=True)
```



Теперь, когда мы объединили наши данные в подходящий для модели формат, пора приступить к обучению классификатора изображений с их помощью.

Обучение модели и ее использование для чистки данных

Пришло время использовать те же строчки кода, что и в главе 1, но теперь уже для обучения классификатора медведей. У нас не так уж много данных для

решения этой задачи (максимум по 150 изображений каждого вида), поэтому для обучения модели мы будем использовать `RandomResizedCrop`, стандартный размер изображения 224 пикселя и предустановленную функцию `aug_transforms`:

```
bears = bears.new(
    item_tfms=RandomResizedCrop(224, min_scale=0.5),
    batch_tfms=aug_transforms())
dls = bears.dataloaders(path)
```

Теперь можно создать `Learner` и тонко настроить его обычным способом:

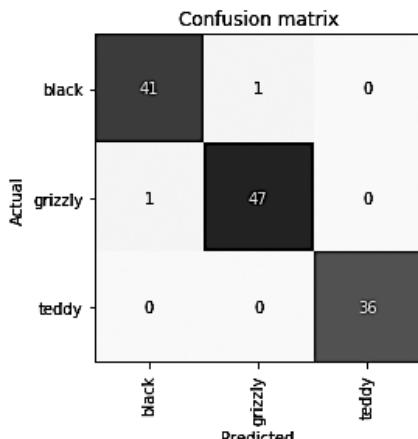
```
learn = cnn_learner(dls, resnet18, metrics=error_rate)
learn.fine_tune(4)
```

epoch	train_loss	valid_loss	error_rate	time
0	1.235733	0.212541	0.087302	00:05

epoch	train_loss	valid_loss	error_rate	time
0	0.213371	0.112450	0.023810	00:05
1	0.173855	0.072306	0.023810	00:06
2	0.147096	0.039068	0.015873	00:06
3	0.123984	0.026801	0.015873	00:06

Посмотрим, относятся ли ошибки нашей модели к тому, что она считает некоторых гризли плюшевыми медведями (это было бы очень опасно!) или черными, а может, наоборот, и т. д. Можно создать *матрицу ошибок*:

```
interp = ClassificationInterpretation.from_learner(learn)
interp.plot_confusion_matrix()
```



Строки здесь отражают всех черных, гризли и плюшевых медведей нашего датасета соответственно. Столбцы представляют изображения, которые модель определила как черных, гризли и плюшевых медведей соответственно. Из этого следует, что диагональ матрицы показывает изображения, которые были классифицированы верно, остальные же ее ячейки отражают неверную классификацию. Это один из многих способов, которыми fastai позволяет вам просматривать результаты обучения модели. Вычисляется матрица (конечно же!) на основе контрольной выборки. При использовании цветокодирования целью является закрашивание всех ячеек белым, за исключением диагонали, где нам нужен темно-синий. Заметьте, что наш классификатор почти не ошибается!

Но полезно взглянуть, где именно происходят те немногие его ошибки, чтобы понять, является причиной датасет (например, изображения, не являющиеся медведями, или ошибочно размеченные) или модель (возможно, она неверно обрабатывает изображения, полученные при необычном освещении, под другим углом и пр.). Для этого можно отсортировать изображения по их потерям.

Потери — это число, которое увеличивается, когда модель ошибается (особенно если при этом она еще и уверена в некорректном ответе) либо когда не ошибается, но при этом в правильности ответа не уверена. В начале части 2 мы углубленно изучим, как потери вычисляются и используются в процессе обучения. Сейчас же `plot_top_losses` показывает нам изображения датасета с максимальными потерями. Согласно заголовку вывода, каждое изображение размечено по четырем пунктам: прогноз, реальность (целевая метка), потери и вероятность. *Вероятность* здесь — это уровень уверенности, исчисляемый от 0 до 1 и присваиваемый моделью каждому прогнозу:

```
interp.plot_top_losses(5, nrows=1)
```



Этот вывод показывает, что изображением с максимальными потерями является то, которое с высокой степенью уверенности было спрогнозировано как «гризли». Тем не менее согласно поиску Bing оно размечено как «черный». Мы не эксперты по медведям, но даже на наш взгляд эта метка однозначно выглядит неверной, значит, ее следует поменять на «гризли».

Разумнее всего выполнять подобную чистку данных *перед началом обучения* модели. Но в этом случае вы видели, что модель может помочь вам обнаружить

проблемы в данных быстрее и проще. Поэтому мы обычно предпочитаем обу чать сначала быструю и простую модель, а затем использовать ее для помощи в чистке данных.

fastai содержит для этого процесса удобный GUI, называемый `ImageClassifierCleaner`, который позволяет выбирать категорию, а также обучающую и контрольную выборки и в итоге отобразить изображения с максимальными потерями в упорядоченном виде. При этом данный GUI здесь же предлагает меню, с помощью которого можно выбирать изображения для удаления или изменения меток:

```
cleaner = ImageClassifierCleaner(learn)
cleaner
```



Здесь мы видим, что среди «черных медведей» присутствует изображение с двумя животными: черным медведем и гризли. Следовательно, в меню под этим изображением мы выбираем `<Delete>`. `ImageClassifierCleaner` не выполняет за вас удаление или изменение меток, он просто возвращает индексы элементов для изменения. Например, чтобы удалить (`unlink`) все выбранные для этого изображения, нужно выполнить следующую инструкцию:

```
for idx in cleaner.delete(): cleaner.fns[idx].unlink()
```

Если же нам нужно переместить выбранные изображения в другую категорию, то мы выполняем:

```
for idx,cat in cleaner.change(): shutil.move(str(cleaner.fns[idx]), path/cat)
```

СЛОВО СИЛЬВЕЙНУ



Чистка данных и подготовка их для модели — это две сложнейшие задачи для аналитика. Некоторые утверждают, что на это уходит до 90 % всего времени. Библиотека fastai в этом смысле стремится предоставить инструменты, максимально упрощающие процесс.

Мы еще встретимся с примерами чистки данных с помощью модели. Сейчас же, закончив эту процедуру, можно обучить нашу модель повторно. Попробуйте сделать это сами и посмотрите, повысится ли итоговая точность.



НЕТ НЕОБХОДИМОСТИ В БОЛЬШИХ ДАННЫХ

После очистки датасета мы обычно наблюдаем в этом задании 100%-ную точность. Такой результат мы видим, даже когда загружаем меньше 150 изображений для каждой категории. Это с очевидностью доказывает несостоятельность распространенного мнения, что *для работы с глубоким обучением вам необходимо огромное количество данных!*

Теперь, когда мы обучили нашу модель, посмотрим, как можно развернуть ее для использования на практике.

Превращение модели в онлайн-приложение

Разберем, что нужно для преобразования полученной модели в рабочее онлайн-приложение. Мы ограничимся созданием базового рабочего прототипа, так как в книге не предусмотрено обучение всем деталям разработки веб-приложений.

Использование модели для вывода

После создания модели, которая вас устраивает, нужно сохранить ее, чтобы затем скопировать на сервер, где вы будете использовать ее в работе. Помните, что модель состоит из двух частей: *архитектуры* и обученных *параметров*. Простейший способ — это сохранить и то и другое, так как в этом случае при дальнейшей загрузке модели вы будете уверены, что данные компоненты соответствуют друг другу. Для сохранения этих двух частей используйте метод `export`.

Этот метод сохраняет даже определение того, как создавать `DataLoaders`. Это важно, потому что в противном случае вам придется повторно прописывать, как преобразовывать данные для использования модели в продакшне. Fastai по умолчанию автоматически использует для вывода `DataLoader` вашей контрольной выборки, поэтому аугментация данных применена не будет, что, как правило, и требуется.

При вызове `export` fastai сохранит файл под названием `export.pkl`:

```
learn.export()
```

Убедимся, что этот файл есть, использовав метод `ls`, который fastai добавляет в Python-класс `Path`:

```
path = Path()  
path.ls(file_exts='.pkl')  
  
(#1) [Path('export.pkl')]
```

Этот файл потребуется там, где вы будете развертывать приложение. Сейчас же создадим простое приложение внутри нашего блокнота.

Когда мы используем модель для получения прогнозов, а не обучения, мы называем это *выводом*. Для создания файла, обучающегося из экспортированного, мы используем `load_learner` (в этом случае это не обязательно, поскольку у нас уже есть рабочий `Learner` в блокноте; но так мы демонстрируем вам весь процесс от и до):

```
learn_inf = load_learner(path/'export.pkl')
```

Когда мы делаем вывод, то обычно получаем прогнозы только для одного изображения за раз. Чтобы это сделать, передайте имя файла в `predict`:

```
learn_inf.predict('images/grizzly.jpg')  
('grizzly', tensor(1), tensor([9.0767e-06, 9.9999e-01, 1.5748e-07]))
```

В итоге мы получили три элемента: прогнозированную категорию в том же формате, что был передан изначально (в данном случае это была строка), индекс прогнозированной категории и вероятности каждой категории. Последние два показателя основаны на последовательности категорий в `vocab` (словаре) `DataLoaders`, который представляет собой сохраненный список всех возможных категорий. Во время вывода вы можете обратиться к `DataLoaders` как к атрибуту `Learner`:

```
learn_inf.dls.vocab  
(#3) ['black', 'grizzly', 'teddy']
```

Можно заметить, что если мы обратимся к словарю с помощью индекса целого числа, возвращенного функцией `predict`, то, как и ожидается, получим «гризли». Обратите внимание, что если мы обратимся по индексу к списку вероятностей, то увидим вероятность того, что это гризли, равной примерно 1.

Нам известно, как делать прогнозы из сохраненной модели, поэтому у нас есть все необходимое для начала построений приложения, и все это можно сделать прямо в блокноте Jupyter.

Создание в блокноте приложения на основе модели

Чтобы использовать модель в приложении, мы можем просто рассмотреть метод `predict` как стандартную функцию. Это значит, что создать на основе

модели приложение можно с помощью множества доступных для разработчиков фреймворков и техник.

Однако большинство data scientists не знакомы с миром разработки веб-приложений. Попробуем использовать что-то вам в этом смысле известное: оказывается, что мы можем создать полноценное рабочее веб-приложение с помощью одних только блокнотов Jupyter. Для этого нужны всего две вещи:

- виджеты IPython (`ipywidgets`);
- Voilà.

Виджеты IPython – это компоненты GUI, которые объединяют функциональность JavaScript и Python в браузере и могут быть созданы в блокноте Jupyter. Например, модель для чистки изображений, которую мы недавно использовали, полностью написана с помощью виджетов IPython. Но мы не хотим вынуждать наших пользователей применять Jupyter самостоятельно.

Именно для этого и создана *Voilà* – система, которая делает приложения, состоящие из виджетов IPython, доступными для пользователей без помощи Jupyter. Блокнот *уже* является видом веб-приложения, только зависящим от другого веб-приложения: самого Jupyter. По сути, *Voilà* помогает автоматически преобразовывать сложное веб-приложение, которое мы уже неявно создали (блокнот), в упрощенную и более удобно развертываемую версию, которая функционирует уже как обычное веб-приложение, а не как блокнот.

Но у нас по-прежнему сохраняется преимущество ведения разработки в блокноте, так что с помощью `ipywidgets` мы можем собирать наш GUI шаг за шагом. Мы используем этот подход для создания простого классификатора изображений, и сначала нам понадобится виджет загрузки файлов:

```
btn_upload = widgets.FileUpload()  
btn_upload
```

Upload (0)

Теперь можно получить изображение:

```
img = PILImage.create(btn_upload.data[-1])
```



Для его отображения можно использовать виджет `Output`:

```
out_pl = widgets.Output()
out_pl.clear_output()
with out_pl: display(img.to_thumb(128,128))
out_pl
```



Теперь можно получить прогнозы:

```
pred,pred_idx,probs = learn_inf.predict(img)
```

И использовать `Label` для их показа:

```
lbl_pred = widgets.Label()
lbl_pred.value = f'Prediction: {pred}; Probability: {probs[pred_idx]:.04f}'
lbl_pred
```

```
Prediction: grizzly; Probability: 1.0000
```

Нам понадобится кнопка для запуска классификации. Выглядеть она будет в точности как кнопка `Upload`:

```
btn_run = widgets.Button(description='Classify')
btn_run
```

Нам также нужен обработчик события щелчка, то есть функция, которая будет вызываться при нажатии кнопки. Можно просто скопировать предыдущие строки кода:

```
def on_click_classify(change):
    img = PILImage.create(btn_upload.data[-1])
    out_pl.clear_output()
    with out_pl: display(img.to_thumb(128,128))
    pred,pred_idx,probs = learn_inf.predict(img)
    lbl_pred.value = f'Prediction: {pred}; Probability: {probs[pred_idx]:.04f}'

btn_run.on_click(on_click_classify)
```

Можете проверить работоспособность кнопки, щелкнув по ней. В итоге вы должны увидеть автоматическое обновление изображения и прогноза.

Теперь мы можем поместить их все в вертикальную рамку (`Vbox`), завершая, таким образом, наш GUI:

```
VBox([widgets.Label('Select your bear!'),  
      btn_upload, btn_run, out_pl, lbl_pred])
```



Вот мы и написали весь нужный код для нашего приложения. Преобразуем его в развертываемый вариант.

Превращение блокнота в реальное приложение

Теперь, когда в блокноте все работает, можно переходить к созданию приложения. Для этого начните новый блокнот и добавьте в него только код, необходимый для создания и показа нужных виджетов, а также Markdown для любого текста, который вы хотите отображать. Можете заглянуть в блокнот *bear_classifier* в репозитории книги, где вы увидите простое приложение, которое создали мы.

Далее, если вы этого еще не сделали, установите Voilà, скопировав нижеприведенные строки в ячейку блокнота и выполнив ее:

```
!pip install voila  
!jupyter serverextension enable voila --sys-prefix
```

Ячейки, начинающиеся со знака !, содержат не код Python, а код, переданный в ячейку (bash, Windows PowerShell и т. д.). Если вы уверенно пользуетесь командной строкой, которую мы еще будем подробно обсуждать, то можете просто ввести эти две строки (без !) прямо в терминале. В этом случае первая строка установит библиотеку *voila* и приложение, а вторая подключит его к вашему блокноту Jupyter.

Voilà запускает блокноты так же, как и сервер блокнота Jupyter, но делает и кое-что очень важное: удаляет все вводы ячеек и показывает только вывод (включая iPython widgets) вместе с ячейками Markdown. Так что остается только веб-приложение! Чтобы просмотреть ваш блокнот как веб-приложение Voilà, замените слова *notebooks* в URL браузера на *voila/render*. В результате вы увидите то же содержимое, что и в блокноте, но без каких-либо ячеек кода.

Конечно же, не обязательно использовать Voilà или ipywidgets. Ваша модель – это просто функция, которую вы можете вызывать (`pred, pred_idx, probs = learn.predict(img)`), что позволяет использовать ее с любым фреймворком, размещенным на любой платформе. Вы можете взять созданный в ipywidgets и Voilà прототип и позднее преобразовать его в обычное веб-приложение. Мы показываем вам здесь этот подход, потому что, на наш взгляд, это отличный способ для специалистов по данным и других мало знакомых с веб-разработкой энтузиастов создавать приложения на основе моделей.

Итак, приложение у нас есть – давайте его развернем!

Развертывание приложения

Как вы теперь знаете, для обучения практически любой полезной модели глубокого обучения вам необходим GPU. А нужен ли он вам для использования этой модели в работе? Нет! Вам почти однозначно *не нужен GPU для обслуживания вашей модели в продакшнене*, и тому есть ряд причин:

- Как мы видели, GPU полезны, только когда необходимо выполнять много идентичных задач параллельно. Если вы выполняете, например, классификацию изображений, то чаще всего классифицируете только одно изображение за раз, а обработка одного изображения обычно не требует столько вычислительной мощности, чтобы использование GPU было оправдано с точки зрения энергоресурсов. Поэтому в плане энергоэффективности в таких случаях чаще лучше использовать CPU.
- В качестве альтернативы можно подождать, пока несколько пользователей отправят изображения, а затем объединить их и обработать одновременно с помощью GPU. Но тогда вы вынудите пользователей ждать. Кроме того, для оправданности такого подхода ваш сайт должен быть достаточно посещаем. Если же вам эта функциональность все же нужна, можете использовать такие инструменты, как ONNX Runtime (<https://oreil.ly/nj-6f>) от Microsoft или AWS SageMaker (<https://oreil.ly/ajcap>).
- При использовании GPU для вывода неизбежны существенные сложности. В частности, память GPU потребует осторожного ручного управления, и вам понадобится внимательно продуманная система очередей, чтобы гарантировать обработку только одного пакета изображений за раз.
- Конкуренция на рынке предложений ресурсов CPU существенно выше, чем в случае с GPU, поэтому и стоимость использования CPU-серверов обычно ниже.

Из-за сложности обслуживания GPU возникло множество систем, нацеленных на автоматизацию этого процесса. Тем не менее управление этими системами и их выполнение тоже представляет свои сложности и, как правило, требует

компиляции вашей модели в разные подходящие для того или иного их варианта формы. Обычно лучше не связываться с такой сложностью до тех пор, пока ваше приложение не достигнет достаточного уровня популярности, чтобы в этом возник отчетливый финансовый смысл.

Если же речь идет о вашем начальном прототипе или любом хобби-проекте, который вы хотите показать миру, то их можно без проблем разместить бесплатно. Оптимальное место и способ для этого время от времени будут меняться, так что следите за сайтом книги, чтобы иметь наиболее актуальную информацию. На момент написания этой книги, а именно в начале 2020 года, простейшим (и бесплатным) подходом будет использование Binder. Чтобы опубликовать ваше приложение на этом ресурсе, выполните следующее.

1. Добавьте блокнот в репозиторий GitHub (<http://github.com/>).
2. Вставьте URL этого репозитория в поле для URL ресурса Binder, как показано на рис. 2.4.
3. В раскрывающемся списке меню URL to open выберите URL.
4. В пункте URL to open введите /voila/render/имя.ipynb (заменив имя называнием вашего блокнота).
5. Нажмите кнопку буфера обмена справа внизу, чтобы скопировать этот URL, после чего вставьте его в безопасное место.
6. Нажмите кнопку Launch.

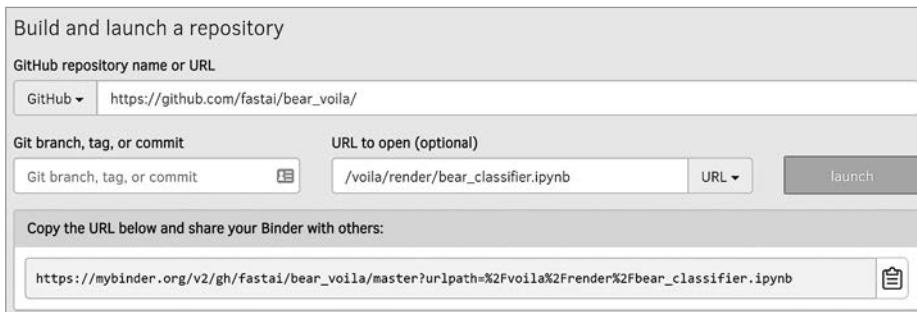


Рис. 2.4. Развёртывание в Binder

При первом выполнении этой процедуры Binder потребуется около пяти минут на создание вашего сайта. За это время он найдет виртуальную машину, которая сможет запускать приложение, выделит хранилище и соберет файлы, необходимые для Jupyter и для представления вашего блокнота в виде веб-приложения.

И наконец, после запуска приложения он перенаправит ваш браузер на страницу вашего нового приложения. Для предоставления к этому приложению доступа другим людям можете просто передавать им скопированный ранее URL.

Остальные (как платные, так и нет) варианты развертывания веб-приложения можете найти на сайте книги (<https://book.fast.ai/>).

Вы также можете захотеть развернуть приложение на мобильных устройствах или таких периферийных вариантах, как Raspberry Pi. Есть много библиотек и фреймворков, позволяющих вам интегрировать модель непосредственно в мобильное приложение. Тем не менее эти подходы обычно требуют выполнения многих дополнительных шагов и рутинного кода, а также не всегда поддерживают все слои PyTorch и fastai, которые могут использоваться в вашей модели. Кроме того, необходимые действия будут также зависеть от типов мобильных устройств, для которых вы делаете развертывание: вам может потребоваться проделать одни действия для запуска на устройствах iOS, другие — для запуска на новейших устройствах Android, третья — для устаревших версий этой ОС, и т. д. Вместо этого мы рекомендуем по возможности развертывать саму модель на сервере и уже потом подключать к ней ваше мобильное или периферийное приложение как веб-сервис.

Такой подход имеет немало плюсов. Начальный процесс установки упрощается, так как вам нужно развертывать только небольшое приложение GUI, которое подключается к серверу для выполнения всей тяжелой нагрузки. Более важно скорее то, что обновление этой ключевой логики может происходить на сервере, не требуя распространения на всех пользователей. У вашего сервера будет намного больше памяти и мощности, чем в большинстве периферийных устройств, и в случае роста требований модели масштабировать эти ресурсы станет куда проще. Оборудование, используемое вами на сервере, также будет более стандартным, а значит, и более поддерживаемым fastai и PyTorch, поэтому компилировать модель в другую форму вам не потребуется.

Но есть и минусы. Ваше приложение потребует использования сетевого подключения, и при каждом вызове модели будет возникать некоторая задержка. (Модели нейронной сети в любом случае требуются некоторое время для запуска, так что такая дополнительная сетевая задержка может стать незаметной для пользователей. В действительности благодаря использованию более производительного оборудования на сервере общая задержка может оказаться даже ниже, чем при локальном выполнении.) Кроме того, если в вашем приложении используются уязвимые данные, то пользователи могут беспокоиться о способе отправки данных на удаленный сервер, поэтому иногда потребность в конфиденциальности будет означать, что вам нужно выполнять модель на периферийном устройстве (этого можно попробовать избежать, используя собственный (on-premise) сервер, защищенный, например, брандмауэром компании). Управление сложностью и масштабированием сервера также может создать дополнительную нагрузку, в то время как при выполнении модели на периферийных устройствах каждый пользователь использует собственные вычислительные ресурсы, что упрощает масштабирование при увеличении числа пользователей (также называемое *горизонтальным масштабированием*).



СЛОВО АЛЕКСИСУ

Я лично наблюдал, как мобильный ландшафт ML менялся у меня на работе. Мы предлагаем приложение для iPhone, использующее компьютерное зрение, и в течение нескольких лет мы выполняли собственные модели компьютерного зрения в облаке. Тогда это был единственный возможный способ, поскольку для работы моделей нужно было много памяти и вычислительных ресурсов, а на обработку входных данных затрачивалось до нескольких минут. Такой подход требовал построения не только моделей, что было весело, но и инфраструктуры, обеспечивающей бесперебойное функционирование конкретного числа «вычислительных рабочих машин», подключение дополнительных машин при увеличении трафика, стабильное хранилище для объемных вводов и выводов, а также дающей возможность приложению iOS узнавать и сообщать пользователям о процессе выполнения их запросов. Сегодня Apple предоставляет API для преобразования моделей с целью эффективного выполнения на устройствах, и большинство iOS-устройств имеют отдельное оборудование для ML, поэтому именно такую стратегию мы используем для наших новых моделей. Но и это по-прежнему нелегко, хотя в нашем случае оно того стоит, так как обеспечивает ускоренную работу пользователям и не вызывает такой обеспокоенности работой серверов. Что лучше подойдет в вашем случае, будет зависеть от того, какой пользовательский опыт вы стремитесь организовать и что считаете более простым в реализации. Если вы хорошо разбираетесь в обслуживании серверов, выбирайте этот вариант. Если для вас не представляет сложности создание нативных мобильных приложений, занимайтесь этим. В этой ситуации можно пойти различными путями.

Рекомендуем по возможности использовать для этих задач сервер на базе CPU до тех пор, пока он будет отвечать вашим требованиям. Если вам удастся разработать достаточно успешное приложение, то у вас появится возможность оправдать вложения в более сложные подходы к развертыванию.

Поздравляем! Вы успешно создали модель глубокого обучения и развернули ее! Теперь самое время подумать о том, что может пойти не так.

Как избежать катастрофы

На практике модель глубокого обучения будет всего лишь частью гораздо более крупной системы. Как мы обсуждали в начале главы, построение продукта для обработки данных требует продумывания всего процесса от и до, начиная от его концепции и заканчивая использованием в работе. В этой книге мы никак не сможем рассмотреть всю сложность управления развернутыми проектами, подразумевающую контроль нескольких версий модели, A/B-тестирование, канареечные релизы, обновление данных (стоит ли продолжать постоянно наращивать датасеты или же лучше периодически удалять часть устаревших

данных?), обработку маркировки данных, отслеживание всех этих процессов, обнаружение отклонений в модели и т. д.

В этом разделе мы проведем краткий обзор некоторых из наиболее важных проблем. Для более подробного ознакомления со сложностями развертывания рекомендуем вам обратиться к работе Эммануэля Амейсина (*Emmanuel Ameyisin Building Machine Learning Powered Applications* (<http://shop.oreilly.com/product/>) («Построение приложения на основе машинного обучения»), выпущенной O'Reilly).

Особенность в том, что понять и протестировать поведение модели глубокого обучения гораздо труднее, чем в случае с большинством кода, который вы пишете. При обычной разработке ПО можно анализировать конкретные шаги, предпринимаемые этим ПО, и внимательно проверять, какие из них соответствуют запланированному вами поведению. Но в нейронной сети поведение не прописывается точно, а определяется стремлением модели соответствовать обучающим данным.

Это может привести к катастрофе! Например, предположим, что мы реально внедрили нашу систему обнаружения медведей, к которой подключают видеокамеры вокруг палаточных городков в национальных парках для предупреждения отдыхающих о приближении медведей. Если бы мы использовали модель, обученную на нашем датасете, то на практике возникло бы очень много сложностей. Например:

- работа с видеоданными, а не изображениями;
- обработкаочных изображений, которых может не быть в датасете;
- работа с изображениями камер низкого разрешения;
- обеспечение достаточно быстрой для эффективного реагирования обратной связи;
- распознавание медведей в положениях, редко наблюдаемых на фотографиях, размещенных людьми онлайн (например, сзади, частично перекрытыми кустами или далеко от камеры).

Большая часть этой проблемы в том, что фотографии, которые люди обычно загружают в интернет, ясно и мастерски отражают свою тематику, что нехарактерно для входных данных, которые система будет получать на практике. Поэтому для создания реально полезной системы нам может потребоваться проделать много собственной работы по сбору и разметке данных.

Это всего лишь один пример более общей проблемы *не соответствующих областей данных*. Иначе говоря, в работе модель может встречать данные, сильно отличающиеся от использованных в обучении. Для этой проблемы пока нет полноценного технического решения. Напротив, нужно быть осторожными при внедрении такого рода технологии.

Осторожными нужно быть и по другим причинам. Одна очень распространенная проблема — это *сдвиг области*, при котором тип встречаемых нашей моделью

данных со временем меняется. Например, страховая компания может использовать модель глубокого обучения как часть своего алгоритма ценообразования с учетом рисков, но со временем типы привлекаемых этой компанией клиентов и типы представляемых ею рисков могут настолько измениться, что исходные обучающие данные полностью утратят актуальность.

Не соответствующие области данных и сдвиг области являются примерами более крупной проблемы: того, что вы никогда не сможете полностью понять все возможные варианты поведения нейронной сети, так как в них используется слишком много параметров. Это естественный недостаток их лучшей черты — гибкости, которая позволяет им решать сложные задачи, в которых мы бы даже не смогли как следует определить предпочтительные варианты решения. Однако есть и хорошие новости, а именно то, что всегда доступны способы снизить эти риски, применяя осторожный и продуманный процесс. Детали в данном случае будут варьироваться в зависимости от деталей решаемой задачи, но мы постараемся продемонстрировать вам высокую уровень подход, в общем виде показанный на рис. 2.5. Надеемся, что он послужит для вас полезным ориентиром.

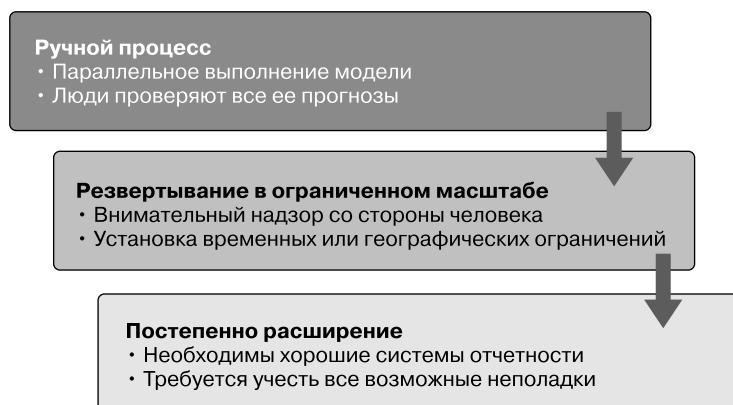


Рис. 2.5. Процесс развертывания

Первым шагом будет по возможности повсеместное внедрение полностью ручного процесса, при котором модель глубокого обучения выполняется параллельно, при этом не влияя на какие-либо действия непосредственно. Человек, участвующий в ручном процессе, должен просматривать выводы модели и проверять, насколько они актуальны. Например, в случае с нашим классификатором медведей смотритель парка мог бы использовать экран, отображающий видеопотоки со всех камер, где любые возможные появления медведей определялись бы красным. Смотритель в данной схеме оставался бы бдительным в той же степени, что и до внедрения модели. В этот период она просто помогает следить за порядком.

Вторым шагом будет попытка ограничить область действия модели и обеспечить ее тщательный контроль со стороны человека. Например, использовать модель на ограниченном пространстве в течение ограниченного времени. Вместо того чтобы внедрять классификатор медведей в каждый национальный парк страны, мы можем выбрать один наблюдательный пункт на период в одну неделю и поручить смотрителю парка проверять каждую объявляемую моделью тревогу.

Затем постепенно увеличивайте область внедрения. При этом убедитесь, что у вас есть хорошие системы отчетности, чтобы быть в курсе любых значительных изменений в действиях по сравнению с ручным процессом. Например, если после внедрения системы в определенном месте число тревог по случаю появления в этом районе медведей удвоится или станет вдвое меньше, вас это должно сильно обеспокоить. Постарайтесь продумать все возможные варианты, в которых ваша система могла бы ошибиться, а затем подумайте о том, какой показатель, отчет или картинка могли бы отразить эту проблему, и убедитесь, что ваша регулярная отчетность эту информацию включает.



СЛОВО ДЖЕРЕМИ

Двадцать лет назад я основал компанию, которая называлась Optimal Decisions. Она использовала ML и оптимизацию для помощи огромным страховым компаниям в определении цен на услуги, влияя на десятки миллиардов долларов страховых расходов. Мы использовали описанные здесь подходы для предотвращения потенциальных проблем в случае, если что-то пойдет не так. Кроме того, прежде чем работать с клиентами для внедрения какого-либо продукта в работу, мы старались симулировать его влияние путем тестирования всей системы на их данных прошлого года. Отправка этих алгоритмов в работу всякий раз сопровождалась нервной напряженностью, но в итоге каждое развертывание оказывалось успешным.

Непредвиденные последствия и петли обратной связи

Одна из серьезнейших сложностей при внедрении модели в том, что она может изменить поведение системы, в которую встраивается. Рассмотрите в качестве примера алгоритм «предиктивной полицейской деятельности», который прогнозирует повышение преступности в конкретных районах, вследствие чего туда направляется больше полицейских, что может привести к увеличению количества преступлений, и т. д. В работе Королевского статистического общества *To Predict and Serve?* (<https://oreil.ly/3YEWH>) («Прогнозировать и служить?») Кристиан Лам (Kristian Lum) и Уильям Исаак (William Isaac) отмечают, что «предиктивная полицейская деятельность соответствует своему названию: прогнозирует будущую деятельность полиции, а не будущие преступления».

В этом случае часть проблемы заключается в том, что при наличии смещения (о котором мы подробно поговорим в следующей главе) *петли обратной связи*

приводят к негативным следствиям в виде все большего искажения. Например, есть опасения, что это уже происходит в США, где наблюдается существенное смещение в показателях арестов на расовой почве.

Согласно данным Американского союза защиты гражданских свобод (<https://oreil.ly/A9ijk>) (ACLU), «несмотря на примерно равные показатели употребления, чернокожих в 3,73 раза чаще арестовывают за марихуану, чем белых». Влияние этого искажения наряду с внедрением алгоритмов предиктивной полицейской деятельности во многих частях США привело к тому, что Бари Уильямс (Bárá Williams) написал в New York Times (<https://oreil.ly/xR0di>): «Та же технология, которая вызывает столько восхищения в моей профессиональной деятельности, используется в силовых структурах такими способами, которые в ближайшие годы могут привести к тому, что мой сын, которому всего семь лет, с большой вероятностью будет профицирован или арестован, а самое страшное в этом то, что причиной может стать просто его раса или место нашего проживания».

Полезным упражнением перед внедрением важной системы ML будет ответ на вопрос: «Что случится, если она заработает действительно очень хорошо?» Другими словами, что, если эффективность прогнозирования будет чрезвычайно велика и его влияние на поведение окажется также очень значительным? На ком это в первую очередь отразится? Как могут потенциально выглядеть самые экстремальные результаты? Откуда вы можете узнать, что именно произошло?

Такие размышления могут помочь разработать более осторожный план внедрения с участием систем постоянного мониторинга и человеческого надзора. Конечно же, человеческий надзор окажется бесполезен, если на него не реагировать, поэтому убедитесь, что у вас есть надежные и устойчивые каналы связи, чтобы нужные люди вовремя оказались в курсе неисправностей и имели возможность их устраниить.

Записывайте!

По опыту наших студентов, один из самых эффективных способов закрепления пройденного материала — записать его. Вряд ли найдется лучший способ проверить собственное понимание темы, кроме как попытаться объяснить ее кому-то еще. Это полезно, даже если вы никогда никому свои записи не показываете, но будет еще лучше, если вы все же будете ими делиться. Поэтому если вы еще этого не делаете, то мы рекомендуем начать вести блог. Теперь, когда вы закончили эту главу и умеете уже не только обучать, но и развертывать модели, вы вполне можете написать свой первый пост о путешествии в страну глубокого обучения. Что вас удивило? Какие вы предвидите возможности для применения глубокого обучения в своей области? Какие ожидаете препятствия?

Рейчел Томас, соосновательница fast.ai, в статье *Why You (Yes, You) Should Blog* (<https://oreil.ly/X9-3L>) («Почему вы (да, именно вы) должны вести блог») написала следующее:

В первую очередь себе молодой я бы посоветовала как можно раньше начать вести блог, и вот почему.

- Это как резюме, только лучше. Я знаю людей, которые благодаря своим постам в блогах получили предложения о работе.
- Блог помогает учиться. Организация полученных знаний всегда помогает синтезировать собственные идеи. Один из способов проверить, понимаешь ли ты что-либо, — это попробовать объяснить это кому-то, для чего очень подходит пост в блоге.
- Благодаря своим постам я получала приглашения на конференции и выступления. Меня даже пригласили на саммит разработчиков TensorFlow (он был потрясающий!) после поста о том, как мне не нравится TensorFlow.
- Нетворкинг. Я встречалась с людьми, отвечающими на мои посты.
- Экономия времени. Всякий раз после нескольких ответов на один вопрос через e-mail вам следует оформить его в виде поста, чтобы в дальнейшем просто делиться им с интересующимися.

Наиболее же важный совет, пожалуй, такой:

Вы как никто другой способны помочь людям, находящимся на шаг позади, так как пройденный материал еще хорошо держится в вашей памяти. Многие эксперты забыли, что значит быть новичком (или находиться на среднем уровне), а также забыли, что тему сложно понять, когда ты слышишь ее впервые. Контекст вашего личного опыта, ваш особый стиль и ваш уровень знаний определят ту специфику, о которой вы будете писать.

В приложении А мы подробно расписали, как настроить свой блог. Если у вас еще его нет, то загляните туда прямо сейчас, так как у нас есть для вас реально хороший способ начать это делать бесплатно, без рекламы — и вы даже сможете использовать Jupyter Notebook!

Вопросник

1. Какие серьезные недостатки есть у текстовых моделей?
2. Какие возможны негативные социальные последствия при использовании моделей для генерации текста?
3. Какую альтернативу автоматизации процесса следует рассмотреть в ситуациях, где модель может совершать ошибки, способные навредить?

4. В обработке какого вида табличных данных глубокое обучение особенно хорошо?
5. Какой главный недостаток непосредственного использования модели глубокого обучения для рекомендательных систем?
6. Из каких шагов состоит трансмиссионный (Drivetrain) подход?
7. Какие шаги трансмиссионного подхода можно использовать в рекомендательной системе?
8. Создайте модель распознавания изображений, используя актуальные для вас данные, и разверните ее в интернете.
9. Что такое `DataLoaders`?
10. Какие четыре пункта нужно сообщить `fastai` для создания `DataLoaders`?
11. За что отвечает параметр `splitter` в `DataBlock`?
12. Как добиться того, чтобы случайное разделение всегда генерировало одинаковую контрольную выборку?
13. Какие буквы обычно используются для обозначения независимых и зависимых переменных?
14. В чем разница между обрезкой, заполнением и сжатием при изменении размера? Что и в каких случаях следует использовать?
15. Что такое аугментация данных? Зачем она нужна?
16. Приведите пример, в котором модель классификации медведей может плохо работать из-за структурных или стилистических отличий в обучающих данных.
17. В чем разница между `item_tfms` и `batch_tfms`?
18. Что такое матрица смещения?
19. Что сохраняет `export`?
20. Как называется ситуация, когда мы используем модель для прогнозирования вместо обучения?
21. Что такое виджеты IPython?
22. Когда следует использовать CPU для развертывания? В каких случаях GPU окажется предпочтительнее?
23. Каковы недостатки развертывания приложения на сервере, а не на клиентском (или периферийном) устройстве, таком как смартфон или PC?
24. Какие три вида проблем могут возникнуть при внедрении системы оповещения о появлении медведей на практике?
25. Что такое не соответствующие области данные?
26. Что такое сдвиг области?
27. Из каких трех шагов состоит процесс развертывания?

Дополнительные задания

1. Подумайте, как трансмиссионный подход отображается на интересующий вас проект или задачу.
2. Когда будет лучше избежать конкретных видов аугментации данных?
3. Проведите в отношении проекта, для которого собираетесь применить глубокое обучение, аналитический эксперимент «Что случится, если все заработает очень хорошо?»
4. Создайте блог и напишите свой первый пост. Например, расскажите, для чего, на ваш взгляд, глубокое обучение может быть полезно в интересующей вас области.

ГЛАВА 3

Этика данных

БЛАГОДАРНОСТЬ ДОКТОРУ РЕЙЧЕЛ ТОМАС

Эта глава была написана совместно с доктором Рейчел Томас, соучредителем fast.ai и директором-основателем центра прикладной этики данных в Университете Сан-Франциско. В этом центре в значительной степени следуют программы, разработанной ею для курса *Introduction to Data Ethics* («Введение в этику данных»).

Как мы уже говорили в главах 1 и 2, иногда модели ML могут ошибаться, содержать баги, а также встречаться с данными, которых ранее не видели, в связи с чем вести себя неожиданным образом. Бывает и так, что они работают в точности как предполагалось, но используются для того, для чего мы бы предпочли вообще никогда их не использовать.

Поскольку глубокое обучение является очень мощным инструментом, который можно использовать для множества задач, становится особенно важным учитывать последствия наших выборов. Философское учение об *этике* — это изучение добра и зла, включая само определение этих слов, распознавание правильных и неправильных поступков, а также понимание связи между этими поступками и их следствиями. Область *этики данных* образовалась уже давно, и в ней работают многие академики. С ее помощью определяется политика во многих юрисдикциях. Ее используют как малые, так и крупные компании, чтобы гарантировать положительное влияние разрабатываемых ими продуктов на общество. Кроме того, к этой области обращаются исследователи, которые желают убедиться в том, что их работа послужит именно во благо, а не наоборот.

Поэтому, практикуя глубокое обучение, вы наверняка окажетесь в ситуации, где этику данных нужно будет учесть. Что же такое этика данных? Это подраздел этики, поэтому начнем с нее.



СЛОВО ДЖЕРЕМИ

В университете философия этики была моим приоритетным направлением (она бы стала темой моей диссертации, если бы я закончил учебу, а не решил все бросить и уйти в реальный мир). На базе этих знаний могу сказать следующее: нет единого мнения о том, что такое хорошо и плохо, существуют ли эти понятия в принципе, как их следует определять, каких людей считать хорошими, а каких — плохими и т. д. Так что от теории многоного ждать не приходится, поэтому здесь мы сфокусируемся не на ней, а на реальных примерах.

Отвечая на вопрос, что такое этика (<https://oreil.ly/nyVh4>), Центр прикладной этики им. Марккулы приводит следующие определения:

- грамотно обоснованные стандарты хорошего и плохого, предписывающие поведение людей;
- изучение развития личных этических стандартов человека.

Здесь нет списка верных ответов, как и списка того, что можно, а что нельзя. Этика сложна и всегда зависит от контекста. В этом вопросе требуется рассматривать точки зрения многих участников. Этика — это мускул, который вам нужно развивать и прокачивать. В этой главе нашей целью будет предоставить несколько ориентиров, которые помогут вам на этом пути.

Выявление этических проблем лучше всего делать в рамках команды. Это единственный способ реально рассмотреть разные точки зрения на ситуацию. Разница прожитого опыта позволяет одним видеть то, чего не видят другие. Подобная работа с командой оказывается полезной для многих «нарашаивающих мышцы» занятий, включая и это.

Эта глава — далеко не единственная часть книги, где мы будем говорить об этике, но будет правильным рассмотреть данную тему более подробно именно в одном месте. Чтобы легче ориентироваться, стоит взглянуть на несколько примеров, и мы выбрали три, которые, на наш взгляд, эффективно отражают некоторые из ключевых тем.

Ключевые примеры этики данных

Начнем с трех специфичных примеров, иллюстрирующих распространенные этические проблемы в технологиях (несколько позже мы изучим эти проблемы более детально).

Процессы оказания помощи

Неисправные алгоритмы в сфере здравоохранения Арканзаса поставили пациентов в затруднительное положение.

Петли обратной связи

Рекомендательная система YouTube помогла вызвать бум интереса к теориям заговоров.

Необъективность

При поиске традиционного афроамериканского имени в Google отображается реклама проверки криминального прошлого.

Для каждого принципа мы приведем как минимум один пример. А вы, в свою очередь, подумайте о том, что могли бы сделать в этих ситуациях, проанализировав возможные сложности и способы их решения.

Баги и оказание помощи: неисправный алгоритм, распределявший медицинские льготы

The Verge исследовали ПО, использованное в более чем половине штатов США, чтобы определить, сколько медицинских услуг получают люди. Результаты были изложены в статье *What Happens When an Algorithm Cuts Your Healthcare* (<https://oreil.ly/5Ziok>) («Что случается, когда алгоритм урезает ваши медицинские льготы»). После применения этого алгоритма в Арканзасе у сотен людей (многие с серьезными нарушениями здоровья) резко сократилось медицинское обслуживание.

Например, у Тэмми Доббс (Tammy Dobbs), женщины с церебральным параличом, нуждающейся в помощнике, чтобы встать с кровати, принять ванну и т. д., время получаемой помощи внезапно сократилось на 20 часов в неделю. Объяснений этому сокращению она не получила. Только в результате судебного разбирательства выяснилось, что в программной реализации алгоритма были ошибки, что негативно сказалось на людях с диабетом или церебральным параличом. Теперь Доббс и многие другие люди, пользующиеся подобными льготами здравоохранения, живут с опасением, что эти льготы могут опять внезапно урезать по необъяснимым причинам.

Петли обратной связи: рекомендательная система YouTube

Петли обратной связи могут возникнуть, когда ваша модель определяет следующий круг получаемых данных. Данные, возвращаемые быстро, искажаются самим ПО.

Например, у YouTube 1,9 миллиарда пользователей, которые просматривают более 1 миллиарда часов видео в день. Алгоритм их рекомендательной системы (разработанный Google), который был спроектирован для оптимизации времени просмотра, определяет около 70 % просматриваемого контента. Но возникла

проблема: он привел к неуправляемым петлям обратной связи, на что New York Times в феврале 2019 года отреагировала заголовком *YouTube Unleashed a Conspiracy Theory Boom. Can It Be Contained?* (<https://oreil.ly/Lt3aU>) («YouTube развязал бум теории заговоров. Можно ли его сдержать?»). Предполагается, что рекомендательные системы прогнозируют контент, который понравится людям, но они также во многом определяют то, какой контент люди вообще увидят.

Предвзятость: «арест» профессора Латаньи Суни

Доктор Латанья Суни — профессор Гарвардского университета и директор университетской лаборатории конфиденциальных данных. В своей работе *Discrimination in Online Ad Delivery* (<https://oreil.ly/1qBxU>) («Дискриминация в онлайн-объявлениях») она описывает, как поиск ее имени в Google приводил к появлению объявлений *Latanya Sweeney, Arrested?* («Латанья Суни арестована?»), несмотря на то что она единственная известная Латанья Суни и арестам никогда не подвергалась. Но погуглив другие имена, такие как *Kirsten Lindquist*, она получила более нейтральные объявления, хотя Кирстен Линдквист была арестована трижды.

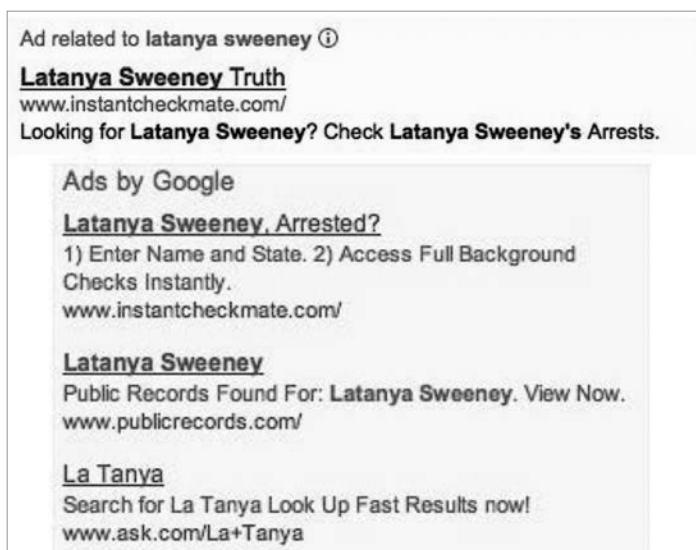


Рис. 3.1. Поиск Google, показывающий объявления с несуществующей записью об аресте профессора Латаньи Суни

Будучи программистом, она подошла к изучению систематически и просмотрела более 2000 имен. В итоге ей удалось выявить четкую закономерность: поиск исторических имен чернокожих людей приводил к выдаче объявлений, пред-

полагающих, что у человека есть судимость, в то время как для традиционно белых имен поиск выдавал более нейтральные варианты.

Это пример предвзятости, которая может существенно повлиять на жизнь человека: например, если работодатель решит поискать в Google информацию о соискателе, то может выясниться, что у него имеется криминальное прошлое, хотя это не так.

Почему это важно?

Одной из вполне естественных реакций на рассмотрение приведенной проблемы будет: «Ну и что? Какое это имеет отношение ко мне? Я специалист по данным, а не политик. Я не отношусь к числу руководителей нашей компании, решающих, что делать. Я просто стараюсь создать максимально успешно прогнозирующую модель».

Такое рассуждение очень логично, но попробуем убедить вас, что каждый, кто занимается обучением модели, обязательно должен учитывать то, как эта модель будет использоваться, а также подумать о том, как сделать это использование исключительно положительным. Вы реально можете на это повлиять, и если этого не сделать, то последствия могут быть очень печальными.

Одним из наиболее ужасающих примеров того, что происходит, когда инженеры сосредоточиваются на разработке технологий любой ценой, является история IBM и нацистской Германии. В 2001 году швейцарский судья постановил, что можно небезосновательно «сделать вывод, что IBM технически содействовала нацистам в совершении преступлений против человечества, облегчая ряд действий, включающих ведение бухгалтерского учета и классификации на устройствах IBM в самих концентрационных лагерях».

Если говорить точнее, то IBM снабжала нацистов продуктами для табулирования данных, которые были необходимы для отслеживания массового истребления евреев и других групп. Руководила этим процессом верхушка компании, реализуя маркетинг, соответствующий задачам Гитлера и его команды. Президент Томас Уотсон (Thomas Watson) лично утвердил выпуск специальных IBM-машин для алфавитного ввода в 1939 году, чтобы помочь организовать депортацию польских евреев. На рис. 3.2 показана встреча Адольфа Гитлера (крайний слева) с генеральным директором IBM Томом Уотсоном-старшим (второй слева), незадолго до того, как Гитлер наградил Уотсона медалью за «Особые заслуги перед Рейхом» в 1937 году.

Но это был не единственный инцидент — участие этой компании было обширным. IBM и ее дочерние компании проводили регулярное обучение персонала и обслуживание оборудования в самих концентрационных лагерях: распечатывали карточки, настраивали и чинили машины, так как они часто выходили из

строя. IBM организовала в своей системе перфокарт категоризацию, отражавшую способ, которым был убит человек, группу, к которой он приписывался, а также логистическую информацию, необходимую для отслеживания людей в огромной системе Холокоста (рис. 3.3). В системе IBM евреям был присвоен код 8: их было убито около 6 000 000. Для цыган использовался код 12 (их нацисты определяли, как «асоциальных», убив более 300 000 в лагере Zigeunerlager, иначе именовавшемся как «лагерь цыган»). Общий вид казни имел код 4, смерть в газовых камерах — 6.



Рис. 3.2. Генеральный директор IBM Том Уотсон-старший на встрече с Адольфом Гитлером

Nummer	Person	Kodierung	Geschlecht	Alter	Stellung	Merkmal	Merkmale	Familie	Körpergröße	Sitzhöhe	Größe	Rassezettel
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				
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				
2	2 2 2 2 2 2 2 2 2 2 2 2 2							2 2 2 2 2 2 2 2 2 2 2 2 2				
3	3 3 3 3 3 3 3 3 3 3 3 3 3							3 3 3 3 3 3 3 3 3 3 3 3 3				
4	4 4 4 4 4 4 4 4 4 4 4 4 4							4 4 4 4 4 4 4 4 4 4 4 4 4				
5	5 5 5 5 5 5 5 5 5 5 5 5 5							5 5 5 5 5 5 5 5 5 5 5 5 5				
6	6 6 6 6 6 6 6 6 6 6 6 6 6							6 6 6 6 6 6 6 6 6 6 6 6 6				
7	7 7 7 7 7 7 7 7 7 7 7 7 7							7 7 7 7 7 7 7 7 7 7 7 7 7				
8	8 8 8 8 8 8 8 8 8 8 8 8 8							8 8 8 8 8 8 8 8 8 8 8 8 8				
9	9 9 9 9 9 9 9 9 9 9 9 9 9							9 9 9 9 9 9 9 9 9 9 9 9 9				

RASSENAMT - 4

1

Рис. 3.3. Перфокарта, использованная IBM в концентрационных лагерях

Конечно же, руководители проекта, инженеры и техники, участвовавшие во всем этом, вели себя как обычно — заботились о семьях, ходили в церковь по воскресеньям, старались выполнять свою работу как можно лучше. Следовали приказам. Маркетологи также делали все что могли для достижения бизнес-целей. Как заметил Эдвин Блэк (Edwin Black), автор *IBM and the Holocaust* (IBM и Холокост) (Dialog Press): «Для слепого технократа средства были важнее результатов. Уничтожение еврейского народа стало даже менее важным, потому что воодушевляющие технические достижения IBM только усиливались фантастическими прибылями, которые можно было получить во времена, когда по всему миру выстраивались очереди за хлебом».

А вот здесь задумайтесь: как бы вы себя почувствовали, узнав, что являетесь частью системы, которая вредит обществу? Смогли бы вы это принять? Что бы вы могли сделать, чтобы этого не произошло? Здесь мы привели примеры крайних ситуаций, но на сегодня есть много негативных социальных последствий, связанных с искусственным интеллектом и машинным обучением, и некоторые мы опишем в текущей главе.

Причем это не только моральное бремя. Иногда инженеры непосредственно расплачиваются за свои действия. Например, первым, кого посадили в тюрьму в результате скандала вокруг Volkswagen, вскрывшего, что эта компания подделала тесты на выбросы систем дизельных двигателей, был не контролировавший проект менеджер и не руководитель из верхушки компании, а инженер Джеймс Лианг (James Liang), который просто делал то, что ему сказали.

Конечно, все не так плохо — если выяснится, что проект, в котором вы участвуете, в итоге послужит на благо даже всего одного человека, вы сможете искренне этому порадоваться.

Ладно, надеюсь, что мы убедили вас в необходимости думать об этических последствиях. Но что же вам нужно делать? Как специалисты, работающие с данными, мы склонны стремиться к улучшению наших моделей путем оптимизации тех или иных метрик. Но оптимизация метрики не обязательно приведет к улучшению результатов. А даже если это и поможет их улучшить, то скорее всего, окажется не единственной важной деталью. Рассмотрите последовательность шагов, происходящих между разработкой исследователем/практиком модели/алгоритма и точкой, в которой эта работа используется для принятия решения. Если мы рассчитываем получить желаемый результат, то необходимо рассматривать этот конвейер как целое.

Обычно между этими крайними точками есть достаточно длинная цепочка действий. Это особенно верно в случаях, когда вы являетесь исследователем, который может даже и не знать, будет ли результат исследования в итоге для чего-либо использован, либо если вы занимаетесь сбором данных, что происходит на еще более ранних стадиях конвейера. Но никто не сможет проинформировать всех участников цепочки о возможностях, ограничениях и деталях

вашей работы лучше, чем вы сами. Хоть и нет «волшебного средства», которое бы могло гарантировать, что ваша работа будет использована во благо, но вовлекаясь в процесс и задавая правильные вопросы, вы хотя бы сможете убедиться в том, что важные аспекты учтены.

Иногда будет правильным отреагировать на просьбу выполнить некую часть работы ответом «нет». Хотя зачастую в ответ мы слышим: «Если это сделаю не я, то сделает кто-то другой». Но подумайте вот о чем: если вас выбрали для этой работы, значит, вы оказались, по их мнению, наиболее подходящим для нее сотрудником. Если же вы за нее не возьметесь, это будет означать, что наиболее подходящий человек над проектом работать не будет. А если первые пять человек, к которым они обратятся, также откажут, то будет еще лучше.

Тесное взаимодействие процессов ML и дизайна продукта

По всей вероятности, вы делаете эту работу, так как рассчитываете на ее дальнейшее использование. Иначе вы просто тратите время. Итак, начнем с предположения, что ваша работа найдет себе применение. В процессе сбора данных и разработки модели вы принимаете множество решений. На каком уровне группировки хранить данные? Какую функцию потерь следует использовать? Какие контрольные и обучающие выборки использовать? Стоит ли сосредоточиться на простоте реализации, скорости вывода или точности модели? Как ваша модель будет обрабатывать не соответствующие области элементы данных? Допускает ли она точную настройку или через какое-то время ее следует повторно обучить с нуля?

Это вопросы, касающиеся не только алгоритма. Они также относятся к дизайну конечного продукта. Но кто бы в итоге ни разрабатывал систему с применением вашей модели — продакт-менеджеры, исполнительные директора, судьи, журналисты или доктора, — никто из них не будет в достаточной степени понимать принятые вами решения, не говоря уже об их изменении.

Например, два исследования показали, что разработанное Amazon ПО для распознавания лиц выдавало неточные (<https://oreil.ly/bL5D9>) и смешанные в расовом отношении результаты (<https://oreil.ly/cDYqz>). В Amazon заявили, что исследователям должны изменить предустановленные параметры, не объяснив при этом, как это могло бы повлиять на необъективные результаты. Более того, выяснилось, что компания Amazon также не инструктировала отделения полиции (<https://oreil.ly/lSOAj>), которые применяли их ПО. Вероятно, существовал большой разрыв между разработчиками этих алгоритмов и штатом Amazon, отвечавшим за составление руководств, переданных полиции.

Недостаток тесного взаимодействия привел к серьезным проблемам как для общества в целом, так и для полиции и Amazon в частности. Выяснилось, что

их система ошибочно сопоставила 28 членов конгресса с фотографиями людей в розыске. (При этом ошибочно сопоставленные конгрессмены большей частью оказались «цветными» людьми, как можно видеть на рис. 3.4.)



Рис. 3.4. Фотографии конгрессменов, сопоставленные программным обеспечением Amazon со снимками людей в розыске

Специалисты по данным должны быть частью междисциплинарной команды. А исследователи должны тесно работать с теми, кто будет впоследствии использовать их разработку. Более того, сами эксперты предметной области могут изучить достаточно, чтобы иметь возможность обучать и отлаживать некоторые модели самостоятельно, — надеемся, что некоторые из них читают эту книгу прямо сейчас!

Современное рабочее пространство — это очень специализированное место. Все стремятся к выполнению точно определенного вида работы. Это особенно актуально для крупных компаний, где бывает сложно знать все элементы пазла. Порой такие компании даже намеренно скрывают общие цели проекта, если знают, что сотрудникам они могут оказаться не по душе. Иногда это делается путем максимально возможного дробления проекта на части.

Мы не хотим сказать, что все это легко. Это трудно, очень трудно. Мы все должны стараться изо всех сил. И нам часто доводилось видеть, как люди, попав в более высокоуровневый контекст таких проектов, стараются разрабатывать междисциплинарные возможности и команды, становясь одними из наиболее ценных и высокооплачиваемых членов организаций. Такой вид работы очень ценится среди старших членов руководства, даже если иногда воспринимается средним руководящим звеном как неудобный.

Темы этики данных

Этика данных — это обширная область, и мы не можем рассмотреть ее всю. Выберем несколько тем, которые, на наш взгляд, особенно важны:

- необходимость в защите прав и ответственности;
- петли обратной связи;
- необъективность;
- дезинформация.

Рассмотрим их по порядку.

Защита прав и ответственность

В сложной системе часто бывает так, что никто не чувствует ответственности за результаты. Хоть это и можно понять, но к положительным результатам такое положение дел не ведет. В ранее приведенном примере с системой здравоохранения Арканзаса, в которой ошибка привела к тому, что люди с церебральным параличом утратили доступ к необходимой помощи, создатель алгоритма обвинил представителей правительства, а они обвинили тех, кто реализовал ПО. Профессор Нью-Йоркского университета Danah Boyd (Danah Boyd) (<https://oreil.ly/KK5Hf>) описала этот феномен так: «Бюрократия нередко использовалась для перекладывания или во избежание ответственности... Современные алгоритмические системы увеличивают бюрократию».

Еще одна причина, по которой защита прав необходима, заключается в том, что данные часто содержат ошибки. Очень важно наличие механизмов для аудита и исправления ошибок. В базе данных о подозреваемых членах банд, управляемой представителями правопорядка Калифорнии, обнаружилось множество ошибок, включая 42 младенца, которые были добавлены в базу, когда им было еще меньше 1 года (двадцать восемь из них были отмечены как «признающие принадлежность к банде»). В этом случае отсутствовала возможность исправления ошибок или удаления из базы людей после добавления. Другой пример — система кредитных отчетов США: обширное исследование кредитных отчетов, проведенное Федеральной торговой комиссией (FTC) в 2012 году, выяснило, что в файлах 26 % потребителей присутствовала по меньшей мере одна ошибка, а у 5 % содержались ошибки, которые могли привести к серьезным последствиям.

Да, процесс корректировки подобных ошибок невероятно медленный и непрозрачный. Когда репортер общественного радио Бобби Аллин (Bobby Allin) (<https://oreil.ly/BUD6h>) выяснил, что он ошибочно попал в список как имеющий судимость за огнестрельное оружие, ему потребовалось «больше дюжины звон-

ков, содействие клерка окружного суда и шесть недель для решения этой проблемы. Причем получилось это только тогда, когда я связался с пиар-отделом компании как журналист».

Как создатели продуктов машинного обучения мы не всегда считаем своей обязанностью понимать, как наши алгоритмы будут в итоге применены на практике. Но это необходимо.

Петли обратной связи

В главе 1 мы объясняли, как алгоритм может взаимодействовать со своей средой, создавая петлю обратной связи, которая приводит к прогнозам, усиливающим предпринимаемые в реальном мире действия, которые, в свою очередь, ведут к еще более выраженным прогнозам в том же направлении. В качестве примера еще раз рассмотрим рекомендательную систему YouTube. Пару лет назад представители Google рассказали, как они внедрили обучение с подкреплением (близко связанное с глубоким обучением, но в нем функция потерь представляет потенциальный результат, полученный спустя большой промежуток времени) для улучшения рекомендательной системы YouTube. Они описали, как использовали алгоритм, который давал рекомендации таким образом, чтобы время просмотра оптимизировалось.

Однако людей, как правило, привлекает противоречивый контент. Это означает, что видео о таких вещах, как теории заговоров, стали рекомендоваться системой все чаще и чаще. Более того, выяснилось, что люди, интересующиеся теориями заговоров, относятся к категории тех, кто смотрит много онлайн-видео. Это привело к их все большему увлечению YouTube. Растущее число любителей теорий заговоров, смотрящих видео на YouTube, привело к тому, что алгоритм все чаще рекомендовал теории заговоров и другой экстремистский контент, вследствие чего все больше экстремистов стали смотреть YouTube и все больше обычных зрителей YouTube стали разделять их экстремистские взгляды, на что алгоритм продолжил рекомендовать еще больше контента такого рода. Система выходила из-под контроля.

Причем данное явление не ограничивалось конкретно этим типом контента. В июне 2019-го *New York Times* опубликовала статью о рекомендательной системе YouTube под заголовком *On YouTube's Digital Playground, an Open Gate for Pedophiles* (<https://oreil.ly/81BEy>) («Цифровая площадка YouTube — открытые ворота для педофилов»). Эта статья начиналась со следующей пугающей истории:

Кристиан К. (Christiane C.) не увидела ничего плохого в том, что ее десятилетняя дочь с другом загрузили видео, на котором они играют в бассейне на заднем дворе... Несколько дней спустя у этого видео оказалась тысяча просмотров. Через какое-то время число выросло до 400 000... Кристиан впоследствии говорила: «Я еще раз посмотрела это видео, и меня взволно-

вало количество его просмотров». А причины для волнения действительно были. Как выяснили исследователи, автоматическая система рекомендаций начала показывать это видео пользователям, которые смотрели другие видео с частично одетыми подростками.

Само по себе каждое видео может быть абсолютно невинным, снятым, например, ребенком. Любые кадры с обнаженными людьми мимолетны и появляются случайно, но если сгруппировать их вместе, то станут отчетливо заметны общие признаки.

Алгоритм рекомендаций YouTube начал организовывать плейлисты для педофилов, выбирая невинные домашние видео, на которых присутствовали частично одетые подростки.

Никто в Google не планировал создавать систему, превращающую семейные видео в порно для педофилов. Что же произошло?

Отчасти проблема здесь в том, что в финансово важной системе метрики играют центральную роль. Как вы уже видели, когда у алгоритма есть метрика для оптимизации, он будет оптимизировать ее показатель всеми возможными путями. Это может вести ко всевозможным пограничным случаям, и люди, взаимодействующие с системой, будут искать, находить, а затем использовать эти пограничные случаи и петли обратной связи в целях собственной выгоды.

Есть признаки того, что именно это и произошло с рекомендательной системой YouTube в 2018 году. *The Guardian* опубликовала статью под названием *How an Ex-YouTube Insider Investigated Its Secret Algorithm* (<https://oreil.ly/yjnPT>) («Как бывший инсайдер YouTube расследовал его секретный алгоритм») о Гийоме Шасло (Guillaume Chaslot), бывшем инженере YouTube, создавшем сайт (<https://algotransparency.org/>) для отслеживания этих проблем. Шасло опубликовал график, показанный на рис. 3.5, следом за публикацией Роберта Мюллера (Robert Mueller) *Report on the Investigation Into Russian Interference in the 2016 Presidential Election* («Отчет о расследовании российского вмешательства в выборы президента 2016 года»).

Освещение отчета Мюллера на канале Russia Today резко вырывается вперед по числу каналов, рекомендовавших его к просмотру. Это предполагает, что государственное российское СМИ Russia Today успешно использовало алгоритм рекомендаций. К сожалению, недостаток прозрачности подобных систем усложняет раскрытие рассматриваемых нами проблем.

Один из обозревателей этой книги, Орельен Жерон (Aurélien Géron), вел команду, работающую над классификацией видеоматериалов в YouTube, с 2013 по 2016 год (задолго до описанных выше событий). Он пояснил, что дело не только в петлях обратной связи, где задействованы люди. Насчет примера с YouTube он сказал нам следующее:

Одним из важных сигналов для классификации основной темы видео является канал, на котором это видео размещено. Например, видео, загруженное на кулинарный канал, скорее всего, будет кулинарным. Но откуда мы знаем, какая у канала тема? Ну... частично просмотрев темы размещенных на нем видео. Вы видите петлю? Например, многие видео содержат описание, указывающее, какая при съемке использовалась камера. В результате некоторые из этих видео могут классифицироваться как видео о «фотографии». Если на канале есть такие ошибочно классифицированные материалы, то он может быть отмечен как канал о «фотографии», что приведет к тому, что будущие видео на нем также будут ошибочно классифицироваться как о «фотографии». Это даже может привести к неконтролируемым вирусным классификациям. Один из способов разорвать такую петлю обратной связи — это классифицировать видео с указанием канала и без него. Тогда при классификации каналов вы сможете использовать только классы, полученные без указания канала. Таким образом петля обратной связи будет разорвана.

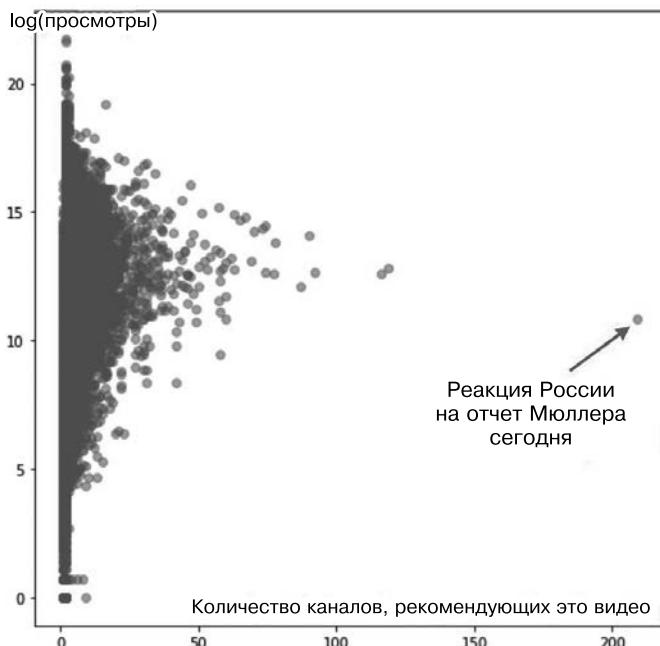


Рис. 3.5. Освещение отчета Мюллера

Есть положительные примеры людей и организаций, стремящихся решить такие проблемы. Эван Эстола (Evan Estola), ведущий инженер ML в Meetup, обсуждал пример (<https://oreil.ly/QfHzT>) мужчин как более заинтересованной группы в технических встречах, чем женщины. Учитывание пола могло привести к тому, что

алгоритмы Meetup стали бы рекомендовать меньше митапов женщинам, вследствие чего меньшее число женщин узнали бы об этих встречах и участвовали в них, что привело бы к рекомендации алгоритмом еще меньшего числа встреч женщинам, и т. д. Поэтому для исключения создания подобной петли обратной связи Эван с командой приняли этическое решение не учитывать в этой части модели пол участников проекта. Воодушевляет видение того, как компания не просто бездумно оптимизирует метрику, а также берет в расчет ее влияние. Как говорит Эван: «Вам необходимо решить, какую функцию в алгоритме не использовать... оптимальный алгоритм может оказаться не самым подходящим для запуска в работу».

В то время как в Meetup решили избежать подобных следствий, Facebook демонстрирует пример безудержного развития петли обратной связи. Подобно YouTube он имеет тенденцию радикализировать пользователей, заинтересованных в одной теории заговоров, знакомя их с другими. Рене Ди Реста (Renee DiResta), исследователь распространения дезинформации, пишет по этому поводу следующее (<https://oreil.ly/svgOt>):

Как только люди присоединяются к одной конспирологической группе, они алгоритмически направляются во множество других. Вступите в группу против вакцинации, и вам предложат группы борющихся с ГМО, отслеживающих химтрейлы, развивающих идею плоской Земли (да, такие существуют) и «лечащих рак природными способами». Вместо того чтобы вытащить пользователя из подобной кроличьей норы, алгоритм рекомендаций проталкивает их еще глубже.

Очень важно помнить, что такое поведение может возникнуть, поэтому нужно либо предвидеть петлю обратной связи, либо при обнаружении первых же признаков этого отклонения предпринимать положительные действия для ее разрыва. Помимо этого, нужно также учитывать *необъективность*, которая, как мы уже кратко говорили в прошлой главе, может взаимодействовать с петлями обратной связи, приводя к очень неблагоприятным последствиям.

Необъективность

Обсуждения смещения в интернете очень быстро превращаются в путаницу. Слово «смещение» может означать очень разные вещи. Статистики, видя, как специалисты по этике данных обсуждают смещение, думают, что те подразумевают статистическое значение этого термина, но это не так. При этом они также не имеют в виду смещения весов или параметров вашей модели.

В этом случае специалисты по этике говорят о принципе смещения из области социологии. Харини Суреш (Harini Suresh) и Джон Гуттаг (John Guttag) из MIT в своей работе *A Framework for Understanding Unintended Consequences of Machine*

Learning (<https://oreil.ly/aF33V>) («Схема для понимания непредвиденных последствий машинного обучения») описывают шесть типов смещения в области ML, которые визуально отражены на рис. 3.6.

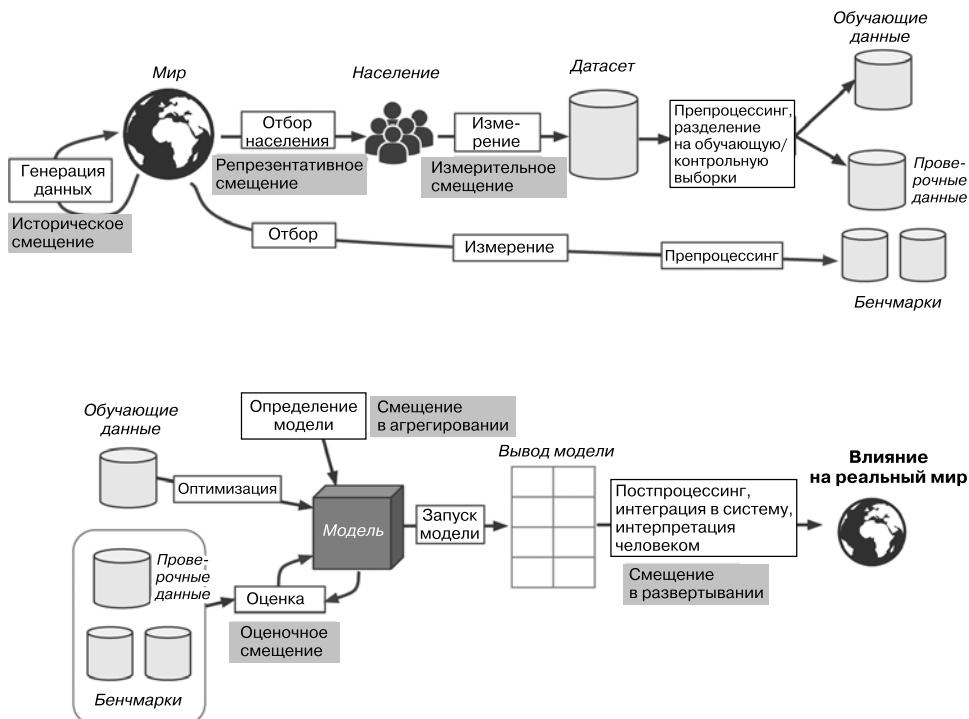


Рис. 3.6. Смещение в машинном обучении может быть вызвано несколькими источниками (схема любезно предоставлена Харини Сурешем и Джоном В. Гуттагом)

Рассмотрим четыре наиболее актуальных вида смещения.

Историческое смещение

Историческое смещение обусловлено тем фактом, что люди предвзяты, процессы основаны на предубеждениях, и общество тоже. Суреш и Гуттаг говорят так: «Историческое смещение — это фундаментальная структурная проблема первого этапа процесса генерации данных, которая может существовать даже при наличии идеальных выборок и признаков».

Вот несколько примеров исторических *расовых предубеждений* в США, взятых из статьи *Racial Bias, Even When We Have Good Intentions, New York Times* (<https://oreil.ly/cBQop>) («Даже благие намерения не устраняют расовые предубеждения»), написанной Сендхилом Муллайнатаном из Чикагского университета.

- Когда докторам показывали идентичные карточки пациентов, катетеризацию сердца (лечебная процедура) они реже рекомендовали темнокожим пациентам.
- При торгах за подержанный автомобиль чернокожим предлагали начальную цену на \$700 выше и шли на гораздо меньшие уступки.
- При обращении по объявлениям сайта Craigslist об аренде жилья темнокожим пользователям отвечали реже, чем белым пользователям.
- Суд присяжных, состоящий исключительно из белых людей, на 16 % чаще выносил обвинительный приговор темнокожим подзащитным, чем белым. Однако если в составе присяжных присутствовал хотя бы один темнокожий, то соотношение обвинительных приговоров выравнивалось.

COMPAS, широко используемый в США для вынесения приговоров и принятия решений об освобождении под залог, является примером важного алгоритма, который при тестировании некоммерческой организацией ProPublica (<https://oreil.ly/1XocO>) отчетливо проявил расовую предвзятость на практике (рис. 3.7).

	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%

Рис. 3.7. Результаты алгоритма COMPAS

Любой датасет, подразумевающий участие людей, может иметь подобный вид смещения: медицинские данные, данные продаж, политические данные, данные о жилье и т. д. Всепроникающая предвзятость, которая кроется в самих людях, ведет к всепроникающему смещению в датасетах. Расовые предубеждения обнаруживаются даже в моделях компьютерного зрения, как видно из примера на рис. 3.8, где показаны автоматически разбитые по категориям фотографии, выложенные в Twitter пользователем Google Photos.

Да, вы все правильно поняли: система Google Photos классифицировала фотографию чернокожей девушки с другом как «горилл!» Этот алгоритмический промах привлек много внимания в медиасреде. «Мы потрясены и искренне сожалеем, что такое произошло, — сообщил представитель компании. — В отношении автоматической разметки изображений еще многое предстоит сделать, и мы стараемся найти способы избежания подобных ошибок в будущем».



Рис. 3.8. Одна из присвоенных меток крайне ошибочна...

К сожалению, исправить проблемы в системах ML при поступлении в них изначально проблемных данных достаточно сложно. Как сообщалось в *The Guardian*, первая попытка Google исправить ситуацию не внушала особого оптимизма (рис. 3.9).

Google's solution to accidental algorithmic racism: ban gorillas

Google's 'immediate action' over AI labelling of black people as gorillas was simply to block the word, along with chimpanzee and monkey, reports suggest



▲ A silverback high mountain gorilla, which you'll no longer be able to label satisfactorily on Google Photos.
Photograph: Thomas Mukoya/Reuters

Рис. 3.9. Первый ответ Google на описанную проблему

Такой вид проблем присущ не только Google. Исследователи из MIT изучили наиболее популярные API компьютерного зрения, проверив их точность. Но при этом они не просто вычислили один показатель точности, а рассмотрели ее в отношении сразу четырех групп, как показано на рис. 3.10.

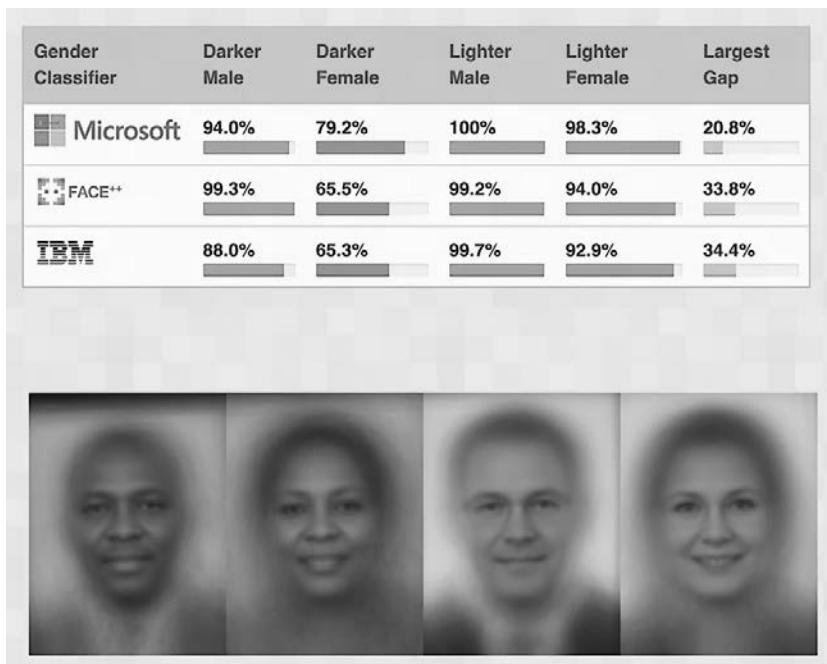


Рис. 3.10. Показатель ошибок в отношении пола и расы в различных системах распознавания лиц

К примеру, в системе IBM показатель ошибок при определении темнокожих женщин составил 34,7 %, в отличие от показателя 0,3 % для светлокожих мужчин, то есть допущено в 100 раз больше ошибок! Некоторые люди некорректно отреагировали на эти эксперименты, заявив, что разница возникла просто из-за того, что компьютерам сложнее распознавать темную кожу. А дальше произошло следующее: после негативной общественной реакции на эти результаты все рассматриваемые компании существенно улучшили свои модели в отношении распознавания темнокожих, вследствие чего, спустя всего один год, результаты для светлокожих и темнокожих стали практически одинаковыми. Возникавшие изначально ошибки показали, что разработчики либо плохо проработали датасеты, не добавив в них достаточное количество темнокожих лиц, либо просто плохо протестировали свои продукты в отношении них.

Джой Буоламвини (Joy Buolamwini), один из исследователей MIT, предупредил: «Мы вступили в эру излишне уверенной, но недостаточно подготовленной ав-

томатизации. Если мы не справимся с разработкой этического и справедливого искусственного интеллекта, то в будущем мы рискуем утратить достижения в области гражданских прав и гендерного равенства».

Частично проблема заключается в систематическом нарушении баланса при создании популярных датасетов, используемых для обучения моделей. В аннотации к работе Шреи Шанкарь и др. (Shreya Shankar et al.) *No Classification Without Representation: Assessing Geodiversity Issues in Open Data Sets for the Developing World* (<https://oreil.ly/VqtOA>) («Нет классификации без репрезентативности: анализ проблем георазнообразия в открытых наборах данных для развивающихся стран») сказано следующее: «Мы анализируем два больших общедоступных датасета изображений для анализа георазнообразия и обнаруживаем, что эти датасеты демонстрируют наблюдаемое смещение в сторону американоцентрической и евроцентрической репрезентативности. Далее мы анализируем классификаторы, обученные на этих датасетах, чтобы понять оказываемое ими влияние и найти существенные различия в относительной производительности для изображений из разных регионов». На рис. 3.11 показана одна из опубликованных в этой работе диаграмм, которая демонстрирует географический состав изображений из двух наиболее важных датасетов для обучения моделей (на момент написания книги информация остается актуальной).

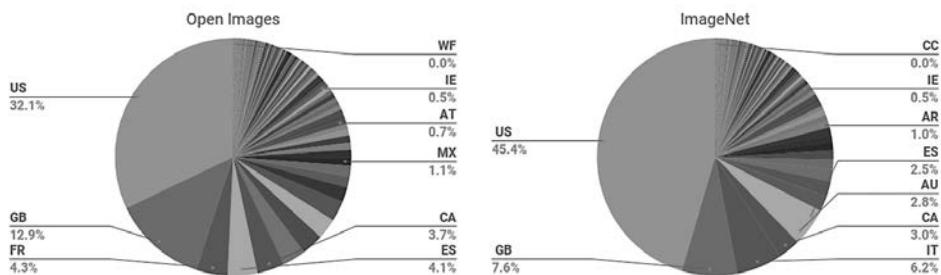


Рис. 3.11. Происхождение изображений в популярных обучающих датасетах

Здесь видно, что подавляющее большинство изображений предоставляют США и другие западные страны. Это ведет к тому, что обученные на ImageNet модели показывают себя хуже в отношении сцен из других стран и культур. Например, исследование показало, что такие модели хуже идентифицируют предметы домашнего быта (такие, как мыло, специи, диваны или кровати) из стран с низкими доходами населения. На рис. 3.12 показано изображение из работы *Does Object Recognition Work for Everyone?* (<https://oreil.ly/BkFjL>) («Для всех ли работает распознавание объектов?») Терранса Де Брис и др. (Terrance DeVries et al.) из отдела по исследованию ИИ в Facebook, отражающее эту идею.



Ground truth: Soap

Nepal, 288 \$/month

Azure: food, cheese, bread, cake, sandwich
 Clarifai: food, wood, cooking, delicious, healthy
 Google: food, dish, cuisine, comfort food, spam
 Amazon: food, confectionary, sweets, burger
 Watson: food, food product, turmeric, seasoning
 Tencent: food, dish, matter, fast food, nutrient



Ground truth: Soap

UK, 1890 \$/month

Azure: toilet, design, art, sink
 Clarifai: people, faucet, healthcare, lavatory, wash closet
 Google: product, liquid, water, fluid, bathroom accessory
 Amazon: sink, indoors, bottle, sink faucet
 Watson: gas tank, storage tank, toiletry, dispenser, soap dispenser
 Tencent: lotion, toiletry, soap dispenser, dispenser, after shave



Ground truth: Spices

Phillipines, 262 \$/month

Azure: bottle, beer, counter, drink, open
 Clarifai: container, food, bottle, drink, stock
 Google: product, yellow, drink, bottle, plastic bottle
 Amazon: beverage, beer, alcohol, drink, bottle
 Watson: food, larder food supply, pantry, condiment, food seasoning
 Tencent: condiment, sauce, flavorer, catsup, hot sauce



Ground truth: Spices

USA, 4559 \$/month

Azure: bottle, wall, counter, food
 Clarifai: container, food, can, medicine, stock
 Google: seasoning, seasoned salt, ingredient, spice, spice rack
 Amazon: shelf, tin, pantry, furniture, aluminium
 Watson: tin, food, pantry, paint, can
 Tencent: spice rack, chili sauce, condiment, canned food, rack

Рис. 3.12. Распознавание объектов в действии

В этом примере мы видим, что распознавание образца мыла низкого качества далеко от точного, и каждый коммерческий сервис по распознаванию изображений с наибольшей вероятностью прогнозирует для него определение «food», то есть «еда».

Вскоре мы также обсудим, что вдобавок к этому большинство исследователей ИИ и разработчиков являются белыми мужчинами. В большую часть рассмотренных нами проектов на этапе тестирования пользователями привлекаются друзья и родственники непосредственной группы разработки продукта. Стоит ли в таком случае удивляться, что мы наблюдаем проблемы, подобные описанным выше.

Похожее историческое смещение также обнаружено в текстах, использованных в качестве данных для моделей обработки естественного языка, что по-разному оказывается на последующих задачах машинного обучения. Например, вплоть до прошлого года Google Translate демонстрировал систематическое смещение в переводе турецкого гендерно-нейтрального местоимения «о» на английский: применительно к профессиям, которые чаще ассоциируются с мужчинами, сервис использовал «он», а в случае с профессиями, чаще ассоциирующимися с женщинами, использовалось местоимение «она» (рис. 3.13).

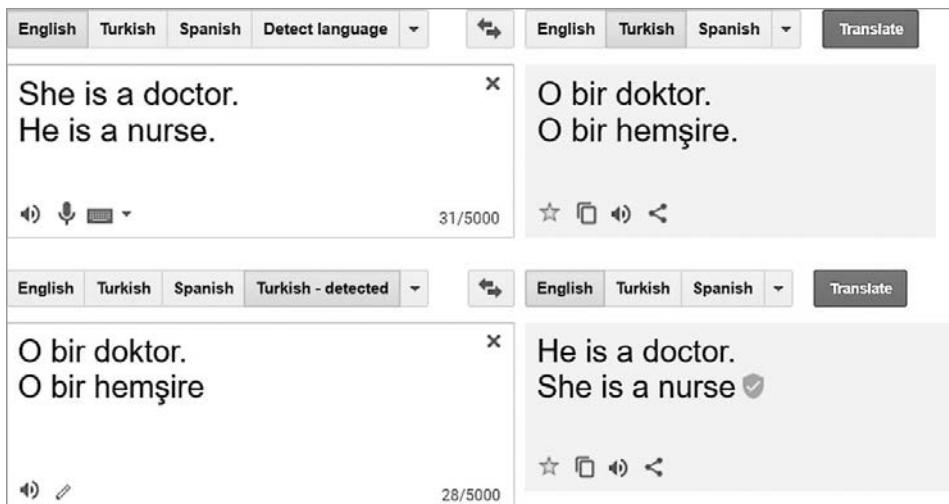


Рис. 3.13. Гендерное смещение в текстовых датасетах

Подобный вид смещения мы также наблюдаем в онлайн-объявлениях. Например, исследование Мухаммеда Али (Muhammad Ali) и др., проведенное в 2019 году, выявило, что даже когда размещающий объявление человек

не предусматривает в нем дискриминации, Facebook будет показывать объявления совершенно различным аудиториям, основываясь на расе и половой принадлежности. Реклама жилья с одинаковым текстом, но сопровождаемая изображением белой или черной семьи, показывалась отличающейся по расовому признаку аудитории.

Смещение измерений

В статье *Does Machine Learning Automate Moral Hazard and Error* (<https://oreil.ly/79Qtn>) («Автоматизирует ли машинное обучение моральную угрозу и ошибки?»), опубликованной в журнале *American Economic Review*, Сенхил Муллайнатан (Sendhil Mullainathan) и Зиад Обермейер (Ziad Obermeyer) рассматривают модель, пытающуюся ответить на вопрос: какие факторы окажутся наиболее актуальными для прогноза инсульта при использовании исторических данных электронной медицинской карты (EHR)? Ниже в порядке убывания значимости перечислены наиболее актуальные с точки зрения модели предикторы:

- инсульт в анамнезе;
- сердечно-сосудистые заболевания;
- случайная травма;
- доброкачественная опухоль молочной железы;
- колоноскопия;
- синусит.

При этом реальное отношение к инсульту имеют только первые два пункта! Если учесть всю рассмотренную нами недавно информацию, будет нетрудно догадаться почему. В действительности мы не диагностировали *инфаркт*, возникающий, когда некоторая область мозга недополучает кислород из-за нарушения кровоснабжения. Мы проанализировали симптомы тех, кто обратился к врачу, прошел соответствующие тесты и получил диагноз «*инфаркт*». На самом деле наличие инсульта — это не единственный показатель, имеющий отношение к приведенному списку, потому что список также относится к тем, кто регулярно посещает врача (что определяется наличием доступа к медицинским услугам, способностью их оплатить, отсутствием расовой или гендерной дискриминации, а также другими факторами). Если вы собираетесь обратиться к врачу при получении *случайной травмы*, то, скорее всего, также обратитесь к нему при инсульте.

Это пример погрешности измерений, которая появляется, когда наши модели совершают ошибки из-за того, что мы измеряем не те вещи, либо измеряем их неверным способом, либо неправильно внедряем полученные измерения в модель.

Смещение агрегирования

Смещение агрегирования происходит, когда модели не агрегируют данные тем образом, который включал бы все подходящие факторы, или когда модель не включает необходимые условия взаимодействия, нелинейности и т. д. Этому особенно подвержены медицинские учреждения. Например, способ лечения диабета зачастую основан на простой одномерной статистике и исследованиях с участием небольших групп разнотипных людей. Результаты же обычно анализируются без учета этнической или гендерной принадлежности. При этом выяснилось, что осложнения, возникающие у пациентов с диабетом, зависят от их этнической группы (<https://oreil.ly/gNS39>), и уровни HbA1c (широко используемые для диагностирования и отслеживания диабета) сложным образом отличаются в зависимости от пола и этнической принадлежности человека (<https://oreil.ly/nR4fx>). Это может приводить к ошибочным диагнозам или неправильному лечению, потому что медицинские решения основываются на модели, которая не включает эти важные переменные и взаимодействия.

Смещение представления

В аннотации к работе *Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting* (<https://oreil.ly/0iowq>) («Смещения в биографических данных: изучение проблемы смещения семантического представления в условиях особой важности»), написанной Марией Де-Артэага (Maria De-Arteaga) и др., отмечается, что в профессиях присутствует гендерное смещение (например, женщины чаще становятся медсестрами, а мужчины — священниками) и что «различия в истинно положительных показателях между полами коррелируют с имеющимся гендерным дисбалансом в профессиях, что может усугубить этот дисбаланс».

Другими словами, исследователи заметили, что модели, прогнозирующие профессии людей, не только *отражают* фактический гендерный дисбаланс среди основного населения, но и *усиливают* его. Такой тип *смещения представления* достаточно распространен, особенно в простых моделях. Если модель обнаруживает в основе отчетливые, ясно видимые связи, то она будет часто предполагать, что эти связи сохраняются постоянно. Как показано на рис. 3.14, взятом из приведенной работы, для профессий, в которых процент женщин выше, модель имеет тенденцию завышать распространенность этой профессии среди женщин.

К примеру, в обучающем датасете 14,6 % хирургов были женщинами, хотя в прогнозе модели только 11,6 % истинно-положительных показателей оказались женщинами. Таким образом данная модель усиливает смещение, существующее в обучающей выборке.

Теперь, когда мы увидели, что подобные смещения существуют, что мы можем предпринять для их нейтрализации?

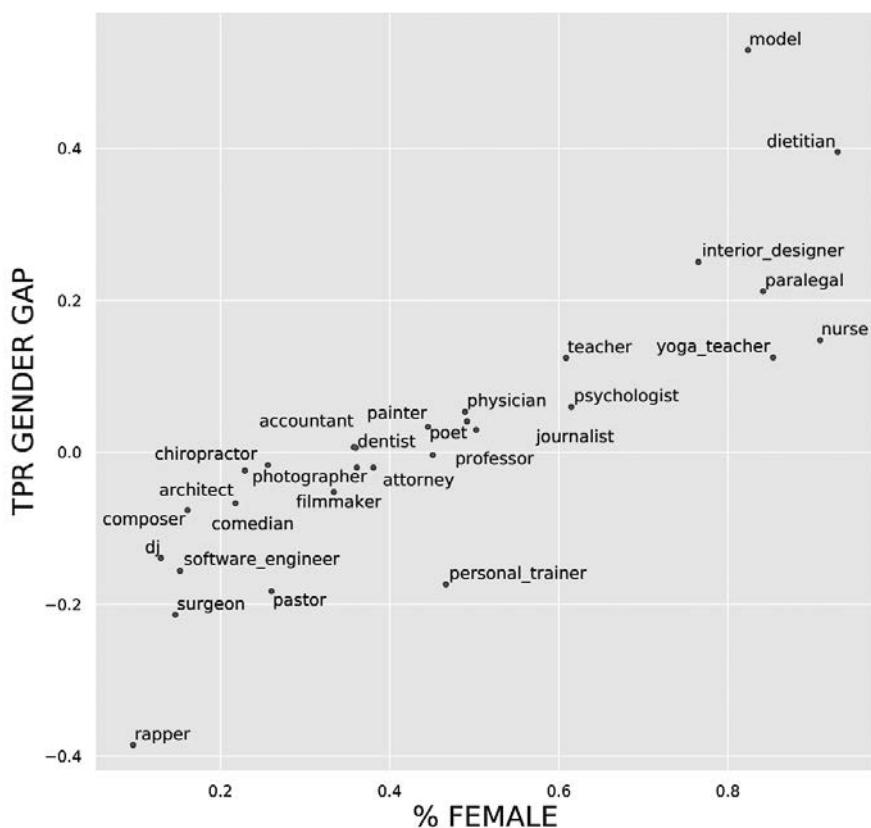


Рис. 3.14. Ошибка модели в прогнозировании профессии, графически отображенная на основе процента занятых в этой профессии женщин

Устранение различных видов смещения

Для смягчения разных видов смещения требуются разные подходы. В то время как смещение репрезентации можно решить сбором более разнообразного датасета, это не поможет в случае с историческим смещением или смещением в измерениях. Все датасеты имеют определенное смещение, по-другому просто не бывает. Многие исследователи в этой области сошлись во мнениях относительно введения лучшего документирования решений, контекста и особенностей того, как и почему был создан конкретный датасет, в каких сценариях его следует использовать и какие он имеет ограничения. Таким образом, те, кто будет использовать определенный датасет, не будут застигнуты врасплох его смещениями и ограничениями.

Мы часто слышим вопрос: «Раз люди необъективны, имеет ли значение предвзятость алгоритма?» Этот вопрос звучит настолько часто, что за ним наверняка

кроется определенный смысл для тех, кто его задает. Но для нас он вовсе не кажется логичным! Независимо от того, звучит ли он логично, важно понять, что алгоритмы (особенно алгоритмы машинного обучения) и люди отличаются. Рассмотрите следующие пункты, относящиеся к алгоритмам ML.

Машинное обучение может создавать петли обратной связи

Из-за петли обратной связи небольшое смещение может в течение короткого промежутка времени увеличиться многократно.

Машинное обучение может усиливать смещение

Человеческая предвзятость может вести к еще большему увеличению смещения в ML.

Алгоритмы и люди используются по-разному

Принимающие решения люди и принимающие решения алгоритмы не допускают взаимной замены на практике. Примеры этого приведены в виде списка ниже.

В технологии — сила

А с ней приходит и ответственность.

Как показал пример с системой здравоохранения Арканзаса, машинное обучение часто реализуется на практике не потому, что ведет к улучшенным результатам, а потому что оно дешевле и более эффективно. Кэти О'Нил (Cathy O'Neill) в своей книге *Weapons of Math Destruction* (Crown) («Орудия математического уничтожения») описала шаблон, в котором привилегированные личности обслуживаются людьми, в то время как бедные обслуживаются алгоритмами. Это всего один из многих вариантов, в которых алгоритмы используются не так, как принимающие решения люди. К другим примерам можно отнести такие:

- Люди склонны предполагать, что алгоритмы объективны или не допускают ошибок (даже если в них предусмотрена возможность перенастройки людьми).
- Алгоритмы внедряются без возможности коррекции.
- Алгоритмы часто используются масштабно.
- Алгоритмические системы недороги.

Даже при отсутствии смещения алгоритмы (в особенности глубокое обучение, поскольку является эффективным и масштабируемым) могут вести к социальным проблемам, например, при использовании для *дезинформации*.

Дезинформация

Понятие *дезинформации* корнями уходит на сотни или даже тысячи лет в прошлое. Оно необязательно относится к введению кого-либо в заблуждение, но чаще используется для внесения дисгармонии и неуверенности, а также для того, чтобы люди оставили попытки поиска истины. Столкнувшись с двумя противодействующими сторонами, человек может усомниться в том, кому или чему можно верить.

Некоторые думают, что дезинформация в прямом смысле означает неверную информацию, или *фейковые* новости, но в реальности дезинформация может зачастую содержать доли истины или полуправду, выдернутую из контекста. Ладислав Биттман (Ladislav Bittman) был офицером разведки в СССР, который перешел на сторону США в 1970-х и 1980-х годах написал несколько книг о роли дезинформации в советских пропагандистских операциях. В книге *The KGB and Soviet Disinformation* («КГБ и советская дезинформация») (издательство Pergamon) он писал: «Большинство кампаний представляют осторожно собранную смесь фактов, полуправд, преувеличений и продуманной лжи».

Не так давно этот факт раскрылся в отношении самих США, когда ФБР подробно описало массивную кампанию по дезинформации на выборах президента 2016 года, связанную с Россией. Понимание дезинформации, использованной в этой кампании, очень поучительно. Например, ФБР выяснило, что в ходе этой российской кампании часто организовывались две искусственные акции «общественного» протеста, по одной для каждой стороны ситуации, причем обе они выступали в одно и то же время. Газета *Houston Chronicle* (<https://oreil.ly/VyCkL>) сообщила об одном из таких странных событий следующее (рис. 3.15):

Группа, назвавшая себя «Сердце Техаса», организовала через соцсети протест, как они заявили, против «исламизации» Техаса. На одной стороне Тревис-стрит я обнаружил 10 протестующих. На другой стороне наблюдалось около 50 контрпротестующих. Но я не смог обнаружить самих организаторов митинга. Никакого «Сердца Техаса». Мне показалось это странным, и я решил упомянуть об этом в статье: какая группа не показывается на собственном событии? Теперь я понимаю, почему так произошло. По-видимому, в тот момент организаторы находились в Санкт-Петербурге, в России. «Сердце Техаса» — это одна из групп интернет-троллей, упомянутая в недавнем обвинительном заключении спецпрокурора Роберта Мюллера (Robert Mueller) в отношении россиян, пытающихся вмешаться в президентские выборы в США.

Дезинформация нередко включает организацию кампаний лжеповедения. Например, с помощью фейковых аккаунтов может искусственно формироваться мнение, что определенной точки зрения придерживаются многие люди. Несмотря на то что большинство из нас привыкли считать себя независимыми

в суждениях, в реальности мы формируемся под влиянием людей из близкого окружения. Онлайн-дискуссии могут влиять на нашу точку зрения или изменять спектр допустимых, на наш взгляд, точек зрения. Люди — это социальные животные, что и определяет нашу подверженность влиянию со стороны окружающих. Все чаще радикализация происходит в онлайне, а значит, влияние исходит от людей в виртуальном пространстве форумов и соцсетей.

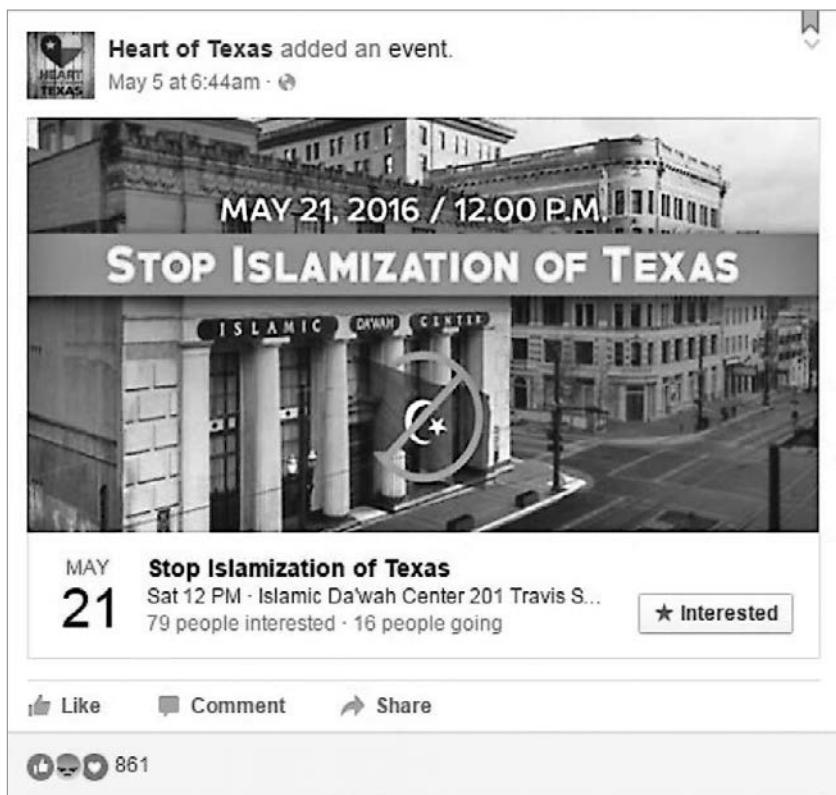


Рис. 3.15. Событие, организованное группой «Сердце Техаса»

Дезинформация посредством автоматически генерируемого текста является особенно выраженной проблемой, так как глубокое обучение существенно увеличило ее возможности. Мы обсудим эту тему подробно, когда приступим к созданию языковых моделей в главе 10.

Один из способов решения проблемы дезинформации состоит в разработке некой легко применяемой формы цифровой подписи наряду с созданием норм, согласно которым нам следует доверять только проверенному на эту подпись контенту. Глава Института Аллена по ИИ Орен Эциони (Oren Etzioni) изложил

такое предложение в статье под названием *How Will We Prevent AI-Based Forgery?* (<https://oreil.ly/8z7wm>) («Как нам предотвратить фальсификацию на основе ИИ?»): «Искусственный интеллект уже способен выдавать высококачественную подделку недорого и автоматизированно, что потенциально ведет к чудовищным последствиям для демократии, безопасности и общества. Угроза ИИ-подделок означает, что нам необходимо действовать в направлении внедрения обязательных цифровых подписей как средств аутентификации цифрового контента».

Хотя мы и не в силах обсудить все этические сложности, формируемые глубоким обучением и алгоритмами в целом, надеемся, что это краткое введение будет полезной стартовой точкой. Теперь пришло время перейти к вопросам выявления и решения этических проблем.

Выявление этических проблем и их решение

Ошибки неизбежны. Поэтому их обнаружение и устранение должно стать частью дизайна любой системы, использующей ML (и многие другие системы тоже). Проблемы, возникающие в рамках этики данных, зачастую сложны и междисциплинарны, но работать над их исправлением необходимо.

Итак, что можно сделать? Это обширная тема, но мы выделим несколько шагов по решению подобного рода проблем:

- анализируйте проект, над которым работаете;
- реализуйте в компании процессы для обнаружения и устраниния этических рисков;
- поддерживайте правильную политику;
- создавайте максимально разносторонние команды сотрудников.

Теперь пройдем по каждому шагу, начиная с анализа проекта.

Анализ проекта

Легко упустить важные моменты, когда рассматриваешь этические последствия своей работы. В такой ситуации очень помогает постановка правильных вопросов. Рейчел Томас (Rachel Thomas) рекомендует рассмотреть такие вопросы, которые можно задавать при разработке проекта.

- Должны ли мы вообще это делать?
- Какое смещение присутствует в данных?
- Можно ли проверить код и данные?

- Каковы показатели ошибок для различных подгрупп?
- Какова точность основанной на правилах простой альтернативы?
- Какие существуют процессы для обработки возражений или ошибок?
- Насколько разносторонняя команда это создала?

Эти вопросы помогут вам выявить важнейшие проблемы и возможные альтернативы, которые проще понимать и контролировать. Помимо постановки правильных вопросов, также важно учесть методы и процессы, которые нужно реализовать.

Здесь необходимо подумать о том, какие данные вы собираете и храните. Нередко бывает так, что данные в итоге используются не для того, для чего предназначались изначально. Например, IBM начала продавать оборудование и услуги задолго до Холокоста, включая помощь в ведении переписи населения Германии, организованной Адольфом Гитлером в 1933 году, которая эффективно выявила, что в стране проживает намного больше евреев, чем до этого считалось. Аналогичным образом данные переписи населения в США во время Второй мировой войны использовались для облавы на американцев японского происхождения (которые были гражданами США) с целью их интернирования. Важно распознать, как собранные данные и изображения могут быть использованы в качестве оружия в дальнейшем. Профессор Колумбийского университета Тим Ву (Tim Wu) (<https://oreil.ly/6L0QM>) написал: «Вы должны предположить, что любые личные данные, которые хранит Facebook или Android, являются данными, которые постараются заполучить правительства разных стран или украсть воры».

Какие процессы нужно реализовать

Центр им. Марккулы выпустил доклад *An Ethical Toolkit for Engineering/Design Practice* (<https://oreil.ly/vDGGC>) («Набор этических инструментов для разработки и проектирования»), помимо прочего описывающий конкретные практики для применения в вашей компании, включая плановую профилактику появления этических рисков (аналогично тестированию кибербезопасности на угрозу проникновения), расширение этического кругозора добавлением точек зрения различных участников и учет потенциальных действий злоумышленников (возможное злоупотребление, кражу, ошибочное толкование, взлом, уничтожение или использование созданного вами продукта в качестве оружия).

Даже если у вас недостаточно разносторонняя команда, вы все равно можете попробовать самостоятельно представить точки зрения более широкой группы, рассматривая вопросы, подобные приведенным Центром Марккулы.

- Чьи интересы, желания, навыки, опыты и ценности мы не обсудили фактически?

- На кого в конечном счете окажет влияние наш продукт? Что мы сделали для защиты их интересов? Откуда нам эти интересы известны — спрашивали ли мы о них?
- На каких людей или какие группы будет оказано косвенное, но существенное влияние?
- Кто может использовать этот продукт так, как мы не ожидали, и в каких целях это может быть сделано?

Этические призмы

Еще один полезный ресурс, предоставленный центром Марккулы, — это работа *Conceptual Frameworks in Technology and Engineering Practice* (<https://www.scu.edu/ethics-in-technology-practice/ethical-lenses/>) («Концептуальные основы технологии и практической инженерии»). В ней рассматривается, как с помощью основополагающих этических призм можно выявить конкретные проблемы, а также излагаются следующие аспекты и ключевые вопросы.

Права

В каком варианте наилучшим образом учитываются права всех участников?

Справедливость

В каком варианте люди рассматриваются равноправно?

Утилитарность

Какой вариант создаст наибольшие блага при минимальном вреде?

Общее благо

Какой вариант наилучшим образом послужит сообществу в целом, а не отдельным его членам?

Добродетель

Какой вариант ведет меня к желаемой мной модели человеческого поведения?

Рекомендации центра подразумевают углубленное погружение в каждую из этих перспектив, включая рассмотрение проекта через призму его *последствий*.

- На кого проект окажет непосредственное влияние, а на кого — косвенное?
- Принесут ли эффекты в совокупности больше блага, чем вреда, и какого именно рода будут эти блага или вред?

- Подумали ли мы обо всех смежных типах вреда/пользы (психологическом, экологическом, политическом, моральном, когнитивном, эмоциональном, институциональном и культурном)?
- Как этот проект может повлиять на будущие поколения?
- Не подвергаются ли наименее влиятельные слои населения наибольшим рискам получения вреда от этого проекта? Не получат ли наиболее влиятельные максимальную выгоду?
- Адекватно ли мы рассмотрели «двуцелевое использование» и непредвиденные последующие эффекты?

Альтернативой этому будет *деонтологическая* перспектива, фокусирующаяся на базовых понятиях добра и зла.

- Какие права других людей и обязанности по отношению к ним мы должны уважать?
- Как этот проект может повлиять на чувство достоинства и независимости каждого участника?
- Какие соображения доверия и справедливости соответствуют этому дизайну/проекту?
- Подразумевает ли этот проект противоречивые моральные обязанности перед другими людьми или конфликт прав его участников? Как в этом случае можно расставить приоритеты?

Один из лучших способов сформировать полноценные продуманные ответы на поставленные вопросы — опрашивать *разнообразных в культурно-личностном отношении* людей.

Сила разнообразия

Согласно проведенному Element AI анализу (<https://oreil.ly/sO09p>) число женщин — исследователей ИИ составляет всего 12 %. Аналогичным образом выглядит статистика, когда речь заходит о расе и возрасте. Если у всех участников команды будет схожий жизненный опыт, то такая команда будет склонна к однобокому рассмотрению этических рисков. В *Harvard Business Review* (HBR) были опубликованы несколько исследований, показывающих множество преимуществ неоднородных команд, включая следующие.

- *How Diversity Can Drive Innovation* (<https://oreil.ly/WRFSm>) («Как неоднородность может способствовать инновации»).
- *Teams Solve Problems Faster When They're More Cognitively Diverse* (<https://oreil.ly/vky5b>) («Когнитивно разнообразные команды быстрее справляются с задачами»).

- *Why Diverse Teams Are Smarter* (<https://oreil.ly/SFVBF>) («Почему неоднородные команды умнее»).
- *Defend Your Research: What Makes a Team Smarter? More Women* (<https://oreil.ly/A1A5n>) («Отстаивайте свое исследование: что делает команды умнее? Больше женщин»).

Разнообразие может вести к более раннему выявлению проблем и рассмотрению более широкого спектра их решений. Например, Трейси Чоу (Tracy Chou) была начинающим инженером в Quora. Впоследствии она описала свой опыт (<https://oreil.ly/n7WSn>), рассказав, как выступала за добавление функции, которая позволила бы блокировать троллей и других вредителей. Чоу вспоминает: «Мне не терпелось проработать эту функцию, так как я лично ощущала враждебные высказывания и оскорблении на этом сайте (вполне вероятно, что это происходило по гендерным причинам)... Но если бы я не имела такого личного взгляда, то, возможно, команда Quora не внесла бы уже на ранних этапах реализации проекта в список приоритетов разработку кнопки блокирования». Преследования часто вынуждают людей из социально отчужденных групп покидать онлайн-платформы, поэтому такая функциональность была важна для поддержки здоровой атмосферы сообщества Quora.

Очень важно понимать тот факт, что женщины уходят из технологической индустрии вдвое чаще мужчин. Согласно данным Harvard Business Review (<https://oreil.ly/ZIC7t>) из индустрии уходит 41 % женщин, в то время как для мужчин этот показатель составляет всего 17 %. Анализ более 200 книг, официальных документов и статей выявил, что причиной ухода является то, что «с ними обращаются несправедливо, меньше платят, и вероятность продвижения по службе ниже, чем у их коллег-мужчин».

Исследования подтвердили наличие факторов, затрудняющих процесс продвижения женщин по карьерной лестнице. При оценке их работы женщины получают более расплывчатую обратную связь и повышенную личностную критику, в то время как мужчинам в таких случаях даются конструктивные рекомендации, привязанные к результатам работы (что более полезно). Женщины часто не допускаются до творческой и новаторской деятельности, а также не получают более перспективных заданий «для роста», которые могли бы помочь им продвинуться по службе. Одно из исследований обнаружило, что голоса мужчин воспринимаются более убедительными, основанными на фактах, чем женские, даже при озвучивании одних и тех же фраз.

Статистика подтверждает, что наличие наставника помогает развивать свою карьеру мужчинам, но не женщинам. Причина в том, что женщинам под видом наставничества чаще дают советы о том, как нужно меняться и получать больше знаний самостоятельно. При получении же наставничества мужчинами это означает публичное признание их авторитета. Угадайте, какой из вариантов больше поможет в продвижении по службе?

До тех пор, пока одни квалифицированные женщины покидают сферу технологий, обучение других написанию кода не решит проблемы недостатка разносторонних кадров в этой области. Инициативы по выравниванию гендерного баланса часто подразумевают привлечение именно белых женщин, хотя цветные представительницы этого пола испытывают в подобных ситуациях даже больше сопутствующих трудностей. В ходе серии интервью (<https://oreil.ly/t5C6b>), взятых у 60 цветных женщин, работающих в области исследования STEM (науки, технологии, инженерии и математики), выяснилось, что 100 % из них испытывают дискриминацию.

Проблема еще и в том, что в технологической индустрии нарушен сам процесс найма. Одно из исследований, отражающее это нарушение, было проведено Triplebyte (<https://oreil.ly/2Wtw4>) — компанией, которая помогает получить инженерам-программистам работу в организациях, проводя в рамках общего процесса стандартизированное техническое собеседование. У этой компании есть удивительный датасет: результаты того, как 300 инженеров прошли экзамен, объединенные с результатами того, как эти инженеры выступили на собеседовании в различных компаниях. Главный вывод исследования Triplebyte (<https://oreil.ly/2Wtw4>) заключается в том, что «компании часто ищут программистов такого типа, которые мало связаны с тем, чем компания фактически занимается или в чем нуждается. Вместо этого такие сотрудники отражают чаще сложившуюся культуру компании и опыт ее основателей».

Это становится проблемой для тех, кто старается вступить в область глубокого обучения, поскольку большинство ИИ-групп современных компаний были основаны учеными. Эти группы склонны искать людей, похожих на них самих, то есть людей, которые смогут решать сложные математические задачи и понимать профлексику.

Но это открывает большие возможности для компаний, которые готовы выйти за границы статусов и происхождения, сосредоточившись именно на результатах.

Справедливость, ответственность и прозрачность

Профессиональное общество для ученых компьютерных наук ACM проводит конференции, посвященные этике данных, под названием *Conference on Fairness, Accountability and Transparency* (Конференция о справедливости, ответственности и прозрачности) (ACM FAccT), что ранее вписывалось в акроним FAT, но теперь сокращается до менее двусмысленного FAccT. В Microsoft также есть группа, работающая над Fairness, Accountability, Transparency и Ethics (справедливостью, ответственностью, прозрачностью и этикой) в ИИ (FATE). В этом разделе для обозначения справедливости, ответственности и прозрачности мы будем использовать акроним FAccT.

FAccT – это призма, которую люди использовали для рассмотрения этических вопросов. Один из полезных в этом отношении ресурсов – бесплатная онлайн-книга *Fairness and Machine Learning: Limitations and Opportunities* (<https://fairmlbook.org/>) («Справедливость и машинное обучение: возможности и ограничения»), написанная Солоном Барокасом (Solon Barocas) и др. Эта книга «показывает перспективу машинного обучения, рассматривающего справедливость в качестве центральной проблемы, а не второстепенной». Однако она также предупреждает, что «область охвата этого критерия намеренно заужена... Уменьшение значимости этического аспекта машинного обучения может быть соблазнительным для технологов и предпринимателей, так как позволит фокусироваться на технических мерах воздействия, обходя более глубинные вопросы применения силы и предполагаемой в этом случае ответственности. Мы предостерегаем об опасности подобного соблазна». Вместо того чтобы представить обзор подхода FAccT к этике (что намного лучше сделано в книгах наподобие приведенной), мы сосредоточимся на проблемах такого преуменьшения значимости.

Отличный способ рассмотреть, насколько этическая призма полноцenna, – придумать пример, в котором эта призма и наша собственная этическая интуиция дают разнящиеся результаты. Ос Кейес и др. (Os Keyes et al.) провели подобное графическое исследование в своей работе *A Mulching Proposal: Analysing and Improving an Algorithmic System for Turning the Elderly into High-Nutrient Slurry* (https://oreil.ly/_qug9) («Предложение по мульчированию¹: анализ и улучшение алгоритмической системы для превращения пожилых людей в высокопитательную супензию»). В аннотации сказано:

Этические последствия алгоритмических систем активно обсуждались как в HCI, так и в более широком сообществе людей, заинтересованных в технологическом проектировании, разработке и идейных принципах. В этой работе мы исследуем применение одной выдающейся этической схемы – Справедливость, Ответственность и Прозрачность – в отношении предлагаемого алгоритма, который решает различные социальные проблемы, связанные с продовольственной безопасностью и старением населения. Используя разные стандартизованные формы алгоритмического аудита и оценки, мы существенно увеличиваем соответствие данного алгоритма схеме FAT, получая в итоге более этичную и благотворительную систему. Мы также обсуждаем, как это может послужить в качестве ориентира для других исследователей или практиков, стремящихся к получению от алгоритмических систем более положительных с этической точки зрения результатов.

¹ Мульчирование в садоводстве – поверхностное покрытие почвы мульчей (англ. *mulch*) для ее защиты и улучшения свойств. – *Примеч. ред.*

В этой работе достаточно спорное предложение («превращение пожилых людей в высокопитательную супензию») и его результаты («существенно увеличиваем соответствие алгоритма схеме FAT, получая в итоге более этичную и благодетельную систему», мягко говоря... расходятся!)

В философии, и особенно в философии этики, это один из самых эффективных инструментов: сначала придумать процесс, определение, набор вопросов и т. д., который решает задачу. Затем постараться придумать пример, в котором это конкретное решение приводит к предложению, которое никто бы не посчитал приемлемым, и, отталкиваясь от него, доработать изначально найденные решения.

До сих пор мы фокусировались на том, что вы и ваша организация можете делать. Но порой действий отдельного человека или организации недостаточно. Иногда правительства также должны рассматривать последствия осуществляемых ими политических курсов.

Роль политики

Мы часто общаемся с людьми, которые хотят, чтобы проблемы можно было решить с помощью технических или проектных изменений. Например, использовать технический подход для устранения необъективности данных или внедрить руководство к проектированию, чтобы реализуемая технология получилась не столь завлекающей. Несмотря на то что подобные меры могут быть полезны, их не будет достаточно для решения основополагающих проблем, которые и привели к текущему положению вещей. Например, до тех пор пока создание вызывающих «привыкание» технологий будет приносить прибыль, компании продолжат этим заниматься, не обращая внимания на то, что это оказывает побочные эффекты в виде продвижения теорий заговоров и загрязнения информационной экосистемы. Несмотря на то что отдельные разработчики могут стараться подстраивать дизайн продуктов, мы не увидим заметной разницы, пока не изменятся лежащие в их основе стимулы для получения прибыли.

Эффективность регулирования

Чтобы увидеть, какие явления могут заставить компании предпринять конкретные действия, рассмотрите два примера поведения Facebook. В 2018 году в ходе расследования ООН выяснилось, что Facebook сыграл «определенную роль» в продолжающемся геноциде рохинджа, этнического меньшинства Мьянмы, которые, как сказал генеральный секретарь ООН Антонио Гуттерес (Antonio Guterres), «являются одними из, если не самыми дискриминированными людьми в мире». Местные активисты предупреждали руководство Facebook,

что их платформа, начиная с 2013 года использовалась для разжигания вражды и подстрекательства к насилию. В 2015 году Facebook предупредили, что компания может сыграть ту же роль в Мьянме, что и радиовещание в процессе геноцида в Руанде (когда было убито около миллиона людей). Несмотря на это, к концу 2015 года Facebook нанял всего четырех подрядчиков, говорящих на бирманском языке.

Один из приближенных к событиям человек сказал: «Это не ошибка по недосмотру. Масштаб проблемы был значительный, и это уже было очевидно». Марк Цукерберг на слушании в Конгрессе пообещал нанять «десятки» людей для решения проблемы геноцида в Мьянме (в 2018 году, спустя несколько лет после его начала, когда, помимо прочего, с августа 2017 года уже было сожжено не менее 288 деревень на севере штата Ракхайн).

С описанной ситуацией резко контрастирует то, как быстро Facebook нанял 1200 людей в Германии (https://oreil.ly/q_8Dz) для избежания высоких штрафов (вплоть до 50 млн евро) в соответствии с введенным в этой стране законом против разжигания ненависти. Очевидно, что в данном случае Facebook отреагировал скорее на угрозу финансового штрафа, чем на систематическое истребление этнического меньшинства.

В своей статье о проблемах конфиденциальности (<https://oreil.ly/K5YKf>) Мацей Цегловски (Maciej Cegłowski) проводит параллель с движением за защиту окружающей среды:

Этот нормативный акт был настолько успешен в Первом мире, что мы рискуем забыть, какой была жизнь до него. Удушающий смог, от которого сегодня гибнут тысячи людей в Джакарте и Дели, когда-то был символом Лондона (<https://oreil.ly/pLzU7>). Река Кайахога в Огайо регулярно подвергалась возгоранию из-за загрязнения (<https://oreil.ly/qrU5v>). К особенно ужасающим примерам непредвиденных последствий можно отнести добавление тетраэтилового свинца в бензин, что в течение 50 лет приводило к повышению уровня насильтственных преступлений по всему миру. Ни одно из перечисленных пагубных явлений не удалось бы предотвратить, если бы людей попросили голосовать против них «из своего кошелька» или внимательно оценивать политику в отношении окружающей среды, проводящуюся компаниями, которым они доверили свой бизнес, или вообще перестать использовать рассматриваемые технологии. Для их устранения потребовалось скоординированное, а иногда и технически сложное регулирование, выходящее за рамки юрисдикций. В некоторых случаях, например, при запрете коммерческих хладагентов (<https://oreil.ly/o839J>), разрушающих озоновый слой, для введения технических норм потребовалось международное согласие. Мы же сейчас находимся в положении, когда нам необходимо аналогичное изменение взглядов в отношении нашего закона о конфиденциальности.

Права и политика

Чистый воздух и чистая питьевая вода являются общественными благами, которые практически невозможно защитить с помощью единичных рыночных решений. Для этого потребуются скоординированные действия по регулированию. Аналогичным образом, многие негативные воздействия, происходящие из-за непредвиденных последствий неправильного использования технологии, затрагивают общественные блага, примером чего является загрязнение информационной среды или снижение уровня общественной конфиденциальности. Слишком часто неприкосновенность личной жизни определяется как индивидуальное право, хотя внедрение повсеместного отслеживания людей оказывает прямое влияние на общество (что не утратило бы актуальности, даже если бы некоторым людям удалось этого избежать).

Многие из наблюдаемых нами в технологической среде проблем являются проблемами соблюдения человеческих прав, например, когда смещенный алгоритм рекомендует назначить чернокожим подзащитным более длительные тюремные сроки, когда конкретные объявления вакансий показываются только молодым или когда полиция использует распознавание лиц для идентификации протестующих. Подходящим средством для решения проблем соблюдения человеческих прав, как правило, является закон.

При этом нам требуется не только изменить нормативные и правовые акты, но и добиться этического поведения отдельных людей. Изменение индивидуального поведения не способно решить проблему неурегулированных стимулов получения прибыли, внешние следствия экономической деятельности (когда корпорации получают огромные прибыли, перекладывая издержки и негативные последствия на плечи общества) или системные сбои. Однако и закон не всегда может урегулировать все пограничные случаи, поэтому важно, чтобы отдельные разработчики ПО, а также специалисты по данным были подготовлены к принятию этических решений на практике.

Автомобили: исторический прецедент

Перед нами возникают сложные проблемы, и для них нет простых решений. Это может обескураживать, но мы испытываем надежду, вспоминая другие крупномасштабные вызовы, с которыми людиправлялись на протяжении своей истории. Одним из примеров является движение за повышение безопасности автомобилей, рассмотренное в виде исследования в работе *Datasheets for Datasets* (https://oreil.ly/nqG_r) («Таблицы данных для наборов данных»), написанной Тимнитой Гебру (Timnit Gebru), а также в подкасте о проектировании 99 % Invisible (<https://oreil.ly/2HGPd>). В ранних моделях автомобилей не было ремней безопасности, зато были металлические ручки, которые могли бы стать опасными при аварии, обычные стеклянные окна, разбивавшиеся опасным образом, а также

нескладывающиеся рулевые колонки, которые протыкали водителей. При этом автомобильные компании отказывались даже обсуждать безопасность как проблему, которую они могут помочь решить. Было широко распространено мнение, что машины — это просто машины, и причина проблем кроется в самих людях.

Активистам по защите потребительской безопасности и правозащитникам понадобился не один десяток лет, чтобы в рамках общественной дискуссии начало обсуждаться то, что компании все же несут определенную ответственность, которая должна регулироваться нормами. После изобретения складывающейся рулевой колонки ее реализация откладывалась несколько лет, поскольку финансовой мотивации не было. Главная автомобильная компания General Motors даже наняла частных детективов, чтобы те нарыли грязные факты из жизни правозащитника Ральфа Надера (Ralph Nader). Введение требований по наличию ремней безопасности, манекенов для краш-тестов и складных рулевых колонок оказалось самой ощутимой победой. И только в 2011 году автомобильные компании были обязаны начать использовать манекены, представляющие тело среднестатистической женщины, а не мужчины. До этого женщины на 40 % чаще мужчин получали травмы во время аварий схожего характера. Это наглядный пример того, как предвзятость, политика и технологии ведут к серьезным последствиям..

Резюме

У тех, кто имеет за плечами опыт работы с двоичной логикой, недостаток ясных ответов в этике может поначалу вызывать недовольство. Но результат влияния нашей работы на мир, включая непреднамеренные последствия и ее использование в качестве оружия злоумышленниками, является наиболее важным вопросом, который нам необходимо учитывать. Даже несмотря на то, что здесь не бывает легких ответов, есть определенные подводные камни, которых нужно избегать, а также практики, которым нужно следовать, приближаясь к более этическому поведению.

Многие люди (включая нас!) ищут более удовлетворительные и четкие ответы на вопросы о том, как предотвратить вредоносное влияние технологий. Однако, учитывая сложный, далеко идущий и междисциплинарный характер рассматриваемых проблем, простых решений быть не может. Джюлия Ангвин (Julia Angwin), бывший старший репортер ProPublica, сосредоточенная на проблемах алгоритмического смещения и повсеместного отслеживания (в 2016 году участвовавшая в исследовании алгоритма рецидивизма COMPAS, который помог оживить область FAccT), в своем интервью 2019 года (<https://oreil.ly/o7FpP>) сказала следующее:

Я совершенно уверена, что для решения любой проблемы сначала ее необходимо диагностировать и что мы до сих пор находимся именно на этой стадии.

Если вы подумаете о рубеже столетий и переходе в век индустриализации, то можно вспомнить, не знаю, лет тридцать эксплуатации детского труда, неограниченные рабочие часы и ужасные условия производства. В итоге потребовалось очень много журналистских разоблачений и правозащитных действий, чтобы диагностировать и хоть в какой-то степени осознать эти проблемы, после чего уже действия активистов помогли изменить законы. Я чувствую, что мы как будто находимся на втором этапе индустриализации, в этот раз информационной... Своей задачей я считаю пытаться сделать максимально ясными ее обратные стороны и диагностировать их очень точно, чтобы суметь впоследствии решить. Это тяжелая работа, которую должно делать гораздо больше людей.

Обнадеживает, что, по мнению Ангвин, мы все еще находимся в стадии диагностирования: если вы чувствуете, что не до конца поняли описанные проблемы, то это естественно и нормально. Никто еще не изобрел от них «лекарство», но при этом жизненно необходимо продолжать работать, чтобы лучше понять и решить стоящие перед нами проблемы.

Один из рецензентов нашей книги Фред Монро (Fred Monroe) раньше работал в хедж-фондах. Прочитав эту главу, он сказал нам, что многие из рассмотренных в ней проблем (распространенность данных, существенно отличающихся от использованных для обучения модели, влияние петель обратной связи на модель после ее развертывания и масштабирования и пр.) также являлись ключевыми при построении прибыльных моделей для ведения торгов. Действия, которые нужно совершить для учета социальных последствий, во многом будут пересекаться с тем, что нужно делать для учета последствий в отношении организаций, рынка и клиентов. Поэтому тщательный анализ этических аспектов также поможет продумать, как создать продукт успешным по всем направлениям.

Вопросник

1. Даёт ли этика список «верных ответов»?
2. Как может работа с людьми, имеющими разный опыт, помочь в рассмотрении этических вопросов?
3. Какую роль сыграла IBM в нацистской Германии? Почему эта компания участвовала именно так? Почему в этом участвовали ее сотрудники?
4. Как на ситуацию повлиял первый отправленный в тюрьму человек в связи с дизельным скандалом вокруг компании Volkswagen?
5. Какая проблема обнаружилась в базе данных о подозреваемых членах преступных банд, обслуживаемой правоохранительными структурами Калифорнии?

6. Почему рекомендательный алгоритм YouTube рекомендовал видео частично одетых детей педофилам, несмотря на то что такая функция разработчиками из Google в него не закладывалась?
7. Какие проблемы с центрированием наблюдаются в метриках?
8. Почему ресурс Meetup.com не включил учитывание пола в свою рекомендательную систему для технических встреч?
9. Назовите шесть типов необъективности в машинном обучении, определенных Сурешом и Гуттагом.
10. Приведите два примера исторической расовой предвзятости смещения в США.
11. Откуда взято большинство фотографий датасета ImageNet?
12. Почему в работе «Автоматизирует ли машинное обучение моральную угрозу и ошибки?» синусит был отнесен к предикторам инсульта?
13. Что такое репрезентативное смещение?
14. Чем машины отличаются от людей с точки зрения их использования для принятия решений?
15. Дезинформация — это то же, что и фейковые новости?
16. Почему дезинформация, распространяемая путем автоматической генерации текста, является особенно важной проблемой?
17. Какие пять этических призм описывает центр им. Марккулы?
18. В каких случаях политика является подходящим инструментом для решения проблем в области этики данных?

Дополнительные задания

1. Прочтите статью *What Happens When an Algorithm Cuts Your Healthcare* (<https://oreil.ly/5Ziok>) («Что происходит, когда алгоритм урезает ваши медицинские льготы»). Как подобных проблем можно избежать в будущем?
2. Поиските дополнительную информацию о рекомендательной системе YouTube и ее социальном влиянии. Считаете ли вы, что рекомендательные системы должны всегда иметь петли обратной связи с негативными результатами? Какие меры мог бы предпринять Google для их избежания? А правительство?
3. Прочтите работу *Discrimination in Online Ad Delivery* (<https://oreil.ly/jgKpM>) («Дискриминация в онлайн-объявлениях»). Считаете ли вы, что Google должен отвечать за случившееся с доктором Суини? Каков был бы правильный ответ на это?

4. Как может междисциплинарная команда помочь избежать негативных последствий?
5. Прочтите работу *Does Machine Learning Automate Moral Hazard and Error? (Healthcare)* («Автоматизирует ли машинное обучение моральную угрозу и ошибки?»). Какие, на ваш взгляд, должны быть предприняты меры для выявленных в этой работе проблем?
6. Прочтите статью *How Will We Prevent AI-Based Forgery?* (<https://oreil.ly/6mqe4>) («Как нам предотвратить фальсификацию на основе ИИ?»). Как вы думаете, сработал бы предложенный Эциони подход? Почему?
7. Заполните раздел «Анализ проекта» на с. 152.
8. Подумайте, могла бы ваша команда быть более разносторонней в культурно-личностном плане? Если да, то какие в этом помогут методы?

Глубокое обучение на практике: итог!

Поздравляем! Вы добрались до конца первого раздела книги. В этом разделе мы постарались показать, на что способно глубокое обучение и как его можно использовать для создания реальных приложений и продуктов. Вы сможете получить гораздо больше пользы от книги, если выделите время и протестируете все, чему успели здесь научиться. Если вы и так это делали, то вы молодец! А если нет, то сейчас самое время поэкспериментировать.

Если вы еще не были на сайте книги (<https://book.fast.ai/>), то зайдите на него прямо сейчас. Очень важно, чтобы вы все настроили для работы с блокнотами. Успех в области глубокого обучения полностью упирается в практику, поэтому вам необходимо заниматься обучением моделей. Так что не откладывайте начало работы с блокнотами, если еще ими не занялись. При этом не забывайте обращаться к сайту для получения всех важных обновлений и оповещений.

Область глубокого обучения очень быстро меняется, но так как изменить слова книги мы не в силах, рекомендуем обращаться именно к сайту, где вы сможете найти всю последнюю информацию.

Убедитесь, что выполнили следующие шаги.

1. Подключитесь к предпочтительному GPU-серверу Jupyter из списка рекомендованных на сайте.
2. Самостоятельно запустите первый блокнот.
3. Загрузите изображение, найденное в первом блокноте, а затем попробуйте несколько других видов изображений, наблюдая, что произойдет.

4. Выполните второй блокнот, собрав собственный датасет на основе придуманных вами поисковых запросов.
5. Подумайте о том, как вы можете использовать глубокое обучение в собственных проектах, рассмотрев при этом, какие виды данных могут вам походить, а также вероятные проблемы и способы их решения.

В следующей части мы не просто будем смотреть, как глубокое обучение работает на практике, а узнаем, как и почему оно вообще работает. Понимание того, «как и почему», необходимо и для практиков, и для исследователей, потому что в этой относительно новой области почти каждый проект требует определенного уровня кастомизации и отладки. Чем лучше вы поймете основы глубокого обучения, тем качественнее будут ваши модели. Эти основы не столь важны для руководителей, продакт-менеджеров и т. д. (хотя все равно нужно их знать, так что можете смело продолжать читать!), но при этом они критически важны для всех, кто занимается обучением и развертыванием моделей самостоятельно.

ЧАСТЬ II

**Понимание приложений
на базе fastai**

ГЛАВА 4

Обучение классификатора цифр: взгляд изнутри

Изучив, как происходит обучение различных моделей, в главе 2, теперь заглянем внутрь и получше разберемся, что именно там происходит. Начнем с компьютерного зрения, чтобы ознакомиться с основными инструментами и принципами глубокого обучения.

Если говорить точнее, то мы обсудим роль массивов и тензоров, а также трансформации — мощной техники для их эффективного использования. Объясним стохастический градиентный спуск (SGD), механизм для обучения путем автоматического обновления весов. Затем обсудим выбор функции потерь для нашей базовой задачи по классификации и роль мини-пакетов. А еще опишем математику, выполняемую базовой нейронной сетью, и в завершение объединим все это вместе.

В последующих главах мы подробно рассмотрим другие типы приложений и поймем, как обобщаются изученные нами принципы и инструменты. В этой главе мы заложим фундамент, и это делает ее одной из сложнейших глав из-за взаимной зависимости всех рассматриваемых принципов. Можно провести аналогию с аркой, в которой для сохранения общей структуры все камни должны быть на местах. Можно продолжить эту аналогию и сказать, что арка является мощной конструкцией, способной служить основой для других элементов. Но при этом для ее сборки требуется терпение.

Давайте начнем. Первым шагом будет изучение способа представления изображений в компьютере.

Пиксели: основа компьютерного зрения

Для понимания происходящего в модели компьютерного зрения сначала нужно разобраться с процессом обработки изображения компьютером. В этом случае для экспериментов мы используем один из наиболее известных датасетов компьютерного зрения, MNIST (<https://oreil.ly/g3RDg>). MNIST содержит изображения рукописных цифр, которые были собраны Национальным институтом стандар-

тов и технологий и упорядочены в единый датасет машинного обучения Яном Лекуном (Yann Lecun) и его коллегами. Лекун использовал MNIST в 1998 году в LeNet-5 (<https://oreil.ly/LCNEEx>), первой компьютерной системе для демонстрации особенно полезного распознавания последовательностей рукописных цифр. Это был один из важнейших прорывов в истории ИИ.

НАСТОЙЧИВОСТЬ И ГЛУБОКОЕ ОБУЧЕНИЕ

История глубокого обучения написана с настойчивостью и трудолюбием несколькими исследователями. В 1990-х и 2000-х нейронные сети вышли из области интереса, и лишь небольшая группа исследователей продолжила трудиться над их совершенствованием. Троє из них: Ян Лекун, Йошуа Бенджио (Yoshua Bengio) и Джейфри Хинтон (Geoffrey Hinton) — в 2018 году получили высшую награду в области computer science, премию Тьюринга (считающуюся «Нобелевской премией в computer science»), за триумфальный прорыв в условиях глубокого скептицизма и отсутствия интереса у более широкого сообщества машинного обучения и статистики.

Хинтон рассказывал, как академические работы, показывающие намного лучшие результаты, чем любые ранее опубликованные труды, отвергались ведущими журналами и конференциями только потому, что использовали нейронные сети. Работа Лекуна над сверточными нейронными сетями, которые мы начнем изучать в следующем разделе, показала, что эти модели могут читать рукописный текст, чего ранее достичь не удавалось. Но его прорыв был проигнорирован большинством исследователей, несмотря на то что использовался в коммерческих целях для чтения 10 % чеков в США!

Помимо этих трех лауреатов премии Тьюринга, над продвижением общества к его текущей точке прогресса усердствовали и многие другие исследователи, например Юрген Шмидхубер (Jurgen Schmidhuber) (который, по мнению многих, также заслужил премию Тьюринга) разработал множество важных идей, в том числе работая вместе со своим студентом Зеппом Хохрайтером (Sepp Hochreiter) над архитектурой долгой краткосрочной памяти (LSTM) (широко используемой для распознавания речи и других задач по моделированию текста, см. пример IMDb в главе 1). Наиболее же важное открытие сделал Пол Вербос (Paul Werbos), который в 1974 году изобрел алгоритм обратного распространения ошибки для нейронных сетей, технику, показанную в этой главе и повсеместно используемую в обучении нейронных сетей (Werbos, 1994) (<https://oreil.ly/wWIWp>). Его разработка практически полностью игнорировалась на протяжении десятилетий, но сегодня она считается наиболее важной основой передового искусственного интеллекта.

В этом есть урок для каждого из нас! В течение своего путешествия по миру глубокого обучения вы столкнетесь со многими препятствиями, как техническими, так и (еще более сложными) со стороны окружающих людей, которые не будут верить в ваш успех. В этом случае есть один гарантированный способ провала — перестать пробовать. Мы заметили, что единственная общая черта среди всех студентов fast.ai, выросших до уровня первоклассных практиков, — это настойчивость.

В этом начальном уроке мы просто попытаемся только создать модель, способную классифицировать любое изображение как тройку или семерку. Итак, скачаем образец MNIST, содержащий изображения только этих цифр:

```
path = untar_data(URLs.MNIST_SAMPLE)
```

Содержимое каталога можно просмотреть с помощью `ls`, метода, добавленного `fastai`. Этот метод возвращает объект особого `fastai`-класса `L`, у которого помимо всей функциональности встроенного в Python метода `list` есть еще и дополнительная. Одна из его удобных функций в том, что при вводе, прежде чем перечислять сами элементы, он отображает их общее количество (если присутствует более десяти элементов, то он показывает только первые несколько):

```
path.ls()
(#9) [Path('cleaned.csv'), Path('item_list.txt'), Path('trained_model.pkl'), Path(
'models'), Path('valid'), Path('labels.csv'), Path('export.pkl'), Path('history.cs
v'), Path('train')]
```

Датасет MNIST использует распространенную среди ML-датасетов структуру: раздельные каталоги для обучающей и контрольной (и/или тестовой) выборок. Давайте заглянем в обучающую:

```
(path/'train').ls()
(#2) [Path('train/7'), Path('train/3')]
```

Здесь есть каталог с цифрой 3 и каталог с цифрой 7. На языке машинного обучения мы говорим, что в этом датасете «3» и «7» являются *метками* (или целями). Давайте заглянем в один из этих каталогов (используя `sorted` для получения одинаково упорядоченного списка файлов):

```
threes = (path/'train'/'3').ls().sorted()
sevens = (path/'train'/'7').ls().sorted()
threes
(#6131) [Path('train/3/10.png'), Path('train/3/10000.png'), Path('train/3/10011.
png'), Path('train/3/10031.png'), Path('train/3/10034.png'), Path('train/3/10042.
png'), Path('train/3/10052.png'), Path('train/3/1007.png'), Path('train/3/10074.
png'), Path('train/3/10091.png')...]
```

Каталог вполне ожидаемо заполнен файлами изображений. Давайте откроем один из них. Мы видим изображение рукописной цифры 3:

```
im3_path = threes[1]
im3 = Image.open(im3_path)
im3
```



Здесь мы используем класс `Image` из *Python Imaging Library* (PIL), являющейся наиболее широко используемым пакетом Python для открытия, просмотра и управления изображениями. Jupyter знаком с изображениями PIL, поэтому показывает их нам автоматически.

В компьютере все представлено в виде чисел. Для того же, чтобы просмотреть числа, составляющие это изображение, нужно преобразовать его в массив *NumPy*

или *тензор PyTorch*. Например, вот как выглядит раздел изображения, преобразованный в массив NumPy:

```
array(im3)[4:10,4:10]

array([[ 0,  0,  0,  0,  0],
       [ 0,  0,  0,  0, 29],
       [ 0,  0,  0, 48, 166,224],
       [ 0, 93, 244, 249, 253,187],
       [ 0, 107, 253, 253, 230, 48],
       [ 0,  3, 20, 20, 15,  0]], dtype=uint8)
```

Значение 4:10 указывает на то, что мы запросили строки, начиная с индекса 4 (включительно) и до 10, далее то же самое для столбцов. NumPy отсчитывает индексы сверху вниз и слева направо, значит, запрошенная часть расположена близко к верхнему левому углу изображения. А вот аналогичная процедура для тензора PyTorch:

```
tensor(im3)[4:10,4:10]

tensor([[ 0,  0,  0,  0,  0],
       [ 0,  0,  0,  0, 29],
       [ 0,  0,  0, 48, 166,244],
       [ 0, 93, 244, 249, 253,187],
       [ 0, 107, 253, 253, 230, 48],
       [ 0,  3, 20, 20, 15,  0]], dtype=torch.uint8)
```

Мы можем разбить массив, чтобы выбрать только ту часть, где находится верхняя часть цифры, и затем использовать Pandas DataFrame для цветокодирования значений с помощью градиента, что отчетливо покажет нам, как изображение создается из значений пикселей:

```
im3_t = tensor(im3)
df = pd.DataFrame(im3_t[4:15,4:22])
df.style.set_properties(**{'font-size':'6pt'}).background_gradient('Greys')
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	0	0	0	0	0	29	150	195	254	255	254	176	193	150	96	0	0	
2	0	0	0	48	166	224	253	253	234	196	253	253	253	253	233	0	0	
3	0	93	244	249	253	187	46	10	8	4	10	194	253	253	233	0	0	
4	0	107	253	253	230	48	0	0	0	0	0	192	253	253	156	0	0	
5	0	3	20	20	15	0	0	0	0	0	43	224	253	245	74	0	0	
6	0	0	0	0	0	0	0	0	0	0	249	253	245	126	0	0	0	
7	0	0	0	0	0	0	0	14	101	223	253	248	124	0	0	0	0	
8	0	0	0	0	0	11	166	239	253	253	253	187	30	0	0	0	0	
9	0	0	0	0	0	16	248	250	253	253	253	232	213	111	2	0	0	
10	0	0	0	0	0	0	0	43	98	98	208	253	253	253	187	22	0	

Вы можете видеть, что фоновые белые пиксели хранятся как число 0, черные представлены числом 255, а оттенки серого выражаются промежуточными значениями между ними. Все изображение содержит 28 пикселей по горизонтали и 28 — по вертикали, итого 784 пикселей. (Это намного меньше, чем фото с телефона, которое состоит из миллионов пикселей, но для наших обучающих экспериментов такой размер будет очень удобен. Тем не менее вскоре мы дойдем и до больших цветных изображений).

Итак, теперь, когда вы поняли, как видит изображение компьютер, вспомним изначальную цель: создать модель, способную распознавать тройки и семерки. Что можно сделать, чтобы компьютер с этой задачей справился?



ОСТАНОВИТЕСЬ И ПОДУМАЙТЕ!

Прежде чем читать дальше, подумайте, как компьютер может распознать эти две цифры. На какие признаки он имеет возможность посмотреть? Как он может определить эти признаки? Как он может их объединить? Обучение более эффективно, когда вы стараетесь решить задачу самостоятельно, а не просто читаете чьи-то ответы. Поэтому рекомендуем на несколько минут отвлечься от книги, взять листок бумаги с ручкой и набросать свои предположения.

Первая попытка: сходство пикселей

Итак, вот первая мысль: как насчет того, чтобы найти среднее значение для каждого пикселя троек, а затем семерок? Так мы получим две группы средних значений, определяющих то, что мы можем назвать как «идеальные» 3 и 7. Затем для распознавания изображения как одной, так и другой цифры, мы будем смотреть, на какую из этих двух идеальных цифр оно больше всего похоже. Как минимум это уже будет лучше, чем ничего, и вполне сможет послужить нам хорошей базовой моделью.



ТЕРМИН: БАЗОВАЯ МОДЕЛЬ

Простая модель, в достаточно высоком качестве которой вы уверены. Она должна быть легка в реализации и тестировании, чтобы в дальнейшем вы могли проверять каждую усовершенствованную идею, убеждаясь, что они лучше нее. Если не начать с построения приемлемого базового образца, то в дальнейшем будет трудно узнать, насколько успешны создаваемые вами особо сложные модели. Хорошим же подходом к построению такого образца будет начать с размышления о том, как сделать его максимально простым в реализации, аналогично тому, как мы делали недавно. Еще одним неплохим вариантом будет поискать людей, которые решали задачи, аналогичные вашей, чтобы скачать и затем выполнить их код в отношении вашего датасета. В идеале следует попробовать оба этих подхода.

Шаг 1 для нашей простой модели — получить средние значения пикселей для двух наших групп. В процессе мы узнаем многие утонченные хитрости численного программирования Python.

Создадим тензор, содержащий все собранные в стек тройки. Нам уже известно, как создавать тензор, содержащий одно изображение. Для формирования тензора со всеми изображениями директории мы сначала используем генерацию списка Python, создав простой список тензоров, содержащих по одному изображению.

Мы будем использовать Jupyter для выполнения проверок — это поможет убедиться, что возвращается нужное число элементов.

```
seven_tensors = [tensor(Image.open(o)) for o in sevens]
three_tensors = [tensor(Image.open(o)) for o in threes] len(three_tensors),
len(seven_tensors)
(6131, 6265)
```



ГЕНЕРАТОР СПИСКОВ

Генераторы списков и словарей — это прекрасная возможность Python. Многие программисты, включая авторов этой книги, используют их ежедневно: эти функции являются частью «идиоматического Python». Но для программистов, приходящих из других языков, они могут показаться незнакомыми. В интернете можно найти много отличных уроков по их использованию, так что мы не станем заострять внимание на этом моменте. Вот краткое объяснение и пример, чтобы вы могли начать ими пользоваться. Генератор списка выглядит так: `new_list = [f(o) for o in a_list if o>0]`. Он возвращает каждый элемент из `a_list`, который больше 0, после его передачи в функцию `f`. Это выражение состоит из трех частей: коллекции, которую вы перебираете (`a_list`), необязательного фильтра (`if o>0`) и некоторого действия для каждого элемента (`f(o)`). Такой способ короче в записи, но также оказывается в разы быстрее альтернативных способов создания списков с помощью цикла.

Мы также убедимся на одном из изображений, что оно отображается нормально. Поскольку теперь у нас есть тензоры (которые Jupyter по умолчанию будет выводить как значения), а не изображения PIL (которые Jupyter по умолчанию будет показывать как изображения), нам для их отображения нужно использовать fastai-функцию `show_image`:

```
show_image(three_tensors[1]);
```



Для каждой позиции пикселя нужно вычислить среднее значение его интенсивности по всем изображениям. Для этого мы сначала объединим все изображения этого списка в один трехмерный тензор. Самый распространенный способ описать подобный тензор — это назвать его *тензором ранга 3*. Часто нам приходится обь-

единять отдельные тензоры коллекции в один тензор. Неудивительно, что PyTorch содержит функцию `stack`, которую мы как раз можем использовать для этой цели.

Некоторые операции в Python, такие как получение среднего значения, требуют *приведения* целочисленных тензоров в к типу `float`. Поскольку нам это понадобится дальше, сразу приведем наш собранный в стек тензор к этому типу. Преобразование в PyTorch выполняется простым написанием имени типа, к которому нужно привести имеющийся, и рассмотрением его в качестве метода.

Как правило, когда изображения представлены числами с плавающей запятой, значения пикселей расположены между 0 и 1, поэтому здесь мы выполним деление на 255:

```
stacked_sevens = torch.stack(seven_tensors).float()/255
stacked_threes = torch.stack(three_tensors).float()/255
stacked_threes.shape

torch.Size([6131, 28, 28])
```

Наиболее важным атрибутом тензора можно назвать его *форму*. Она говорит вам о длине каждой оси. В нашем случае мы видим, что у нас есть 6131 изображение, каждое размером 28×28 пикселей. В этом тензоре нет ничего конкретного, что говорило бы о том, что первая ось соответствует количеству изображений, вторая — их высоте, а третья — ширине: семантика тензора полностью зависит от нас и от того, как мы его построим. Фактически это просто куча чисел в памяти.

Длина в форме тензора — это его ранг:

```
len(stacked_threes.shape)

3
```

Очень важно, чтобы вы запомнили и практиковали эти элементы терминологии: *ранг* — это количество осей, или измерений, тензора; *форма* — это размер каждой оси тензора.



СЛОВО АЛЕКСИСУ

Будьте внимательны, так как термин «измерение» иногда используется и в другом значении. Учитывайте, что мы живем в «трехмерном пространстве», где физическое местоположение может быть описано вектором `v` с длиной 3. Но согласно PyTorch атрибут `v.ndim` (который выглядит как «число измерений» `v`) равен одному, а не трем! Почему? Потому что `v` — это вектор, являющийся тензором ранга 1, то есть у него есть всего одна ось (даже если у этой оси длина 3). Другими словами, иногда измерение используется в отношении размера оси («пространство трехмерно»), в то время как в других случаях оно используется для ранга или числа осей («матрица имеет два измерения»). Когда я начинаю путаться, мне помогает перевод всех инструкций в термины «ранг», «ось» и «длина», которые не вызывают противоречий.

Мы также можем получить ранг тензора напрямую с помощью `ndim`:

```
stacked_threes.ndim
```

```
3
```

В завершение мы можем определить, как выглядит идеальная 3. Для этого мы вычисляем среднее всех тензоров изображений, взяв среднее значение по измерению 0 нашего сборного тензора ранга 3. 0 – это измерение, которым индексируются все изображения.

Другими словами, так мы вычислим среднее значение каждого пикселя (его позиции) во всех изображениях. В результате получится по одному значению для каждой позиции пикселя, или одно общее изображение. Вот:

```
mean3 = stacked_threes.mean(0)  
show_image(mean3);
```



Согласно данному датасету, это идеальная цифра 3! (Вам это может не понравиться, но именно так выглядит цифра 3 максимального качества.) Здесь мы видим, что она очень темна в местах, где все изображения подтверждают, что так и должно быть, но становится дымчатой и размытой в местах, где с этим согласны не все изображения.

Сделаем то же самое для семерок, но теперь совместим все шаги вместе, чтобы сэкономить время:

```
mean7 = stacked_sevens.mean(0)  
show_image(mean7);
```



Теперь выберем произвольную тройку и измерим ее *удаленность* от «идеальных цифр».



ОСТАНОВИТЕСЬ И ПОДУМАЙТЕ!

Как бы вы вычислили степень сходства конкретного изображения с каждой из наших идеальных цифр? Как и прежде, рекомендуем отложить книгу и набросать размышления на бумаге. Исследование показывает, что запоминание и понимание улучшаются, когда в процессе обучения вы решаете задачи, экспериментируете и самостоятельно прорабатываете новые идеи.

Вот образец 3:

```
a_3 = stacked_threes[1]
show_image(a_3);
```



Как определить удаленность от идеальной тройки? Мы не можем просто сложить разности между пикселями этого изображения и идеальной цифры. В этом случае одни разности будут положительными, а другие — отрицательными, компенсируя первые. Это, в свою очередь, приведет к ситуации, в которой общая разность между идеальным изображением и изображением, местами слишком темным, а местами слишком светлым, может отражаться как нулевая. Такой результат окажется ошибочным.

Чтобы этого избежать, специалисты по данным в этом контексте используют два основных способа измерения удаленности.

- Взять среднее *абсолютного значения* разностей (абсолютное значение — это функция, заменяющая отрицательные значения положительными). Это называется *средней абсолютной разностью, или нормой L1*.
- Получить среднее *квадратов* разностей (которое делает все положительным), а затем извлечь из него *квадратный корень* (обратив начальное возведение в квадрат). Это называется *среднеквадратичной ошибкой (RMSE), или нормой L2*.



НИЧЕГО СТРАШНОГО, ЕСЛИ ВЫ ПОДЗАБЫЛИ МАТЕМАТИКУ

В этой книге мы предполагаем, что вы прошли курс старшей школы и хоть кое-что из него помните, — но все мы что-то да забываем! Все зависит от того, чем вы занимались раньше. Возможно, вы забыли, что такое *квадратный корень* или непосредственно принцип его действия. Ничего страшного! Всякий раз, когда вы сталкиваетесь с математическим понятием, которое не объяснено в книге, остановитесь и поищите информацию. Убедитесь, что поняли его основной смысл, как оно работает и почему мы его используем. Одним из лучших мест, где можно освежить знания, является Академия Хана, которая дает отличное введение в тему квадратных корней (<https://oreil.ly/T7mxH>).

Попробуем оба описанных варианта:

```
dist_3_abs = (a_3 - mean3).abs().mean()
dist_3_sqr = ((a_3 - mean3)**2).mean().sqrt()
dist_3_abs, dist_3_sqr

(tensor(0.1114), tensor(0.2021))
dist_7_abs = (a_3 - mean7).abs().mean()
```

```
dist_7_sqr = ((a_3 - mean7)**2).mean().sqrt()  
dist_7_abs,dist_7_sqr  
(tensor(0.1586), tensor(0.3021))
```

В обоих случаях удаленность рассматриваемой тройки от ее «идеала» меньше, чем от идеальной семерки, значит, наша простая модель будет давать в этом случае правильный прогноз.

В PyTorch уже есть обе эти *функции потерь*. Вы можете найти их в модуле `torch.nn.functional`, который команда PyTorch рекомендует импортировать как `F` (и который по умолчанию доступен под этим именем в `fastai`):

```
F.l1_loss(a_3.float(),mean7), F.mse_loss(a_3,mean7).sqrt()  
(tensor(0.1586), tensor(0.3021))
```

Здесь `MSE` обозначает среднеквадратичную ошибку, а `L1` относится к стандартному математическому термину *среднее абсолютное значение* (в математике оно называется *нормой L1*).

СЛОВО СИЛЬВЕЙНУ



Наглядное отличие между нормой $L1$ и среднеквадратичной ошибкой (`MSE`) в том, что последняя ярче отражает серьезные ошибки, чем первая, но при этом менее выразительна в отношении мелких.

СЛОВО ДЖЕРЕМИ



Когда я впервые столкнулся с $L1$, то искал информацию о том, что это вообще такое. Я выяснил, что это *векторная норма*, использующая *абсолютное значение*. Потом я решил выяснить, что такое «векторная норма», и начал читать: *при рассмотрении векторного пространства V над полем F вещественных или комплексных чисел норма V будет любой функцией с неотрицательными значениями p : $V \rightarrow [0, +\infty)$, обладающей следующими свойствами:* для всех $a, b \in F$ и всех $u, v \in V$, $p(u + v) \leq p(u) + p(v)$... и тут я читать перестал. «Мда... Мне никогда не понять математику!» — подумал я в тысячный раз. С тех пор я уяснил, что при каждом появлении таких математических терминов на практике я могу заменить их небольшим кусочком кода. Например, потеря $L1$ просто равна $(a - b).abs().mean()$, где a и b являются тензорами. Мне кажется, что математики просто мыслят не так, как я... По ходу книги я дам для каждого математического термина подходящий для него элемент кода, попутно объяснив понятным языком, что вообще происходит.

Мы только что произвели различные математические операции над тензорами PyTorch. Если вы уже занимались численным программированием в NumPy, то можете заметить, что тензоры похожи на массивы NumPy. Разберем две эти важные структуры данных.

Массивы NumPy и тензоры PyTorch

NumPy (<https://numpy.org/>) — это самая широко используемая библиотека для научного и численного программирования в Python. Она предоставляет аналогичную PyTorch функциональность и API, но при этом не поддерживает использование GPU или вычисление градиентов, что является критически важным в глубоком обучении. Исходя из этого, мы будем везде по возможности использовать вместо массивов NumPy тензоры PyTorch.

(Обратите внимание, что fastai добавляет в NumPy и PyTorch ряд функций, чтобы повысить их сходство. Если какой-либо код из этой книги не заработает у вас на компьютере, это может означать, что вы забыли включить в начало блокнота следующую строку: `from fastai.vision.all import *`.)

Но что вообще такие массивы и тензоры и почему это должно вас волновать?

Python весьма медлителен по сравнению со многими языками. Все быстрое в Python, NumPy или PyTorch, как правило, является оберткой для компилированного объекта, написанного (и оптимизированного) в другом языке, в частности C. В действительности *массивы NumPy и тензоры PyTorch могут выполнять вычисления в тысячи раз быстрее, чем чистый Python*.

Массив NumPy — это многомерная таблица данных, содержащая элементы одного типа. Поскольку он может быть абсолютно любого типа, то существуют массивы массивов, где углубленные массивы могут иметь разные размеры, — такая структура называется *ступенчатым массивом*. Под «многомерной таблицей» мы подразумеваем, к примеру, список (одно измерение), таблицу или матрицу (два измерения), таблицу таблиц или куб (три измерения) и т. д. Если эти элементы все простого типа (целые числа или числа с плавающей запятой), NumPy будет хранить их в памяти в виде компактной структуры данных C. Именно здесь NumPy и проявляет себя. В нем есть большое количество операторов и методов, которые могут выполнять над этими компактными структурами вычисления с той же скоростью, что и оптимизированный C, потому что написаны они именно в оптимизированном C.

Тензор PyTorch — это почти то же, что и массив NumPy, но с небольшим ограничением, раскрывающим дополнительные возможности. Сходство в том, что он является многомерной структурой данных с элементами одного типа. Ограничение заключается в невозможности использовать любой стандартный тип и в необходимости использовать для всех компонентов один базовый численный тип. В результате тензор не настолько гибок, как настоящий массив массивов. Например, тензор PyTorch не может формировать ступени — он всегда представляет многомерную прямоугольную структуру правильной формы.

Подавляющее большинство методов и операторов, поддерживаемых NumPy в отношении этих структур, также поддерживаются в PyTorch, но тензоры PyTorch

имеют дополнительные возможности. Одна из возможностей заключается в том, что эти структуры могут существовать в GPU, и в этом случае их вычисление оптимизируется для GPU и может выполняться гораздо быстрее (при наличии большого количества значений для обработки). Кроме того, PyTorch может автоматически вычислять производные этих операций, включая комбинации операций. Как вы увидите, без этого функционала было бы невозможно глубокое обучение на практике.



СЛОВО СИЛЬВЕЙНУ

Если вы не знакомы с языком C, то не волнуйтесь: он совсем вам не понадобится. Если кратко, то это низкоуровневый (низкоуровневый означает более схожий с языком, который использует компьютер) язык, который намного быстрее Python. При программировании в Python постараитесь как можно реже писать циклы, заменяя их командами, работающими над массивами и тензорами непосредственно.

Возможно, самым важным навыком Python-программиста является способность эффективно использовать API массива/тензора. Чуть позже мы покажем много приемов, а пока приведем краткую сводку того, что нужно знать.

Чтобы создать массив или тензор, передайте список (или список списков, или список списков списков и т. д.) в `array` или `tensor`:

```
data = [[1,2,3],[4,5,6]]
arr = array (data)
tns = tensor(data)

arr # numpy
array([[1, 2, 3],
       [4, 5, 6]])

tns # pytorch
tensor([[1, 2, 3],
       [4, 5, 6]])
```

Все последующие операции показаны для тензоров, но для массивов NumPy синтаксис и результаты будут идентичными.

Вы можете выбрать строку (обратите внимание, что, подобно спискам в Python, тензоры индексируются с 0, то есть 1 относится ко второй строке/столбцу):

```
tns[1]
tensor([4, 5, 6])
```

или столбец, используя `:` для указания *всей первой оси* (иногда мы называем измерения тензоров/массивов *осями*):

```
tns[:,1]  
tensor([2, 5])
```

Их можно комбинировать с Python-синтаксисом среза (`[start:end]`, исключив `end`), чтобы выбрать часть строки или столбца:

```
tns[1,1:3]  
tensor([5, 6])
```

Также можно использовать стандартные операторы, такие как `+`, `-`, `*` и `/`:

```
tns+1  
tensor([[2, 3, 4],  
       [5, 6, 7]])
```

Тензор имеет тип:

```
tns.type()  
'torch.LongTensor'
```

и будет автоматически его менять при необходимости. Например, с `int` на `float`:

```
tns*1.5  
tensor([[1.5000, 3.0000, 4.5000],  
       [6.0000, 7.5000, 9.0000]])
```

Итак, насколько хороша наша базовая модель? Чтобы это вычислить, нам нужно определить метрику.

Вычисление метрик с помощью бродкастинга (Broadcasting)

Вспомните, что *метрика* — это число, вычисляемое на основе прогнозов модели и верных меток датасета и отражающее качество модели. Мы можем, к примеру, использовать одну из функций, которые видели в предыдущем разделе, а именно функцию вычисления среднеквадратичной ошибки или средней абсолютной ошибки, взяв в результате их среднее значение для всего датасета. Тем не менее ни одно из этих чисел не будет в достаточной степени понятно большинству людей, поэтому обычно на практике для моделей классификации мы используем в качестве метрики *точность*.

Необходимо вычислять метрику в отношении *контрольной выборки*. Таким образом мы избегаем случайного переобучения, в результате которого наша мо-

дель будет давать хорошие показатели только для обучающих данных. В случае с моделью определения сходства пикселей такой риск практически отсутствует, поскольку в ней нет обученных компонентов, но мы все равно используем контрольную выборку, чтобы соблюсти стандартный подход и быть готовыми ко второй попытке в дальнейшем.

Для создания контрольной выборки нам нужно полностью исключить некоторые данные из обучения, чтобы модель не видела их вообще. Как выясняется, создатели датасета MNIST уже сделали это за нас. Помните, что у нас был целый отдельный каталог `valid?` Как раз для этого он и нужен.

Сначала создадим тензоры для наших троек и семерок из этого каталога. Эти тензоры будем использовать для вычисления метрики, измеряющей качество нашей первичной модели, определяющей удаленност от идеального изображения:

```
valid_3_tens = torch.stack([tensor(Image.open(o))
                            for o in (path/'valid'/'3').ls()])
valid_3_tens = valid_3_tens.float()/255
valid_7_tens = torch.stack([tensor(Image.open(o))
                            for o in (path/'valid'/'7').ls()])
valid_7_tens = valid_7_tens.float()/255
valid_3_tens.shape, valid_7_tens.shape

(torch.Size([1010, 28, 28]), torch.Size([1028, 28, 28]))
```

Будет здорово выработать привычку проверять формы по ходу работы. Здесь мы видим два тензора, один из которых представляет контрольную выборку троек, состоящую из 1010 изображений размером 28×28 , а второй — контрольную выборку семерок, которая содержит 1028 изображений такого же размера.

В итоге нужно написать функцию `is_3`, которая будет решать, написана на случайном изображении цифра 3 или 7. Для этого ей потребуется определить, к какому из «идеальных» изображений окажется ближе это случайное. Мы же, в свою очередь, должны для этого определить понятие *удаленности*, то есть функцию, вычисляющую удаленност между двумя изображениями.

Можно написать простую функцию, которая вычисляет среднюю абсолютную ошибку, используя выражение, похожее на то, которое мы писали в предыдущем разделе:

```
def mnist_distance(a,b): return (a-b).abs().mean((-1,-2))
mnist_distance(a_3, mean3)

tensor(0.1114)
```

Это то же значение, которое мы ранее вычисляли, определяя удаленност между этими двумя изображениями: идеальной тройки `mean_3` и случайной тройки `a_3`; обе являются тензорами одного изображения с формой $[28, 28]$.

Но чтобы вычислить метрику для общей точности, понадобится вычислить удаленность от идеальной тройки для *каждого* изображения в контрольной выборке. Как же произвести такое вычисление? Мы могли бы написать цикл, проходящий по всем тензорам отдельных изображений, собранных в тензор контрольной выборки, `valid_3_tens`, который имеет форму `[1010, 28, 28]` и представляет 1010 изображений. Но есть способ получше.

Произойдет кое-что интересное, когда мы возьмем ту же функцию вычисления удаленности, написанную для сравнения двух изображений, но передадим ей в качестве аргумента `valid_3_tens`, тензор, который представляет контрольную выборку троек:

```
valid_3_dist = mnist_distance(valid_3_tens, mean3)
valid_3_dist, valid_3_dist.shape

(tensor([0.1050, 0.1526, 0.1186, ..., 0.1122, 0.1170, 0.1086]),
 torch.Size([1010]))
```

Она не будет жаловаться на несовпадение форм, а вернет показатель удаленности для каждого изображения в виде вектора (то есть тензора 1 ранга) длиной 1010 (число троек в контрольной выборке). Как это произошло?

Взгляните еще раз на нашу функцию `mnist_distance`, и вы увидите, что в ней есть вычитание `(a-b)`. Магия здесь в том, что когда PyTorch пытается выполнить простую операцию вычитания между двумя тензорами разных рангов, он использует *бродкастинг*, то есть автоматически расширяет тензор меньшего ранга до размера тензора старшего ранга. Бродкастинг — это важная возможность, существенно упрощающая написание кода тензора.

После бродкастинга, когда два тензора уже имеют одинаковый ранг, PyTorch применят свою обычную логику для тензоров одного ранга: выполняет операцию для каждого соответствующего их элемента и возвращает результат в виде тензора. Например:

```
tensor([1,2,3]) + tensor([1])

tensor([2, 3, 4])
```

Итак, в данном случае PyTorch обрабатывает `mean3`, тензор второго ранга, представляющий одно изображение, как если бы это были 1010 копий одного и того же изображения, а затем вычитает каждую из этих копий из каждой тройки нашей контрольной выборки. Какая, на ваш взгляд, у этого тензора будет форма? Попробуйте выяснить это сами, прежде чем посмотреть ответ:

```
(valid_3_tens-mean3).shape

torch.Size([1010, 28, 28])
```

Мы вычисляем удаленность между идеальной тройкой и каждой из 1010 троек контрольной выборки, для каждого изображения размером 28×28 , получая форму [1010, 28, 28].

Процесс реализации транслирования содержит два важных момента, которые делают его ценным не только в плане выразительности, но и в плане производительности.

- PyTorch в действительности не копирует `mean3` 1010 раз. Он *представляет*, что это был тензор такой формы, но не выделяет для него дополнительной памяти.
- Все вычисления он производит на языке C (или, если вы используете GPU, в CUDA, эквиваленте C в GPU), в десятки тысяч раз быстрее, чем на чистом Python (и до миллионов раз быстрее с помощью GPU!).

Это относится ко всем поэлементным операциям, транслированию и функциям, выполняемым в PyTorch. *Это самая важная техника, необходимая вам для создания эффективного PyTorch кода.*

Далее в `mnist_distance` мы видим метод `abs`. Теперь вы можете догадаться, какую роль в отношении тензора он играет. Функция применяет этот метод к каждому отдельному элементу тензора и возвращает тензор результатов (что как раз и означает *поэлементное* применение метода). В результате мы получаем 1010 матриц абсолютных значений.

И наконец, наша функция вызывает `mean((-1, -2))`. Кортеж $(-1, -2)$ представляет диапазон осей. В Python -1 означает последний элемент, а -2 — предпоследний. Поэтому в данном случае эта инструкция сообщает PyTorch, что мы хотим получить среднее, охватывающее значения, проиндексированные двумя последними осями тензора. Две последние оси — это горизонтальное и вертикальное измерения изображения. После этого у нас остается только первая ось тензора, индексирующая все изображения, в связи с чем итоговый размер и получился (1010). Другими словами, для каждого изображения мы усреднили интенсивность всех его пикселей.

На протяжении всей книги и особенно в главе 17 мы еще многое узнаем о бродкастинге и неоднократно его рассмотрим на примерах.

Мы можем использовать `mnist_distance`, чтобы выяснить, является ли изображение тройкой, используя следующую логику: если удаленность от рассматриваемой цифры до идеальной тройки меньше, чем удаленность от идеальной семерки, тогда это тройка. Функция будет автоматически выполнять транслирование и применяться поэлементно, точно так же, как все функции и операторы PyTorch:

```
def is_3(x): return mnist_distance(x, mean3) < mnist_distance(x, mean7)
```

Протестируем ее на нашем примере:

```
is_3(a_3), is_3(a_3).float()  
(tensor(True), tensor(1.))
```

Обратите внимание, что в случае преобразования логического ответа в число с плавающей запятой мы получаем `1.0` для `True` и `0.0` для `False`.

Благодаря бродкастингу мы можем также протестировать ее для всей контрольной выборки троек:

```
is_3(valid_3_tens)  
tensor([True, True, True, ..., True, True, True])
```

Теперь мы можем вычислить точность для каждой тройки и семерки, получив среднее значение этой функции для всех троек и ее обратное значение для всех семерок:

```
accuracy_3s = is_3(valid_3_tens).float().mean()  
accuracy_7s = (1 - is_3(valid_7_tens).float()).mean()  
  
accuracy_3s, accuracy_7s, (accuracy_3s+accuracy_7s)/2  
(tensor(0.9168), tensor(0.9854), tensor(0.9511))
```

Выглядит неплохо. Получаемая нами точность выше 90 % как для троек, так и для семерок, и мы видели, как можно удобно определять метрику с помощью трансляирования. Но давайте взглянем правде в глаза: тройки и семерки очень отличаются внешне, да и классифицируем мы пока что только две из десяти возможных цифр. Все это говорит о том, что далее нам потребуется добиться лучших показателей!

А для того чтобы показатели улучшить, пора опробовать систему, выполняющую реальное обучение: ту, которая может автоматически изменять саму себя для улучшения производительности. Другими словами, настало время поговорить о процессе обучения и SGD.

Стochastic gradient descent

Помните приведенный нами в главе 1 способ, которым Артур Сэмюэл описал реализацию машинного обучения?

Предположим, мы организуем некоторые автоматические средства проверки эффективности любого текущего назначения веса с позиции фактической производительности и предоставляем механизм изменения назначения веса для максимального увеличения этой производительности. Не нужно

углубляться в детали такой процедуры, чтобы увидеть, что она полностью автоматизирована, и понять, что машина, запрограммированная таким образом, будет учиться на собственном опыте.

Как мы уже говорили, это ключ, который позволит получить постоянно совершенствующуюся модель — модель, способную учиться. Но наше решение для определения сходства пикселей пока этого не делает. У нас нет никаких методик назначения весов или способа улучшения, основанного на тестировании эффективности этого назначения. Другими словами, мы не можем улучшить используемый механизм определения сходства пикселей изменением набора параметров. Чтобы воспользоваться таким преимуществом глубокого обучения, сначала нужно представить нашу задачу так, как определил Сэмюэл.

Вместо того чтобы стараться найти сходство между изображением и его «идолом», мы можем просмотреть каждый отдельный пиксель и придумать для всех них веса так, чтобы максимальные ассоциировались с наиболее темными пикселями для определенной категории. Например, пиксели, расположенные ближе к правому нижнему углу, вряд ли будут активироваться для семерки, значит, для этой цифры они должны иметь низкие веса. Для цифры 3 их активация потребуется, значит, и веса в этом случае должны быть высокими. Это можно представить в виде функции, установив значения весов для каждой возможной категории: например, вероятности того, что наблюдается цифра 3:

```
def pr_three(x,w) = (x*w).sum()
```

Здесь мы предполагаем, что x — это изображение, представленное в виде вектора, иначе говоря, одной длинной строки, составленной из всех строк этого изображения. Мы также предполагаем, что веса — это вектор w . При наличии такой функции нужен только способ обновления весов, приближающих их к наилучшему значению. С таким подходом мы можем повторять этот шаг много раз, постепенно улучшая веса до тех пор, пока они не станут максимально подходящими.

Нужно найти конкретные значения для вектора w , благодаря которым результат функции будет высоким для изображений с тройками и низким для других. Поиск наилучшего вектора w — это способ поиска наилучшей функции распознавания троек. (Так как мы пока не используем глубокую нейронную сеть, то ограничены в возможностях нашей функции, но чуть позже мы устраним это ограничение.)

Если говорить точнее, то для превращения этой функции в классификатор машинного обучения потребуется выполнить следующие этапы.

1. *Инициализировать* веса.
2. Использовать эти веса для каждого изображения, чтобы *прогнозировать*, является оно 3 или 7.

3. На основе полученных прогнозов вычислить, насколько точна модель (*узнать потери*).
4. Вычислить *градиент*, который в отношении каждого веса измеряет, как его изменение может повлиять на потери.
5. На основе выполненных вычислений сдвинуть все веса на один шаг.
6. Вернуться к шагу 2 и *повторить* процесс.
7. Продолжать до тех пор, пока вы не решите *остановить* процесс (например, когда модель станет достаточно точной или вы просто не захотите ждать дольше).

Эти семь этапов, показанные на рис. 4.1, являются ключевыми в обучении всех моделей глубокого обучения. То, что глубокое обучение полностью полагается именно на них, чрезвычайно удивительно и даже парадоксально. Не верится, что этот процесс может решать настолько сложные задачи. Но позже вы увидите, что он действительно с этим справляется.

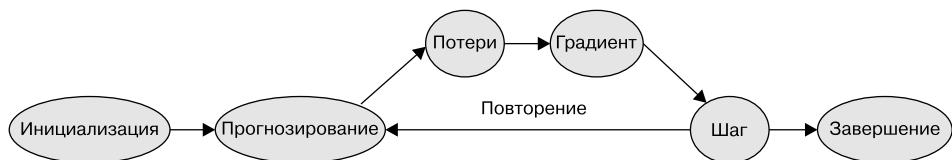


Рис. 4.1. Процесс градиентного спуска

Для выполнения каждого из этих этапов есть множество способов, которые мы будем изучать на протяжении оставшейся части книги. Эти нюансы представляют огромное значение для практиков глубокого обучения, но на деле оказывается, что основной подход для каждого этапа следует определенным базовым принципам. Вот несколько пояснений.

Инициализация

Мы инициализируем параметры со случайными значениями, что может в некоторой степени удивлять. Конечно же, мы можем выбрать и другие варианты, например инициализировать их как процент количества активаций данного пикселя для данной категории, но поскольку нам уже известен стандартный процесс улучшения весов, выясняется, что их инициализация со случайными значениями работает отлично.

Потери

Это понятие Сэмюэл подразумевал, когда говорил о тестировании эффективности любого текущего назначения веса с позиции его текущей производительности. Нам нужна функция, которая при высокой производительности модели будет возвращать малое число (обычно для оценки модели мы

принимаем малые потери как хороший показатель, а большие — как плохой, хотя это просто условность).

Шаг

Для выяснения направления, в котором следует изменять веса, можно просто поэкспериментировать, то есть просто попробовать немного их увеличить и посмотреть, поникаются потери или растут. Как только вы находитите верное направление, далее можно менять значение весов еще чуть больше или меньше до тех пор, пока не найдете оптимальную величину. Тем не менее это будет весьма медленно. Как мы скоро увидим, магия дифференциального исчисления (*calculus*) позволяет выяснить, в каком направлении и на сколько примерно изменить каждый вес, не прибегая ко всем перечисленным шагам. Для этого нам потребуется вычислить градиенты. Такое решение представляет собой простую оптимизацию производительности, и мы могли бы получить точно такие же результаты, используя более медленный ручной процесс.

Остановка

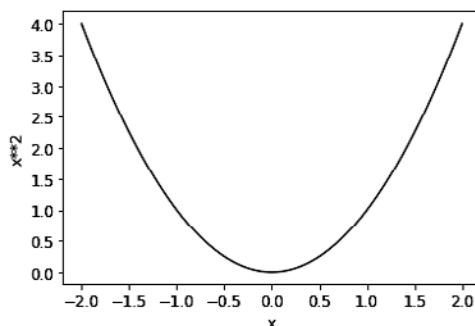
Как только мы решили, сколько эпох обучать модель (несколько предложений на эту тему приводилось в списке ранее), мы это решение применяем. Для нашего классификатора цифр мы бы продолжили обучение либо до момента, когда точность модели начнет падать, либо до истечения выделенного времени.

Прежде чем применять эти этапы к нашей задаче классификации изображений, проиллюстрируем их на простом примере. Сначала определим очень простую квадратичную функцию и представим, что она является нашей функцией потерь, а x — ее параметром:

```
def f(x): return x**2
```

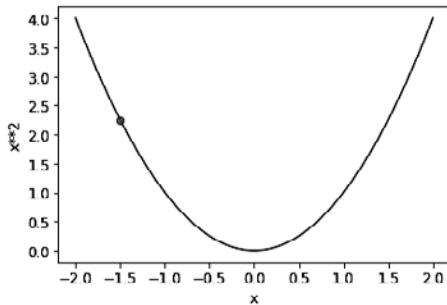
Вот график для этой функции:

```
plot_function(f, 'x', 'x**2')
```

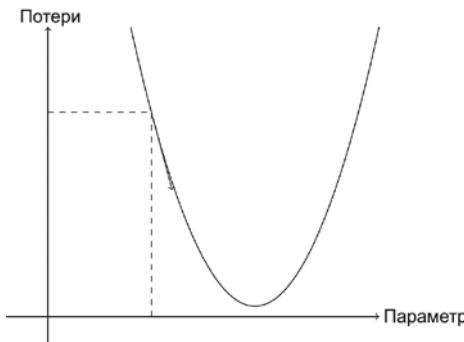


Последовательность описанных этапов начинается с выбора случайного значения для параметра и вычисления значения потерь:

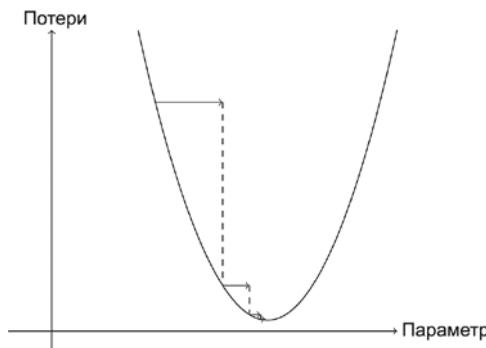
```
plot_function(f, 'x', 'x**2')
plt.scatter(-1.5, f(-1.5), color='red');
```



Далее мы смотрим, что произойдет, если мы немного увеличим или уменьшим параметр: произведем *корректировку*. Это просто наклон в определенной точке:



Можно немного изменить вес по направлению наклонной, вычислить потери и величину корректировки, повторив это несколько раз. В итоге мы дойдем до самой нижней точки нашей кривой:



Основная идея восходит еще к Исааку Ньютону, который определил, что мы можем оптимизировать таким способом случайные функции. Независимо от того, насколько сложными становятся наши функции, этот базовый подход градиентного спуска существенно изменяться не будет. Единственные небольшие изменения, которые мы еще увидим, будут заключаться в применении удобных способов для ускорения процесса путем нахождения оптимальных шагов.

Вычисление градиентов

Единственный магический шаг во всем процессе — это вычисление градиентов. Как уже говорилось, мы используем в качестве оптимизации производительности дифференциальное исчисление. Это позволяет более точно вычислить, увеличивающиеся или уменьшающиеся потери, когда мы скорректируем параметры вверх или вниз. Иначе говоря, градиенты покажут нам, насколько сильно нам нужно изменить каждый вес для улучшения модели.

Из курса матанализа вы можете помнить, что *производная* функции сообщает, какое изменение ее параметров приведет к изменению результатов. Если не помните — не волнуйтесь. Многие из нас все забыли, только окончив школу или вуз. Но чтобы продолжить, понадобится некоторое интуитивное понимание производной. Если эта тема вам непонятна, то в Академии Хана можно пройти курс уроков, посвященных простым производным (<https://oreil.ly/nyd0R>). Вам не потребуется знать, как вычислять их самостоятельно, важно лишь разобраться, что это такое.

Суть производной в том, что ее можно вычислить для любой функции, такой как квадратичная, которую мы видели в предыдущем разделе. В свою очередь, производная тоже является функцией, которая вычисляет не значение, а изменение. Например, производная квадратичной функции при значении 3 говорит нам, насколько быстро изменяется эта функция при значении 3. Если говорить еще конкретнее, то вы можете вспомнить, что градиент определяется как *подъем/пробег*, то есть как изменение значения функции, деленное на изменение значения параметра. Когда известно, как будет меняться функция, то известно, что нужно делать для ее уменьшения. Это ключевой момент в машинном обучении: наличие способа изменять параметры функции, уменьшая ее результат. Дифференциальное исчисление предоставляет нам вычислительное сокращение, а именно производную, позволяющую непосредственно вычислять градиенты функций.

Среди прочего важно понимать, что у функции есть много весов, которые нужно скорректировать, поэтому при вычислении производной мы будем получать не одно число, а множество — градиент для каждого веса. Но математической сложности в этом процессе нет. Вы можете вычислить производную по отношению к одному весу и рассмотреть остальные как константы, после чего повторить эту процедуру для каждого веса. Именно так и вычисляются все градиенты для всех весов.

Мы только что говорили о том, что вычислять градиенты самим не потребуется. Как же так? Это несколько удивительно, но PyTorch способен автоматически вычислять производную практически любой функции, причем очень быстро. В большинстве случаев это будет не дольше чем любая производная функция, которую вы можете написать сами. Вот пример.

Сначала выберем значение тензора, для которого нужны градиенты:

```
xt = tensor(3.).requires_grad_()
```

Обратите внимание на особый метод `requires_grad_`? Это магическое заклинание, с помощью которого мы сообщаем PyTorch, что хотим вычислить градиенты в отношении этой переменной с этим значением. По сути, мы прикрепляем к переменной тег, чтобы PyTorch в дальнейшем понимал, как вычислять градиенты других непосредственных вычислений над этой переменной.



СЛОВО АЛЕКСИСУ

Этот API может несколько запутать, если вы пришли из математики или физики, потому что в их контекстах «градиент» функции — это просто еще одна функция (то есть производная), поэтому вы можете ожидать, что связанные с градиентом API будут возвращать новую функцию. Но в глубоком обучении «градиент» обычно является *значением* производной функции при определенном значении аргумента. PyTorch API также делает ударение на аргумент, а не на функцию, чьи градиенты вы фактически вычисляете. Поначалу это может показаться шагом назад, но это просто другой подход.

Теперь мы вычислим функцию с этим значением. Обратите внимание, что PyTorch выводит не только вычисленное значение, но также и примечание, указывающее, что у него есть функция градиента, которую при необходимости он использует для вычисления градиентов:

```
yt = f(xt)
yt
tensor(9., grad_fn=<PowBackward0>)
```

В завершение мы попросим PyTorch вычислить градиенты:

```
yt.backward()
```

В данном случае `backward` означает *обратное распространение ошибки*. Этот термин относится к процессу вычисления производной каждого слоя. Как это конкретно происходит, мы увидим в главе 17, когда будем вычислять градиенты глубокой нейронной сети с самого начала. Этот процесс называется *обратным проходом*, и он противоположен *прямому проходу*, при котором вычисляются активации. Было бы несколько проще, если бы `backward` назывался просто

`calculate_grad`, но специалисты по глубокому обучению очень уж любят вносить особые термины по всякому поводу.

Теперь можно просмотреть градиенты, проверив атрибут `grad` тензора:

```
xt.grad
tensor(6.)
```

Если вы помните математику в старшей школе, то производной от x^{**2} будет $2*x$. Мы имеем $x=3$, следовательно, градиент будет $2*3=6$, что PyTorch вычислит самостоятельно.

Теперь повторим предыдущие шаги, но уже используем для функции аргумент вектора:

```
xt = tensor([3.,4.,10.]).requires_grad_()
xt
tensor([ 3.,  4., 10.], requires_grad=True)
```

Мы также добавим в нашу функцию `sum`, чтобы она могла принять вектор (то есть тензор ранга 1) и вернуть скаляр (то есть тензор ранга 0):

```
def f(x): return (x**2).sum()
yt = f(xt)
yt
tensor(125., grad_fn=<SumBackward0>)
```

Наши градиенты будут равны $2*x_t$, как мы и ожидали!

```
yt.backward()
xt.grad
tensor([ 6.,  8., 20.])
```

Эти градиенты сообщают только наклон функции, но не говорят, на какую величину нужно скорректировать параметры. Тем не менее они дают нам примерное представление: очень большой наклон может означать, что нужно произвести дополнительные корректировки, если же наклон мал, это может говорить о приближении к оптимальному значению.

Определение шагов скорости обучения

Решение о том, как изменять параметры на основе значений градиентов, является важной частью процесса глубокого обучения. Почти все подходы начинаются с простой идеи умножения градиента на некоторое небольшое значение, называемое *скоростью обучения* (LR). Как правило, эта скорость представлена числом между

0,001 и 0,1, хотя можно использовать любое. Нередко люди выбирают для него значение, просто перепробовав несколько разных и отыскав то, которое приводит к наилучшему результату обучения модели (чуть позже мы покажем вам подход получше, который называется *поиском скорости обучения*). После выбора скорости обучения можно скорректировать параметры с помощью простой функции:

```
w -= w.grad * lr
```

Это называется *пошаговым изменением* параметров посредством *шага оптимизации*.

Если вы выберете слишком медленную скорость обучения, то это может привести к необходимости выполнения большого количества шагов (рис. 4.2).

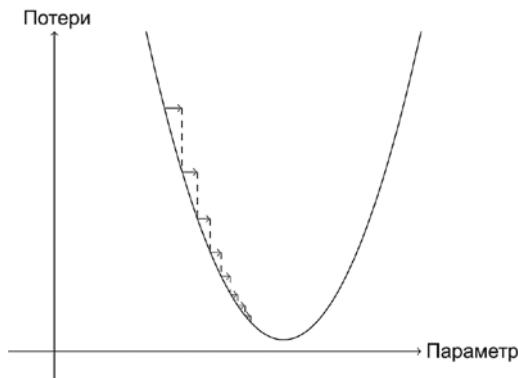


Рис. 4.2. Градиентный спуск с низкой LR

При этом выбор слишком большой скорости обучения потенциально еще хуже, так как может привести к ухудшению показателя *потерь* (рис. 4.3).

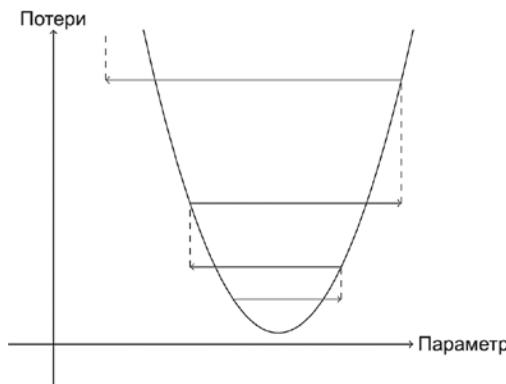


Рис. 4.3. Градиентный спуск с высокой LR

При слишком высокой скорости показатель потерь также может «перескакивать» туда-сюда вместо того, чтобы плавно отклоняться. На рис. 4.4 показано, как это приводит к необходимости выполнения множества шагов для успешного обучения.

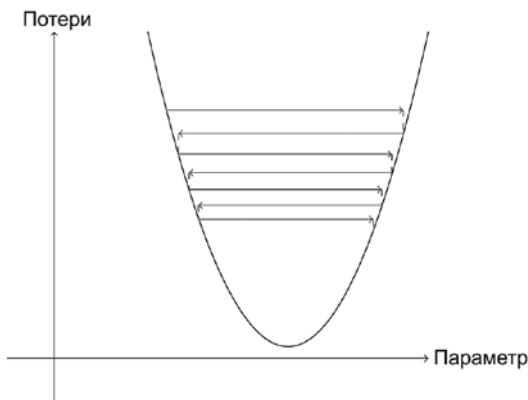


Рис. 4.4. Градиентный спуск с прыгающей LR

Используем все это в сквозном примере SGD.

Сквозной пример SGD

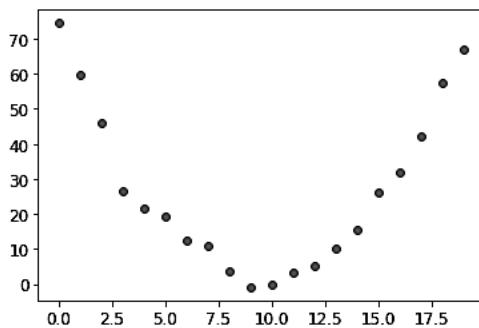
Мы видели, как использовать градиенты для минимизации потерь. Пора взглянуть на пример SGD и разобраться, как нахождение минимума можно использовать для обучения модели, повышая тем самым степень ее соответствия данным.

Начнем с примера простой синтетической модели. Представьте, что вы измеряете скорость вагонетки на американских горках, пересекающую верхнюю точку подъема. Вначале ее скорость будет высокой, но затем при подъеме вверх начнет снижаться, достигнув минимума в верхней точке. После пересечения этой точки скорость снова начнет возрастать, так как вагонетка поедет вниз. Нужно построить модель изменения скорости с течением времени. Выглядеть это может так:

```
time = torch.arange(0,20).float(); time

tensor([ 0.,  1.,  2.,  3.,  4.,  5.,  6.,  7.,  8.,  9.,  10.,  11.,  12.,  13.,
        14.,  15.,  16.,  17.,  18.,  19.])

speed = torch.randn(20)*3 + 0.75*(time-9.5)**2 + 1
plt.scatter(time,speed);
```



Мы добавили немного случайного шума, поскольку ручное измерение недостаточно точно. Это означает, что ответить на вопрос «Какая скорость была у вагонетки?» будет нелегко. Используя SGD, мы можем попытаться найти функцию, соответствующую нашим наблюдениям. Мы не можем рассматривать каждый вариант функции, поэтому предположим, что она квадратичная, то есть имеет форму $a*(time^{**2}) + (b*time) + c$.

Нам нужно провести четкое разграничение входных данных этой функции (время, когда мы измеряем скорость вагонетки) и ее параметров (значениями, определяющими, какую квадратичную функцию мы пробуем). Итак, соберем параметры в один аргумент, разделив тем самым ввод, t , и параметры, $params$, в сигнатуре функции:

```
def f(t, params):
    a,b,c = params
    return a*(t**2) + (b*t) + c
```

Другими словами, мы ограничили задачу нахождения наилучшей подходящей к данным функции до поиска именно ее *квадратичного* варианта. Это значительно упрощает задачу, поскольку каждая квадратичная функция полностью определяется всего тремя параметрами: a , b и c . Поэтому для нахождения наилучшей функции нам нужно найти только наилучшие значения для a , b и c .

Если мы можем решить эту задачу для перечисленных параметров, то сможем применить этот же подход и к другим, уже более сложным функциям с большим числом параметров, таким как нейронная сеть. Сначала найдем параметры для f , а затем вернемся и сделаем то же самое для датасета MNIST с нейронной сетью.

Для начала нам нужно определить, что мы подразумеваем под «наилучшей». Определим же мы это путем выбора *функции потерь*, которая будет возвращать значение на основе прогноза и цели, где более низкие значения функции будут соответствовать «лучшим» прогнозам. Для непрерывных данных стандартно используется *среднеквадратичная ошибка (MSE)*:

```
def mse(preds, targets): return ((preds-targets)**2).mean()
```

Теперь пройдемся по всем семи шагам процесса.

Шаг 1: инициализация параметров

Сначала мы инициализируем параметры со случайными значениями и с помощью `requires_grad` сообщаем PyTorch, что хотим отслеживать их градиенты:

```
params = torch.randn(3).requires_grad_()
```

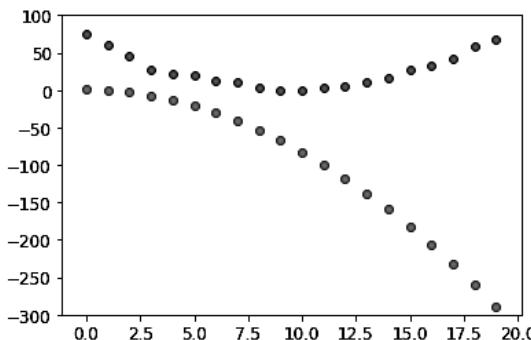
Шаг 2: вычисление прогнозов

Теперь мы вычисляем прогнозы:

```
preds = f(time, params)
```

Создадим небольшую функцию, чтобы посмотреть, насколько близки полученные прогнозы к целям:

```
def show_preds(preds, ax=None):
    if ax is None: ax=plt.subplots()[1]
    ax.scatter(time, speed)
    ax.scatter(time, to_np(preds), color='red')
    ax.set_ylim(-300,100)
show_preds(preds)
```



Не очень-то близко — наши случайные параметры предполагают, что вагонетка в итоге поедет назад, поскольку скорости получились отрицательные.

Шаг 3: вычисление потерь

Вычисление потерь производится так:

```
loss = mse(preds, speed)
loss

tensor(25823.8086, grad_fn=<MeanBackward0>)
```

Теперь наша цель — улучшить результаты, для чего нам понадобится знать градиенты.

Шаг 4: вычисление градиентов

Далее находим градиенты, или примерное значение того, насколько нужно изменить параметры:

```
loss.backward()  
params.grad  
  
tensor([-53195.8594, -3419.7146, -253.8908])  
  
params.grad * 1e-5  
  
tensor([-0.5320, -0.0342, -0.0025])
```

Эти градиенты можно использовать для улучшения параметров. Потребуется выбрать скорость обучения (в следующей главе мы рассмотрим, как это делать на практике, а сейчас просто возьмем $1e-5$ или 0,00001):

```
params  
  
tensor([-0.7658, -0.7506, 1.3525], requires_grad=True)
```

Шаг 5: корректировка весов

Теперь нужно обновить параметры на основе только что полученных градиентов:

```
lr = 1e-5  
params.data -= lr * params.grad.data  
params.grad = None
```



СЛОВО АЛЕКСИСУ

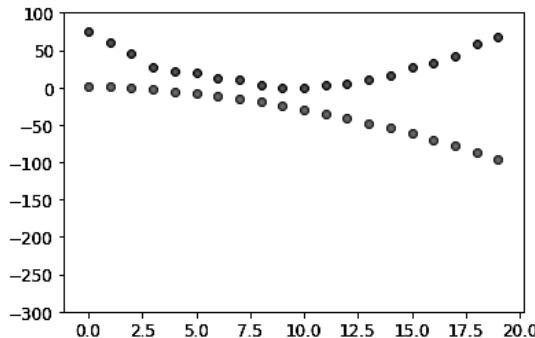
Понимание этой части процесса зависит от усвоения предыдущей информации. Для вычисления градиентов мы вызываем для `loss` метод `backward`. Но `loss` были вычислены с помощью `msc`, который следом получил в качестве ввода `preds`, а он был вычислен с помощью `f`, получившей на вводе `params`, являющиеся объектом, для которого мы изначально вызвали `required_grads` — первичный вызов, который теперь позволяет вызывать `backward` для `loss`. Эта цепочка вызовов функций представляет математическую композицию функций и позволяет PyTorch использовать дифференцирование сложной функции для вычисления этих градиентов.

Посмотрим, улучшились ли потери:

```
preds = f(time, params)  
mse(preds, speed)  
  
tensor(5435.5366, grad_fn=<MeanBackward0>)
```

А также рассмотрим график:

```
show_preds(preds)
```



Нужно повторить этот процесс несколько раз, чтобы создать функцию для выполнения одного шага:

```
def apply_step(params, prn=True):
    preds = f(time, params)
    loss = mse(preds, speed)
    loss.backward()
    params.data -= lr * params.grad.data
    params.grad = None
    if prn: print(loss.item())
    return preds
```

Шаг 6: повторение процесса

Теперь мы этот процесс повторяем. Совершая циклы и внося улучшения, мы надеемся получить хороший результат:

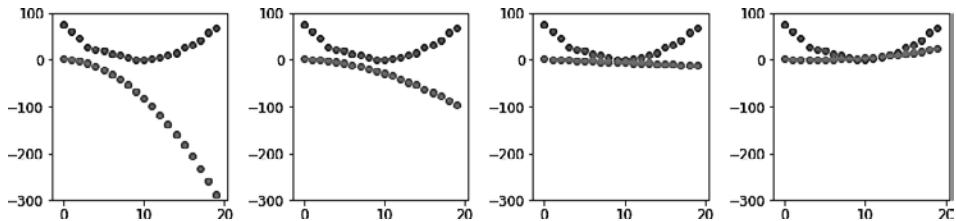
```
for i in range(10): apply_step(params)

5435.53662109375
1577.4495849609375
847.3780517578125
709.22265625
683.0757446289062
678.12451171875
677.1839599609375
677.0025024414062
676.96435546875
676.9537353515625
```

Потери снижаются, что нам и было нужно. Но рассмотрение только этих показателей потерь скрывает тот факт, что каждое повторение подразумевает проверку совершенно новой квадратичной функции на пути к нахождению наилучшей.

Можно представить этот процесс визуально, если вместо вывода функции потерь начертить график функции на каждом шаге. Тогда мы сможем увидеть, как ее форма изменяется в сторону наилучшей квадратичной функции для наших данных:

```
_ ,axs = plt.subplots(1,4,figsize=(12,3))
for ax in axs: show_preds(apply_step(params, False), ax)
plt.tight_layout()
```



Шаг 7: остановка

Мы случайно решили остановиться после десяти эпох обучения. Но как уже говорилось, на практике решение об остановке принимается на основе наблюдения за изменением метрик, а также за изменением потерь при обучении и контроле.

Подведение итогов темы градиентного спуска

Теперь, когда вы увидели, что происходит на каждом шаге, пора еще раз взглянуть на графическое представление процесса градиентного спуска (рис. 4.5) и вкратце его резюмировать.

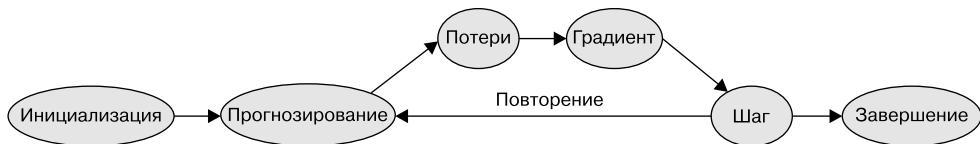


Рис. 4.5. Процесс градиентного спуска

В его начале веса модели могут быть случайными (*обучение с нуля*) или переходить из предварительно обученной модели (*перенос обучения*). В первом случае получаемый нами вывод не будет отражать то, что нам нужно, и даже во втором случае, скорее всего, предварительная модель не будет достаточно хороша в отношении конкретной поставленной нами задачи. Так что этой модели потребуется *выучить* более подходящие веса.

Начнем со сравнения выводов, производимых моделью, с нашими целями (у нас есть размеченные данные, поэтому нам известно, какой результат должна

выдавать модель) с помощью *функции потерь*, которая возвращает число, которое нужно сделать минимальным путем улучшения весов. Для этого мы берем несколько элементов данных (например, изображения) из обучающей выборки и передаем их модели. Далее мы сравниваем соответствующие цели с помощью функции потерь, при этом получаемый нами показатель говорит, насколько были ошибочны прогнозы. После этого мы слегка изменяем веса, чтобы этот показатель немного улучшить.

Чтобы выяснить, как изменить веса, чтобы несколько снизить потери, мы вычисляем *градиенты*. (На самом деле позволяем PyTorch это сделать за нас.) Рассмотрим аналогию. Представьте, что вы заблудились в горах, а машину припарковали у подножия в самой нижней точке. Чтобы найти к ней обратный путь, можно блуждать в произвольных направлениях, но это вряд ли особо поможет. Поскольку вы знаете, что ваш автомобиль находится в нижней точке, то будете спускаться вниз. Делая шаг за шагом в направлении самого крутого спуска горы, вы постепенно доберетесь до нужного места. Мы используем величину градиента (то есть крутизну наклона) для определения длины необходимого шага. В частности, мы умножаем градиент на выбранное число, характеризующее *скорость обучения*, чтобы решить, какой длины шаг нужно делать. После этого мы *повторяем* процесс до тех пор, пока не достигнем нижней точки, которой окажется место парковки автомобиля, где мы и сможем *остановиться*.

Все, что мы только что видели, можно перенести напрямую в датасет MNIST, за исключением функции потерь. Теперь давайте посмотрим, как можно определить хорошую цель обучения.

Функция потерь MNIST

У нас уже есть все x , представляющие зависимые переменные, то есть сами изображения. Мы конкатенируем их все в один тензор, а также преобразуем их из списка матриц (тензор ранга 3) в список векторов (тензор ранга 2). Это можно сделать с помощью PyTorch-метода `view`, который изменяет форму тензора, не меняя его содержимого. `-1` — это особый параметр для `view`, означающий, что нужно сделать ось такой длины, чтобы уместить все данные:

```
train_x = torch.cat([stacked_threes, stacked_sevens]).view(-1, 28*28)
```

Нам нужны метки для каждого изображения, и в качестве таких меток мы используем `1` для троек и `0` для семерок:

```
train_y = tensor([1]*len(threes) + [0]*len(sevens)).unsqueeze(1)
train_x.shape, train_y.shape
(torch.Size([12396, 784]), torch.Size([12396, 1]))
```

В PyTorch `Dataset` необходим для возврата (x, y) при индексировании. Python предоставляет функцию `zip`, которая при совмещении с `list` обеспечивает простой способ получения этой функциональности:

```
dset = list(zip(train_x,train_y))
x,y = dset[0]
x.shape,y

(torch.Size([784]), tensor([1]))
valid_x = torch.cat([valid_3_tens, valid_7_tens]).view(-1, 28*28)
valid_y = tensor([1]*len(valid_3_tens) + [0]*len(valid_7_tens)).unsqueeze(1)
valid_dset = list(zip(valid_x,valid_y))
```

Теперь нам нужен (изначально произвольный) вес для каждого пикселя (это шаг *инициализации*):

```
def init_params(size, std=1.0): return (torch.randn(size)*std).requires_grad_()
weights = init_params((28*28,1))
```

Функция `weights*pixels` не будет достаточно гибкой, так как во всех случаях, когда пиксели равны 0, она тоже равна 0 (то есть ее отрезок равен 0). Вы можете припомнить из школьной программы по математике, что формулой для строки является $y=w*x+b$. Нам по-прежнему нужна переменная b , которую мы также инициализируем со случайным числом:

```
bias = init_params(1)
```

В нейронных сетях часть w уравнения называется *весами*, а b называется *смещением*. Вместе веса и смещение составляют *параметры*.



ТЕРМИН: ПАРАМЕТРЫ

Это веса и смещения модели. В уравнении $w*x+b$ веса представлены как w , а смещения как b .

Теперь можно вычислить прогноз для одного изображения:

```
(train_x[0]*weights.T).sum() + bias
tensor([20.2336], grad_fn=<AddBackward0>)
```

Несмотря на возможность использовать цикл `for` для вычисления прогнозов в отношении каждого изображения, это будет очень медленно. Поскольку циклы Python не выполняются на GPU и сам этот язык достаточно медленный, особенно при работе с циклами, нам нужно представить как можно большую часть вычислений модели, используя функции высшего порядка.

Для этого случая есть очень удобная математическая операция, которая вычисляет w^*x для каждой строки матрицы и называется *матричным умножением* (рис. 4.6).

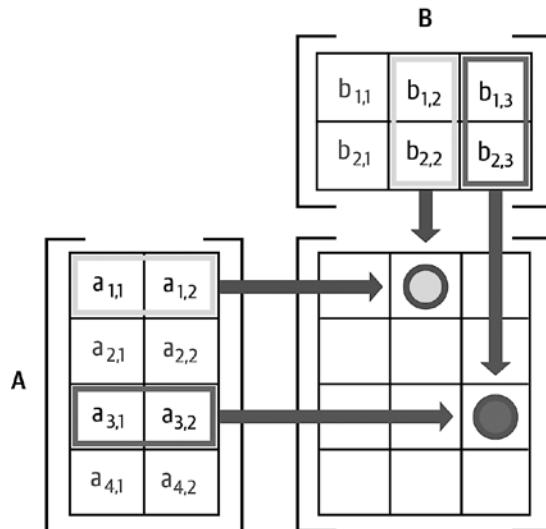


Рис. 4.6. Матричное умножение

На этом изображении две перемноженные матрицы, А и В. Каждый элемент результата их перемножения, который мы будем называть АВ, содержит каждый элемент соответствующей строки А, умноженный на каждый элемент соответствующего ему столбца В, сложенные вместе. Например, строка 1, столбец 2 вычисляется как $a_{1,1} * b_{1,2} + a_{1,2} * b_{2,2}$. Если вам нужно освежить в памяти тему матричного умножения, то предлагаем заглянуть в раздел *Intro to Matrix Multiplication* (<https://oreil.ly/w0XKS>) («Введение в матричное умножение») Академии Хана, поскольку это самая важная математическая операция в глубоком обучении.

В Python матричное умножение представлено оператором @. Давайте его испробуем:

```
def linear1(xb): return xb@weights + bias
preds = linear1(train_x)
preds

tensor([[20.2336,
         [17.0644],
         [15.2384],
         ...,
         [18.3804],
         [23.8567],
         [28.6816]], grad_fn=<AddBackward0>)
```

Первый элемент, как и ожидалось, совпадает с произведенным нами ранее вычислением. Это уравнение, `batch @ weights + bias`, является одним из двух фундаментальных уравнений любой нейронной сети. (Второе — это *функция активации*, с которой мы очень скоро познакомимся.)

Проверим точность. Чтобы определить, представляет вывод 3 или 7, можно просто проверить, превышает ли он 0, а значит, точность для каждого элемента можно вычислить (с помощью бродкастинга, без применения циклов) так:

```
corrects = (preds>0.0).float() == train_y
corrects

tensor([[ True],
       [ True],
       [ True],
       ...,
       [False],
       [False],
       [False]])]

corrects.float().mean().item()
0.4912068545818329
```

Посмотрим, как изменяется точность при небольшом изменении одного из весов:

```
weights[0] *= 1.0001

preds = linear1(train_x)
((preds>0.0).float() == train_y).float().mean().item()

0.4912068545818329
```

Как мы уже видели, для улучшения модели с помощью SGD нужны градиенты, а для вычисления градиентов нужна *функция потерь*, показывающая, насколько хороша наша модель. Это обусловлено тем, что градиенты являются мерой изменения функции потерь при небольших корректировках весов.

Итак, нам нужно выбрать функцию потерь. Очевидным решением будет использовать точность, являющуюся нашей метрикой, еще и в качестве функции потерь. В данном случае мы вычислим прогнозы для каждого изображения и соберем полученные значения для вычисления общей точности, после чего вычислим градиенты каждого веса в отношении этой общей точности.

К сожалению, у нас есть серьезная техническая проблема. Градиент функции — это ее *наклон*, или крутизна, которую можно определить как отношение приращения функции к приращению аргумента, то есть как значение увеличения или уменьшения функции, разделенного на величину изменения нами ввода. Математически это можно написать так:

$$(y_{\text{new}} - y_{\text{old}}) / (x_{\text{new}} - x_{\text{old}})$$

В результате получается хорошее приблизительное значение градиента, когда x_{new} близок к x_{old} , то есть разница между ними очень мала. Однако точность изменяется сильно, только когда прогноз меняется с 3 на 7 или наоборот. Проблема в том, что небольшое изменение весов из x_{old} в x_{new} не будет приводить к изменению прогноза, поэтому $(y_{\text{new}} - y_{\text{old}})$ почти всегда будет 0. Иначе говоря, градиент почти везде будет 0.

Очень малое изменение значения веса часто не будет изменять точность совсем. Это означает, что использование точности в качестве функции потерь не имеет смысла — если мы так поступаем, то в большинстве случаев градиенты будут 0, и модель на основе этого числа обучаться не сможет.



СЛОВО СИЛЬВЕЙНУ

В математическом смысле точность является функцией, которая почти везде постоянна (за исключением порога 0,5), поэтому ее производная почти везде ноль (и бесконечна при пороговом значении). В итоге она дает градиенты, равные нулю или бесконечности, что не имеет пользы для обучения модели.

Нужна такая функция потерь, которая при небольшом улучшении прогнозов на основе весов будет давать немного лучшие потери. Как же именно выглядит «немного улучшенный прогноз»? В данном случае это означает, что если верный ответ 3, то показатель немного увеличивается, а если 7, то немного уменьшается.

Напишем такую функцию. Какую она примет форму?

Функция потерь получает не сами изображения, а прогнозы модели. Поэтому создадим один аргумент, `prds`, со значениями между 0 и 1, где каждое из значений будет прогнозом того, что на изображении наблюдается 3. Это вектор (то есть тензор ранга 1), индексированный в изображениях.

Целью функции потерь является измерение разницы между прогнозируемыми значениями и истинными, которые представлены целями (метками). Следовательно, нужно создать еще один аргумент, `trgts`, со значениями 0 или 1, который будет говорить, является изображение тройкой или нет. Он также будет вектором (то есть еще одним тензором ранга 1), индексированным по изображениям.

В качестве примера предположим, что у нас было три известных нам изображения 3, 7 и 3. Далее предположим, что наша модель с высокой степенью уверенности (0,9) спрогнозировала первое как 3, с гораздо меньшей уверенностью (0,4) второе как 7 и с удовлетворительной уверенностью (0,2), но ошибочно, что последнее является 7. Это бы означало, что наша функция потерь получит данные значения в качестве ввода:

```
trgts = tensor([1,0,1])
prds    = tensor([0.9, 0.4, 0.2])
```

Вот первая попытка функции потерь, измеряющей удаленность `predictions` от `targets`:

```
def mnist_loss(predictions, targets):
    return torch.where(targets==1, 1-predictions, predictions).mean()
```

Мы используем новую функцию, `torch.where(a,b,c)`. Это аналогично применению генератора списков `[b[i] if a[i] else c[i] for i in range(len(a))]` с поправкой на то, что она работает с тензорами на скорости C/CUDA. Говоря простым языком, эта функция будет измерять, насколько удален каждый прогноз от 1, если это должна быть 1, и насколько удален от 0, если это должен быть 0. После этого она будет выводить среднее значение всех этих удаленостей.



ЧИТАЙТЕ ДОКУМЕНТАЦИЮ

Очень важно изучить функции PyTorch, подобные этой, потому что в Python перебор тензоров циклом выполняется на скорости Python, а не C/CUDA. Попробуйте сейчас выполнить команду `help(torch.where)`, чтобы прочесть документацию по этой функции, или, что еще лучше, найдите соответствующую информацию на сайте документации PyTorch.

Проверим ее работоспособность в отношении `prds` и `trgts`:

```
torch.where(trgts==1, 1-prds, prds)
tensor([0.1000, 0.4000, 0.8000])
```

Эта функция возвращает меньшее число, когда прогнозы более точны, когда точные прогнозы более уверены (высокие абсолютные значения) и когда неточные прогнозы сделаны менее уверенно. В PyTorch мы всегда предполагаем, что чем ниже значение функции потерь, тем лучше. Поскольку для получения итоговых потерь нам необходим скаляр, `mnist_loss` получает среднее значение *предыдущего* тензора:

```
mnist_loss(prds,trgts)
tensor(0.4333)
```

Если мы, например, изменим наш прогноз для одной «ложной» цели с 0,2 до 0,8, то потери снижаются, указывая, что такой прогноз более точен.

```
mnist_loss(tensor([0.9, 0.4, 0.8]),trgts)
tensor(0.2333)
```

Одна из проблем с `mnist_loss` в ее текущем виде состоит в том, что в качестве прогнозов она всегда предполагает значения от 0 до 1. В результате нам нужно убедиться, что это именно так, для чего как раз и есть другой тип функции, с которой мы сейчас познакомимся.

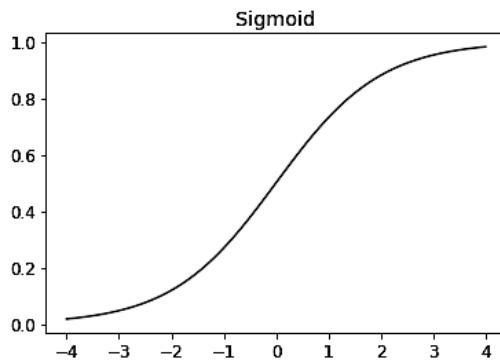
Сигмоида

Функция `sigmoid` всегда выводит число от 0 до 1 и определяется следующим образом:

```
def sigmoid(x): return 1/(1+torch.exp(-x))
```

PyTorch определяет за нас ее ускоренную версию, так что нам для этого ничего делать не нужно. Это очень важная функция в сфере глубокого обучения, поскольку нам часто нужно, чтобы возвращаемые значения были между 0 и 1. Вот как она выглядит:

```
plot_function(torch.sigmoid, title='Sigmoid', min=-4, max=4)
```



Здесь мы видим, что она получает на вводе положительное или отрицательное значение и выводит его сжатую форму в диапазоне от 0 до 1. Кроме того, она является гладкой кривой, которая идет только вверх, что упрощает SGD поиск значимых градиентов.

Давайте обновим `mnist_loss`, чтобы изначально применять ко вводу `sigmoid`:

```
def mnist_loss(predictions, targets):
    predictions = predictions.sigmoid()
    return torch.where(targets==1, 1-predictions, predictions).mean()
```

Теперь можно быть уверенным, что функция потерь заработает, даже если прогнозы не будут вписываться в диапазон от 0 до 1. Требуется только соответствие более высоких прогнозов более высокой степени уверенности.

Определив функцию потерь, нeliшним будет вспомнить, почему мы это вообще сделали. В конце концов, у нас уже была метрика, показывающая общую точность, — зачем определять еще и потери?

Ключевое отличие в том, что метрика служит для человеческого понимания, а потери — для автоматизации обучения. Чтобы запустить этот процесс автома-

тизированного обучения, потери должны быть функцией, имеющей значимую производную. Она не может содержать большие плоские участки и высокие подъемы, но должна быть объективно гладкой. Именно поэтому мы создали функцию потерь, которая будет реагировать на небольшие изменения в уровне уверенности. Это требование означает, что иногда она не будет отражать именно то, к чему мы стремимся, но будет являться компромиссом между нашей целью и функцией, которая может оптимизироваться на основе своего градиента. Эта функция потерь вычисляется для каждого элемента датасета, а затем в конце каждой эпохи все значения потерь усредняются и выдается их итоговое среднее значение.

Метрики же, с другой стороны, представляют собой значимые для нас числа, которые выводятся в конце каждой эпохи, показывая результативность модели. Нам же при оценке качества модели важно научиться обращать внимание именно на метрики, а не на потери.

SGD и мини-пакеты

Теперь, когда у нас есть подходящая для SGD функция, можно рассмотреть подробности следующей фазы процесса обучения, которые относятся к изменению или обновлению весов на основе градиентов. Это будет *шаг оптимизации*.

Для этого нужно вычислить потери для одного или нескольких элементов данных. Сколько же именно нужно? Мы можем либо вычислить потери для всех элементов датасета и взять среднее значение, либо вычислить их только для одного элемента. Тем не менее ни один из этих способов не идеален. Вычисление потерь для всего датасета займет много времени, а при вычислении для одного элемента будет использовано недостаточно информации, что приведет к получению неточного и нестабильного градиента. У вас бы возникли сложности с обновлением весов, если бы вы учитывали только то, насколько это улучшит производительность модели в отношении всего одного элемента.

Поэтому мы идем на компромисс, а именно вычисляем средние потери для нескольких элементов данных за один раз. Такой набор элементов называется *мини-пакетом*. Количество присутствующих в мини-пакете элементов называется *размером пакета*. Чем больше будет этот размер, тем точнее и стабильнее будет оценка градиентов вашего датасета, но для этого потребуется больше времени и вы будете обрабатывать меньшее число мини-пакетов в эпоху. Выбор оптимального размера пакета относится к решениям, которые вы как практик глубокого обучения будете принимать для быстрого и точного обучения модели. В книге мы будем говорить, как правильно в этой ситуации делать выбор.

Еще одна весомая причина использовать мини-пакеты, а не вычислять градиент для отдельных единиц данных в том, что на практике мы всегда производим обучение на ускорителях, а именно GPU. Эти ускорители работают хорошо,

только когда им единовременно дается большой объем работы, поэтому так мы и стараемся делать. Передача мини-пакетов как раз относится к одному из способов это осуществить. Но если вы все же перегрузите их данными, то у них может просто не хватить памяти, так что выбор оптимального количества тоже не так прост!

Как вы видели в процессе рассмотрения аугментации данных в главе 2, мы добиваемся улучшенной генерализации (обобщения), если можем варьировать элементы в процессе обучения. Одним из простых и одновременно эффективных способов варьирования будет решить, какие данные и в какой пакет добавлять. Вместо того чтобы просто упорядоченно нумеровать датасет для каждой эпохи, мы обычно случайным образом его перемешиваем и уже потом составляем мини-пакеты. PyTorch и fastai предоставляют класс, который будет перемешивать данные и собирать их в пакеты за вас. Этот класс называется `DataLoader`.

`DataLoader` может принимать любую коллекцию Python и превращать ее в итератор, разбивающий данные на несколько пакетов:

```
coll = range(15)
d1 = DataLoader(coll, batch_size=5, shuffle=True)
list(d1)

[tensor([ 3, 12,  8, 10,  2]),
 tensor([ 9,  4,  7, 14,  5]),
 tensor([ 1, 13,  0,  6, 11])]
```

Для обучения модели нужна не просто какая-то коллекция Python, а коллекция с независимыми и зависимыми переменными (входными данными и целями модели). Коллекция, содержащая кортежи независимых и зависимых переменных, в PyTorch называется `Dataset`. Вот ее очень простой пример:

```
ds = L(enumerate(string.ascii_lowercase))
ds

(#26) [(0, 'a'), (1, 'b'), (2, 'c'), (3, 'd'), (4, 'e'), (5, 'f'), (6, 'g'), (7, 'h'), (8, 'i'), (9, 'j')...]
```

Когда мы передаем `Dataset` в `DataLoader`, обратно получаем много пакетов, которые сами являются кортежами тензоров, представляющих пакеты независимых и зависимых переменных:

```
d1 = DataLoader(ds, batch_size=6, shuffle=True)
list(d1)

[([tensor([17, 18, 10, 22,  8, 14]), ('r', 's', 'k', 'w', 'i', 'o')),
 (tensor([20, 15,  9, 13, 21, 12]), ('u', 'p', 'j', 'n', 'v', 'm')),
 (tensor([ 7, 25,  6,  5, 11, 23]), ('h', 'z', 'g', 'f', 'l', 'x')),
 (tensor([ 1,  3,  0, 24, 19, 16]), ('b', 'd', 'a', 'y', 't', 'q')),
 (tensor([2, 4]), ('c', 'e'))]
```

Теперь мы готовы к написанию нашего первого обучающего цикла для модели с помощью SGD!

Собрать все вместе

Пора реализовать процесс, который мы наблюдали на рис. 4.1. В коде этот процесс будет реализован для каждой эпохи следующим образом:

```
for x,y in dl:
    pred = model(x)
    loss = loss_func(pred, y)
    loss.backward()
    parameters -= parameters.grad * lr
```

Сначала заново инициализируем параметры:

```
weights = init_params((28*28,1))
bias = init_params(1)
```

`DataLoader` можно создать из `Dataset`:

```
dl = DataLoader(dset, batch_size=256)
xb,yb = first(dl)
xb.shape,yb.shape
(torch.Size([256, 784]), torch.Size([256, 1]))
```

То же самое мы делаем для контрольной выборки:

```
valid_dl = DataLoader(valid_dset, batch_size=256)
```

Теперь создадим для тестирования мини-пакет с размерностью 4:

```
batch = train_x[:4]
batch.shape
torch.Size([4, 784])

preds = linear1(batch)
preds

tensor([[-11.1002],
       [ 5.9263],
       [ 9.9627],
       [-8.1484]], grad_fn=<AddBackward0>)

loss = mnist_loss(preds, train_y[:4])
loss

tensor(0.5006, grad_fn=<MeanBackward0>)
```

Теперь можно вычислить градиенты:

```
loss.backward()
weights.grad.shape, weights.grad.mean(), bias.grad

(torch.Size([784, 1]), tensor(-0.0001), tensor([-0.0008]))
```

Поместим это в функцию:

```
def calc_grad(xb, yb, model):
    preds = model(xb)
    loss = mnist_loss(preds, yb)
    loss.backward()
```

И протестируем:

```
calc_grad(batch, train_y[:4], linear1)
weights.grad.mean(), bias.grad

(tensor(-0.0002), tensor([-0.0015]))
```

Но смотрите, что происходит, когда мы вызываем ее дважды:

```
calc_grad(batch, train_y[:4], linear1)
weights.grad.mean(), bias.grad

(tensor(-0.0003), tensor([-0.0023]))
```

Градиенты изменились! Причина в том, что `loss.backward` *прибавляет* градиенты `loss` к любым сохраненным градиентам. Значит, сначала нужно установить текущие градиенты в ноль:

```
weights.grad.zero_()
bias.grad.zero_();
```



ОПЕРАЦИИ НА МЕСТЕ

Методы PyTorch, чьи имена заканчиваются нижним подчеркиванием, изменяют свои объекты на месте. Например, `bias.zero_` устанавливает для всех элементов тензора `bias` значение 0.

Остался всего один шаг, а именно обновление весов и смещений на основе градиента и скорости обучения. Здесь дадим команду PyTorch не принимать градиенты этого шага — в противном случае возникнет путаница, когда мы попытаемся вычислить производную для следующего пакета. Если мы присвоим к `data` атрибут тензора, PyTorch не будет принимать градиенты этого шага. Вот как в результате выглядит наш простой обучающий цикл для эпохи:

```
def train_epoch(model, lr, params):
    for xb,yb in dl:
```

```
calc_grad(xb, yb, model)
for p in params:
    p.data -= p.grad*lr
    p.grad.zero_()
```

Нужно также проверить результаты, посмотрев на точность в отношении контрольной выборки. Чтобы решить, является вывод 3 или 7, мы проверяем, превышает ли точность 0,5. Вычислить ее показатель для каждого элемента (с помощью транслирования, не циклов) можно так:

```
(preds>0.5).float() == train_y[:4]

tensor([[False],
        [ True],
        [ True],
        [False]])
```

Так мы получаем функцию для вычисления контрольной точности:

```
def batch_accuracy(xb, yb):
    preds = xb.sigmoid()
    correct = (preds>0.5) == yb
    return correct.float().mean()
```

Убедиться, что она работает, можно так:

```
batch_accuracy(linear1(batch), train_y[:4])

tensor(0.5000)
```

А затем объединить пакеты:

```
def validate_epoch(model):
    accs = [batch_accuracy(model(xb), yb) for xb,yb in valid_dl]
    return round(torch.stack(accs).mean().item(), 4)

validate_epoch(linear1)
```

0.5219

Это стартовая точка. Выполним одну эпоху обучения и посмотрим, улучшилась ли точность:

```
lr = 1.
params = weights,bias
train_epoch(linear1, lr, params)
validate_epoch(linear1)
```

0.6883

А затем проведем еще несколько:

```

for i in range(20):
    train_epoch(linear1, lr, params)
    print(validate_epoch(linear1), end=' ')
0.8314 0.9017 0.9227 0.9349 0.9438 0.9501 0.9535 0.9564 0.9594 0.9618 0.9613
0.9638 0.9643 0.9652 0.9662 0.9677 0.9687 0.9691 0.9691 0.9696

```

Неплохо! Мы уже почти достигли той же точности, что и в подходе «пиксельного сходства», создав при этом стандартную основу, на которой сможем продолжать построение. Следующим шагом будет создание объекта, обрабатывающего шаг SGD. В PyTorch он называется *оптимизатором*.

Создание оптимизатора

Поскольку это стандартная основа, PyTorch предоставляет ряд полезных классов, упрощая ее реализацию. Первое, что мы можем сделать, — это заменить функцию `learner` на PyTorch-модуль `nn.Linear`. *Модуль* — это объект класса, который наследует от класса `nn.Module`. Поведение объектов этого класса совпадает с поведением стандартных функций Python в том, что вы можете вызывать их с помощью скобок, на что они будут возвращать активации модели.

`nn.Linear` делает то же, что и `init_params` вместе с `linear`. Он содержит как *веса*, так и *смещения* в одном классе. Вот как можно воспроизвести нашу модель из предыдущего раздела:

```
linear_model = nn.Linear(28*28,1)
```

Каждый модуль PyTorch знает, какие из имеющихся у него параметров могут быть обучены. Они доступны с помощью метода `parameters`:

```
w,b = linear_model.parameters()
w.shape,b.shape

(torch.Size([1, 784]), torch.Size([1]))
```

Эту информацию можно использовать для создания оптимизатора:

```

class BasicOptim:
    def __init__(self,params,lr): self.params,self.lr = list(params),lr

    def step(self, *args, **kwargs):
        for p in self.params: p.data -= p.grad.data * self.lr

    def zero_grad(self, *args, **kwargs):
        for p in self.params: p.grad = None

```

Мы можем создать оптимизатор, передав в него параметры модели:

```
opt = BasicOptim(linear_model.parameters(), lr)
```

Теперь цикл обучения можно упростить:

```
def train_epoch(model):
    for xb,yb in dl:
        calc_grad(xb, yb, model)
        opt.step()
        opt.zero_grad()
```

При этом контрольную функцию менять не нужно:

```
validate_epoch(linear_model)
0.4157
```

Для упрощения поместим наш короткий обучающий цикл в функцию:

```
def train_model(model, epochs):
    for i in range(epochs):
        train_epoch(model)
        print(validate_epoch(model), end=' ')
```

Полученные результаты совпадают с результатами предыдущего раздела:

```
train_model(linear_model, 20)

0.4932 0.8618 0.8203 0.9102 0.9331 0.9468 0.9555 0.9629 0.9658 0.9673 0.9687
> 0.9707 0.9726 0.9751 0.9761 0.9761 0.9775 0.978 0.9785 0.9785
```

fastai предоставляет класс SGD, который по умолчанию делает то же, что и BasicOptim:

```
linear_model = nn.Linear(28*28,1)
opt = SGD(linear_model.parameters(), lr)
train_model(linear_model, 20)

0.4932 0.852 0.8335 0.9116 0.9326 0.9473 0.9555 0.9624 0.9648 0.9668 0.9692
> 0.9712 0.9731 0.9746 0.9761 0.9765 0.9775 0.978 0.9785 0.9785
```

Кроме того, fastai предоставляет Learner.fit, который можно использовать вместо train_model. Для создания Learner сначала нужно создать DataLoaders путем передачи в него обучающего и контрольного DataLoader:

```
dls = DataLoaders(dl, valid_dl)
```

Чтобы создать Learner без использования приложения (такого, как cnn_Learner), нам нужно передать ему все созданные нами в этой главе элементы: DataLoaders, модель, функцию оптимизации (которой будут переданы параметры), функцию потерь и при желании нужные метрики для вывода:

```
learn = Learner(dls, nn.Linear(28*28,1), opt_func=SGD,
                 loss_func=mnist_loss, metrics=batch_accuracy)
```

Теперь можно вызвать `fit`:

```
learn.fit(10, lr=lr)
```

epoch	train_loss	valid_loss	batch_accuracy	time
0	0.636857	0.503549	0.495584	00:00
1	0.545725	0.170281	0.866045	00:00
2	0.199223	0.184893	0.831207	00:00
3	0.086580	0.107836	0.911187	00:00
4	0.045185	0.078481	0.932777	00:00
5	0.029108	0.062792	0.946516	00:00
6	0.022560	0.053017	0.955348	00:00
7	0.019687	0.046500	0.962218	00:00
8	0.018252	0.041929	0.965162	00:00
9	0.017402	0.038573	0.967615	00:00

Мы видим, что в классах PyTorch и fastai нет ничего волшебного. Это просто удобные предварительно упакованные компоненты, немного упрощающие жизнь. (Кроме того, они предоставляют много дополнительной функциональности, которая будет использоваться в последующих главах.)

Теперь с помощью этих классов можно заменить линейную модель на нейронную сеть.

Добавление нелинейности

Пока есть общая процедура для оптимизации параметров функции, которую мы опробовали на нашей скучной функции: простом линейном классификаторе, который ограничен в своих возможностях. Чтобы несколько усложнить его (и получить возможность обрабатывать больше задач), нам нужно добавить между двумя линейными классификаторами что-то нелинейное (то есть отличное от $ax + b$) — так мы получим нейронную сеть.

Вот полное определение простой нейронной сети:

```
def simple_net(xb):
    res = xb@w1 + b1
    res = res.max(tensor(0.0))
    res = res@w2 + b2
    return res
```

Вот и все! Здесь в `simple_net` присутствуют только два линейных классификатора с функцией `max` между ними.

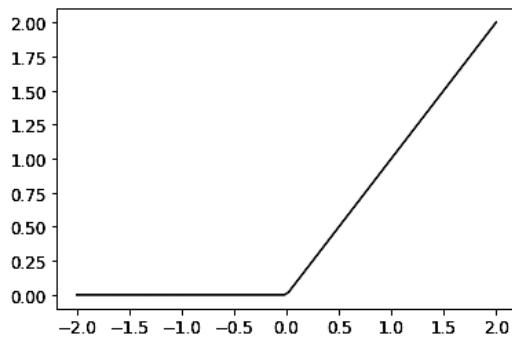
`w1` и `w2` являются тензорами весов, а `b1` и `b2` — тензорами смещений. Это параметры, которые изначально инициализируются случайно, точно так же, как мы делали в предыдущем разделе:

```
w1 = init_params((28*28,30))
b1 = init_params(30)
w2 = init_params((30,1))
b2 = init_params(1)
```

Ключевой момент в том, что `w1` имеет 30 выходных активаций (`w2`, соответственно, должен иметь 30 входных активаций). Это означает, что первый слой может построить 30 различных признаков, каждый из которых будет представлять разный набор пикселей. Вы можете изменить 30 на любое нужное вам значение, тем самым усложнив или упростив модель.

Эта небольшая функция `res.max(tensor(0.0))` называется *блоком линейной ректификации*, также известным как *ReLU*. Думаем, вы согласитесь с тем, что выражение «блок линейной ректификации» звучит слишком мудрено и сложно... Но в действительности это всего лишь `res.max(tensor(0.0))`, где мы замещаем любые отрицательные значения нулем. Эта небольшая функция также доступна в PyTorch в виде `F.relu`:

```
plot_function(F.relu)
```



Основная идея заключается в том, что при использовании большего числа слоев мы можем поручить нашей модели увеличенный объем вычислений, что позволит нам моделировать более сложные функции. Но при этом нет смысла просто добавлять один линейный макет за другим, потому что перемножение компонентов с последующим их многократным сложением можно заменить все тем же перемножением, но уже с однократным сложением. Иными словами, последовательность любого количества линейных слоев можно заменить одним линейным слоем с другим набором параметров.



СЛОВО ДЖЕРЕМИ

В глубоком обучении есть огромное количество терминов, включая «линейный модуль ректификации». При этом сложность большинства из них не превышает написания в одну строку, что мы только что и наблюдали. Правда же в том, что ученым для публикации их работ необходимо использовать максимально сложные и запутанные формулировки. Одним из способов этого является внедрение профессиональных терминов. К сожалению, из-за этого ML становится более отпугивающим и сложным, чем должно быть. Изучение терминологии необходимо, иначе научные работы и обучающие программы не дадут нужного понимания. Но это не значит, что терминология должна отпугивать. Когда вы встречаете неизвестный термин или выражение, то, скорее всего, они будут представлять очень простой принцип.

Но если мы поместим между слоями нелинейную функцию, такую как `max`, то так больше сделать не сможем. Теперь каждый линейный слой будет уже отделен от остальных, что позволит ему выполнять свою индивидуальную работу. Функция `max` особенно интересна тем, что оперирует аналогично простой инструкции `if`.



СЛОВО СИЛЬВЕЙНУ

Математически композиция из двух линейных функций является другой линейной функцией. Это значит, что можно собирать в стек любое нужное количество линейных классификаторов, которые в отсутствие нелинейных функций будут представлять один линейный классификатор.

Как ни странно, можно математически доказать, что эта небольшая функция способна решить любую вычисляемую задачу с произвольно высоким уровнем точности, если найти верные параметры для w_1 и w_2 , а также сделать эти матрицы достаточно большими. Любую произвольно волнистую функцию (*arbitrarily wiggly function*) мы можем аппроксимировать как набор объединенных строк. Чтобы приблизить ее к волнистой функции, нужно просто использовать более короткие строки. Это называется *теоремой универсальной аппроксимации*. Используемые здесь три строки называются *слоями*. Первая и третья представляют собой *линейные слои*, а вторая называется либо *нелинейной функцией*, либо *функцией активации*.

Как и в предыдущем разделе, воспользуемся преимуществами PyTorch, заменив этот код на более простой вариант:

```
simple_net = nn.Sequential(
    nn.Linear(28*28,30),
    nn.ReLU(),
    nn.Linear(30,1)
)
```

`nn.Sequential` создает модуль, который будет поочередно вызывать каждый из перечисленных слоев или функций.

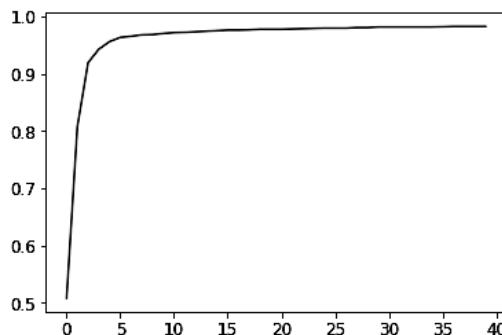
`nn.ReLU` — это модуль PyTorch, который делает то же, что и функция `F.relu`. Большинство функций, встречающихся в модели, совпадают по форме с модулями. Как правило, достаточно просто поменять `F` на `nn` и изменить заглавные буквы. При использовании `nn.Sequential` PyTorch требует от нас применения модульной версии. Так как модули являются классами, нужно их инстанцировать, почему вы и видите в следующем примере `nn.ReLU`.

Поскольку `nn.Sequential` — это модуль, мы можем получить его параметры, в результате чего будет возвращен их полный список из всех содержащихся в нем модулей. Давайте это проверим. Поскольку это более глубокая модель, мы продлим процесс обучения на несколько эпох и снизим саму скорость этого обучения:

```
learn = Learner(dls, simple_net, opt_func=SGD,
                 loss_func=mnist_loss, metrics=batch_accuracy)
learn.fit(40, 0.1)
```

В целях экономии места мы решили не показывать вам все 40 строк вывода. Процесс обучения записан в `learn_recorder`, где таблица вывода хранится в атрибуте `values`, так что мы можем графически отобразить точность в ходе обучения:

```
plt.plot(L(learn.recorder.values).itemgot(2));
```



Мы также можем посмотреть итоговую точность:

```
learn.recorder.values[-1][2]
0.982826292514801
```

Сейчас у нас есть кое-что волшебное.

- Функция, которая при получении верного набора параметров способна решить любую задачу с любым уровнем точности (нейронная сеть).
- Способ нахождения наилучшего набора параметров для любой функции (стохастический градиентный спуск).

Вот почему глубокое обучение способно на такие фантастические вещи. Комбинация простых техник может реально решить любую задачу — и это очень важно осознать. Выглядит все слишком хорошо, чтобы быть правдой, — наверняка все должно быть гораздо сложнее? На мы просто посоветуем вам попробовать всё самим. Мы только что опробовали датасет MNIST, и вы видели результаты. А поскольку мы все делаем с нуля сами (за исключением вычисления градиентов), то вы видите, что за этим не стоит никакой особой магии.

Углубляемся

Ничто не мешает выйти за рамки использования всего двух слоев. Можно добавлять любое желаемое количество при условии размещения между линейными слоями нелинейной функции. Но как вы скоро увидите, чем глубже становится модель, тем сложнее оптимизировать ее параметры на практике. Позже вы познакомитесь с некоторыми простыми, но потрясающе эффективными техниками обучения более глубоких моделей.

Нам уже известно, что одной нелинейной функции с двумя линейными слоями достаточно для аппроксимации любой функции. Зачем же использовать более глубокие модели? Причина в производительности. В углубленных моделях (то есть имеющих больше слоев) нам не обязательно использовать так же много параметров. Как оказывается, можно использовать уменьшенные матрицы с большим числом слоев и получать результаты, превосходящие подход с использованием более крупных матриц при меньшем количестве слоев.

Это означает, что можно обучать модель ускоренно и с меньшими затратами памяти. В 1990-х исследователи были настолько сосредоточены на теореме универсальной аппроксимации, что лишь некоторые из них экспериментировали с использованием более чем одной нелинейной функции. Такой теоретический подход затормозил развитие этой области на несколько лет. Тем не менее некоторые исследователи проводили эксперименты с глубокими моделями и в результате показали, что эти модели могут давать намного лучший результат на практике. В итоге были сформированы выводы, показавшие, почему это произошло. Сегодня уже кажется совершенно необычным использование нейронной сети всего с одной нелинейностью.

Здесь мы показываем, что происходит, когда мы выполняем обучение 18-слойной модели с помощью подхода, использованного нами в главе 1:

```
dls = ImageDataLoaders.from_folder(path)
learn = cnn_learner(dls, resnet18, pretrained=False,
                     loss_func=F.cross_entropy, metrics=accuracy)
learn.fit_one_cycle(1, 0.1)
```

epoch	train_loss	valid_loss	accuracy	time
0	0.082089	0.009578	0.997056	00:11

Точность приблизилась к 100 %! Это большая разница в сравнении с нашим простым примером нейронной сети. Но как вы узнаете из оставшейся части книги, вам потребуется использовать всего несколько простых приемов, чтобы самостоятельно получить аналогичный результат с нуля. Вам уже известны основные ключевые составляющие. (Конечно же, даже зная все возможные приемы, лучше работать с заранее подготовленными классами, предоставляемыми PyTorch и fastai, так как они уберегут вас от необходимости самостоятельно продумывать все мелкие детали.)

Сводка терминов

Поздравляем! Теперь вы знаете, как создавать и обучать глубокую нейронную сеть с нуля. Мы проделали большую работу, чтобы достичь этой точки, но вы можете удивиться, насколько это все в реальности просто.

Сейчас есть отличная возможность определить и рассмотреть некоторые ключевые понятия терминологии.

Нейронная сеть содержит множество чисел, но они бывают только двух типов: вычисляемые числа и параметры, на основе которых они вычисляются. Таким образом возникают два наиболее важных элемента.

Активации

Вычисляемые линейным и нелинейным слоями числа.

Параметры

Числа, которые инициализируются как случайные, а затем оптимизируются (то есть числа, определяющие модель).

Мы часто будем говорить об активациях и параметрах. Помните, что они имеют специфичные значения. Они являются числами. Это не абстрактные понятия, а конкретные числа в модели. Чтобы стать успешным практиком глубокого обучения, нужно привыкнуть рассматривать активации и параметры, графически отображая и тестируя их на соответствие нужному вам поведению.

Все активации и параметры содержатся в *тензорах*, которые являются простыми массивами правильной формы, например матрицами. Матрицы содержат строки

и столбцы, которые называются *осями и измерениями*. Число измерений тензора обозначает его *ранг*. Есть несколько особых тензоров:

- ранг 0: скаляр;
- ранг 1: вектор;
- ранг 2: матрица.

Нейронная сеть содержит некоторое количество слоев, каждый из которых либо *линеен*, либо *нелинеен*. Обычно мы переключаемся между этими типами слоев в нейронной сети. Иногда люди называют линейный слой и следующий за ним нелинейный одним слоем. Да, это может запутывать. В некоторых случаях нелинейный слой также называют *функцией активации*.

Таблица 4.1 приводит сводку ключевых понятий, связанных с SGD.

Таблица 4.1. Таблица глубокого обучения

Термин	Значение
ReLU	Функция, возвращающая 0 для отрицательных чисел, но не меняющая положительные
Мини-пакет	Небольшая группа входных данных и меток, собранных в два массива. Шаг градиентного спуска обновляется для этого пакета, а не для всей эпохи
Прямой проход прогнозов (<i>forward pass</i>)	Применение модели к входным данным и вычисление
Потери	Значение, отражающее, насколько эффективна или неэффективна модель
Градиент	Производная потери в отношении некоторых параметров модели
Обратный проход модели (<i>backward pass</i>)	Вычисление градиентов потерь в отношении ко всем параметрам
Градиентный спуск	Совершение шага в противоположном градиентам направлении с целью улучшения параметров функции
Скорость обучения	Размер шага, совершаемого при применении SGD для обновления параметров модели



НАПОМИНАНИЕ

Если вы решили пропустить главы 2 и 3, чтобы поскорее узнать о внутренних процессах обучения, то вам все же придется вернуться к главе 2 — понимание ее материала скоро вам понадобится!

Вопросник

1. Как на компьютере отображается изображение в оттенках серого? А цветное?
2. Как структурированы файлы и каталоги в датасете MNIST_SAMPLE? Почему именно так?
3. Объясните, как работает подход «пиксельного сходства» при классификации цифр.
4. Что такое генератор списка? Создайте такой, который будет выбирать нечетные числа из списка и удваивать их.
5. Что такое тензор ранга 3?
6. В чем разница между рангом тензора и его формой? Как можно получить ранг, исходя из формы?
7. Что такое RMSE и норма L1?
8. Как можно применить вычисление к тысячам чисел сразу и в несколько тысяч раз быстрее, чем при использовании цикла Python?
9. Создайте тензор 3Х3 или массив, содержащий числа от 1 до 9. Удвойте его. Выберите четыре числа из правой нижней части.
10. Что такое транслирование?
11. Как обычно вычисляются метрики: с помощью обучающей выборки или контрольной? Почему?
12. Что такое SGD?
13. Почему SGD использует мини-пакеты?
14. Из каких семи шагов состоит SGD для машинного обучения?
15. Как мы инициализируем веса модели?
16. Что такое потери?
17. Почему не всегда можно использовать высокую скорость обучения?
18. Что такое градиент?
19. Нужно ли вам знать, как вычислять градиенты самостоятельно?
20. Почему мы не можем использовать в качестве функции потерь точность?
21. Нарисуйте сигмоидную функцию. В чем особенность ее формы?
22. В чем отличие между функцией потерь и метрикой?
23. Какая функция вычисляет новые веса на основе скорости обучения?
24. Что делает класс `DataLoader`?

-
25. Напишите псевдокод, отображающий основные шаги, совершаемые в каждой эпохе для SGD.
 26. Создайте функцию, которая при получении двух аргументов [1, 2, 3, 4] и 'abcd' будет возвращать [(1, 'a'), (2, 'b'), (3, 'c'), (4, 'd')]. В чем особенность структуры полученных на выводе данных?
 27. Какова роль `view` в PyTorch?
 28. Что такое параметры смещения в нейронной сети? Зачем они нам нужны?
 29. Каково назначение оператора `@` в Python?
 30. Что делает метод `backward`?
 31. Зачем нужно обнулять градиенты?
 32. Какую информацию нужно передавать в `Learner`?
 33. Покажите код Python или псевдокод для основных шагов цикла обучения.
 34. Что такое ReLU? Нарисуйте ее график для значений от -2 до +2.
 35. Что такое функция активации?
 36. В чем разница между `F.relu` и `nn.ReLU`?
 37. Теорема универсальной аппроксимации показывает, что любую функцию можно аппроксимировать настолько, насколько необходимо, используя всего одну нелинейность. Зачем же мы обычно используем больше?

Дополнительные задания

1. Создайте собственную реализацию `Learner` с нуля, взяв за основу обучающий цикл, показанный в этой главе.
2. Выполните все шаги этой главы, используя полные датасеты MNIST (для всех цифр, а не только для троек и семерок). Это значимый проект, который потребует от вас немало времени. Вам понадобится проделать самостоятельное исследование, чтобы понять, как преодолеть те или иные возникающие преграды.

ГЛАВА 5

Классификация изображений

Теперь, когда вы понимаете, что собой представляет и для чего используется глубокое обучение, а также знаете, как создавать и развертывать модель, настало время перейти к более подробному изучению. В идеальном мире практикам глубокого обучения не пришлось бы знать каждую деталь внутреннего устройства процессов, но в реальности пока все иначе. Поэтому чтобы заставить модель по-настоящему работать, причем надежно, важно наладить работу множества деталей, а также произвести их проверку. Для этого придется заглянуть внутрь нейронной сети в процессе ее обучения и работы, чтобы выявить возможные проблемы и найти их решения.

С этого момента мы начнем погружаться в изучение механизмов глубокого обучения. Какова архитектура модели компьютерного зрения, модели NLP, табличной модели и т. д.? Как создать архитектуру, соответствующую нуждам конкретной предметной области? Как добиться наилучших результатов от процесса обучения? Как весь этот процесс ускорить? Что нужно менять при изменении датасетов?

Начнем с повторения тех же основных операций, которые рассмотрели в главе 1, но сделаем при этом две вещи:

- улучшим их;
- применим их к более широкому разнообразию видов данных.

Для этого понадобится изучить все элементы глубокого обучения, к которым относятся различные типы слоев, методы регуляризации, оптимизаторы, навык совмещения слоев в архитектуре, техники разметки данных и многое другое. Но мы не собираемся просто вывалить все это на вас. Мы будем постепенно вводить все эти элементы, решая задачи, относящиеся к реализуемым нами проектам.

От собак и кошек к породам домашних животных

В нашей самой первой модели мы учились отличать собак от кошек. Всего несколько лет назад это задание считалось весьма трудным, сегодня же его можно

отнести к очень простым. Мы не сможем показать вам нюансы обучения моделей для решения этой задачи, так как уже получили практически идеальный результат, не озадачиваясь подробностями. Но как оказывается, этот же датасет можно использовать для работы над гораздо более серьезной задачей: определением породы домашнего питомца по его изображению.

В главе 1 мы представляли приложения как уже решенные задачи. Но в реальности все работает иначе. Мы начинаем с датасета, о котором ничего не знаем. После этого нам нужно выяснить, как он собран, как извлекать из него нужные нам данные и как эти данные вообще выглядят. На протяжении оставшейся части книги мы будем показывать вам, как решать эти задачи на практике, включая все промежуточные шаги, необходимые для понимания используемых нами данных и попутной проверки процесса моделирования.

Мы уже скачали датасет Pets и можем получить его путь, используя тот же код, что и в главе 1:

```
from fastai2.vision.all import *
path = untar_data(URLs.PETS)
```

Теперь, если вы хотите понять, как определять вид каждого животного на каждом изображении, нужно разобраться, как эти данные выстроены. Подробные детали структуры данных являются жизненно важным элементом пазла глубокого обучения. Обычно данные предоставляются в одном из следующих видов.

- Отдельными файлами, которые представляют элементы данных, такие как текст документов или изображения. Эти файлы обычно либо организованы по каталогам, либо используют имена, отражающие информацию об их элементах.
- Таблицей данных (например, в формате CSV), в которой каждая строка представляет элемент и может включать имена файлов, обеспечивающие связь между данными таблицы и данными в других форматах, таких как текст документов и изображения.

Однако из этих правил есть исключения, особенно в таких областях, как геномика, где могут встречаться форматы двоичных баз данных или даже сетевые потоки, но в целом подавляющее большинство датасетов, с которыми вы будете работать, используют некую комбинацию двух перечисленных форматов.

Для просмотра содержимого нашего датасета можно использовать метод `ls`:

```
path.ls()
(#3) [Path('annotations'), Path('images'), Path('models')]
```

Здесь мы видим, что этот набор данных содержит каталоги *images* (изображения) и *annotations* (аннотации). На сайте (<https://oreil.ly/xveoN>) этого датасета

поясняется, что каталог *annotations* содержит информацию о том, где находятся животные, а не о том, кем они являются. Мы же в этой главе займемся классификацией, а не локализацией, то есть нас интересует то, кем являются животные, а не их расположение. Поэтому пока использовать каталог *annotations* мы не будем. А сейчас заглянем в директорию *images*:

```
(path/"images").ls()  
  
(#7394) [Path('images/great_pyrenees_173.jpg'), Path('images/wheaten_terrier_46.  
jpg'), Path('images/Ragdoll_262.jpg'), Path('images/german_shorthaired_3.  
jpg'), Path('images/american_bulldog_196.jpg'), Path('images/boxer_188.  
jpg'), Path('images/staffordshire_bull_terrier_173.jpg'), Path('images/basset_  
hound_71.jpg'), Path('images/staffordshire_bull_terrier_37.jpg'), Path('images/  
yorkshire_terrier_18.jpg')...]
```

Большинство функций и методов fastai, возвращающих коллекцию, используют класс под названием `L`. Его можно рассматривать как расширенную версию стандартного типа `list` в Python, упрощающую использование наиболее распространенных операций. К примеру, когда мы отображаем объект этого класса в блокноте, он появляется в приведенном здесь формате. Первое, что показывается, — это число элементов коллекции с префиксом `#`. В предшествующем выводе вы также увидите, что за списком следует символ многоточия. Это означает, что показаны только первые несколько элементов, и это хорошо, потому что мы не хотели бы видеть 7000 имен файлов на экране!

Изучив эти имена файлов, мы можем понять, как они структурированы. Каждое из них содержит название породы, сопровождаемое нижним подчеркиванием (`_`), числом и в конце расширением файла. Нам нужно создать часть кода, извлекающую эти породы из одного `Path`. Это легко сделать, используя блокнот Jupyter, потому что мы можем шаг за шагом создавать рабочие компоненты, а затем использовать их для всего датасета. При этом нужно быть осторожными и не спешить. Например, если внимательно присмотреться, то можно заметить, что некоторые породы содержат несколько слов, поэтому нам нельзя прерывать чтение при встрече первого знака `_`. Выберем одно из имен файлов, на примере которого протестируем создаваемый нами код:

```
fname = (path/"images").ls()[0]
```

Наиболее мощный и гибкий способ извлечения информации из подобных строк — это использование *регулярного выражения*, иначе называемого *регексом* (или, по-русски, *регуляркой*). Регекс является особой строкой, написанной на языке регулярных выражений, которая определяет главное правило для принятия решения, проходит ли другая строка проверку (то есть «совпадает» ли с этим регулярным выражением). Иногда она может также использоваться для извлечения из этой другой строки одной или нескольких частей. Здесь же нужно регулярное выражение, извлекающее из имени файла название породы.

Физически не хватит места подробно описать регулярные выражения, но в интернете есть множество отличных уроков, к тому же многие из вас уже знакомы с этим прекрасным инструментом. Если же нет, то ничего страшного — есть отличная возможность все наверстать. Мы считаем, что эти выражения являются одним из наиболее полезных инструментов в наборе программиста, и многие наши студенты относят эту тему к наиболее интересным и полезным. Поэтому смело гуглите уроки и изучайте материал. Кстати, на сайте книги (<https://book.fast.ai>) также приведен список наиболее ценных, на наш взгляд, руководств¹.



СЛОВО АЛЕКСИСУ

Регулярные выражения не только очень полезны, но также несут в себе весьма любопытную основу. Их называют «регулярными», потому что изначально они являются примерами «регулярного» языка, низшей ступени иерархии Хомского. Это грамматическая классификация, разработанная лингвистом Ноамом Хомски (Noam Chomsky), который также написал «Синтаксические структуры», новаторскую работу по поиску формальной грамматики, лежащей в основе человеческого языка. Это одна из прелестей информатики: может оказаться, что молоток, которым вы пользуетесь каждый день, на самом деле упал с космического корабля.

Когда вы пишете регулярное выражение, то лучше всего начинать с его проверки на одном примере. Используем метод `.findall`, чтобы попробовать регулярку для имени файла объекта `fname`:

```
re.findall(r'(.+)\d+.jpg$', fname.name)
['great_pyrenees']
```

Регулярное выражение извлекает все знаки вплоть до последнего подчеркивания до тех пор, пока последующими знаками не окажутся цифры, а затем расширение JPEG.

Теперь, когда мы убедились, что выражение работает для этого примера, можно использовать его для разметки всего датасета. В fastai есть множество классов, которые помогают в этом процессе. Для разметки с помощью регулярных выражений можно использовать класс `RegexLabeeler`. В этом примере мы задействуем API блока данных, который уже видели в главе 2 (на практике мы почти всегда используем этот API, так как он намного более гибкий, чем простые фабричные методы, которые мы видели в главе 1):

```
pets = DataBlock(blocks = (ImageBlock, CategoryBlock),
                  get_items=get_image_files,
```

¹ Подробнейшую информацию по регулярным выражениям и шаблоны их использования в разных языках программирования можно найти в книге: Фридл Дж. Регулярные выражения. 3-е изд. — СПб.: Питер, 2018. — Примеч. ред.

```
    splitter=RandomSplitter(seed=42),
    get_y=using_attr(RegexLabeller(r'(.+)\d+.jpg$', 'name'),
                     item_tfms=Resize(460),
                     batch_tfms=aug_transforms(size=224, min_scale=0.75))
dls = pets.dataloaders(path/"images")
```

Одна важная часть `DataBlock`, которую мы еще не видели, содержится в следующих двух строках:

```
item_tfms=Resize(460),
batch_tfms=aug_transforms(size=224, min_scale=0.75)
```

Они реализуют стратегию аугментации данных fastai, которую мы называем *подготовкой размера* (*resizing*). Это особый способ выполнения аугментации изображений, целью которого является минимизация утраты данных с сохранением производительности.

Подготовка размера

Нам нужно, чтобы у изображений были одинаковые размеры, что позволит объединить их в тензоры для передачи на обработку в GPU. Нам также нужно минимизировать число отдельных вычислений аугментации. Требования производительности предполагают, что мы по возможности должны совмещать выполняемые аугментацией преобразования, уменьшая их общее число (для уменьшения количества вычислений и количества операций с потерями), и преобразовывать изображения в единые размеры (для их более эффективной обработки GPU).

Сложность же в том, что если выполнять аугментацию после изменения размера изображения, то это может внести в него ложные пустые зоны, ухудшить данные или то и другое сразу. Например, при вращении изображения на 45 градусов углы его новых границ заполняются пустотами, которые модель точно ничему не научат. Кроме того, многие операции вращения и масштабирования потребуют интерполяции для создания пикселей. Но несмотря на то что эти интерполированные пиксели получаются из данных исходного изображения, они имеют более низкое качество.

Для решения этих проблем в подготовке размера подразумеваются два этапа, показанные на рис. 5.1.

1. Увеличение изображения до относительно «больших» размеров, чтобы они оказались существенно больше целевых размеров для обучения.
2. Объединение всех стандартных операций аугментации (включая изменение размера на итоговый целевой) в одну и выполнение совмещенной операции на GPU всего один раз в конце обработки. Это устраняет необ-

ходимость выполнения всех операций по отдельности и многократного интерполирования.

Первый шаг, изменение размера, создает изображения настолько большие, чтобы имелся свободный отступ, позволяющий выполнять дальнейшую аугментацию внутренних областей без создания пустых зон. Это преобразование изменяет размер до квадрата путем крупномасштабной обрезки изображения. В обучающей выборке область обрезки выбирается случайно, а размер обрезки выбирается так, чтобы только включить всю ширину или высоту изображения, но не более. Во втором шаге для аугментации всех данных используется GPU, и все потенциально деструктивные операции выполняются вместе с одной интерполяцией в конце.

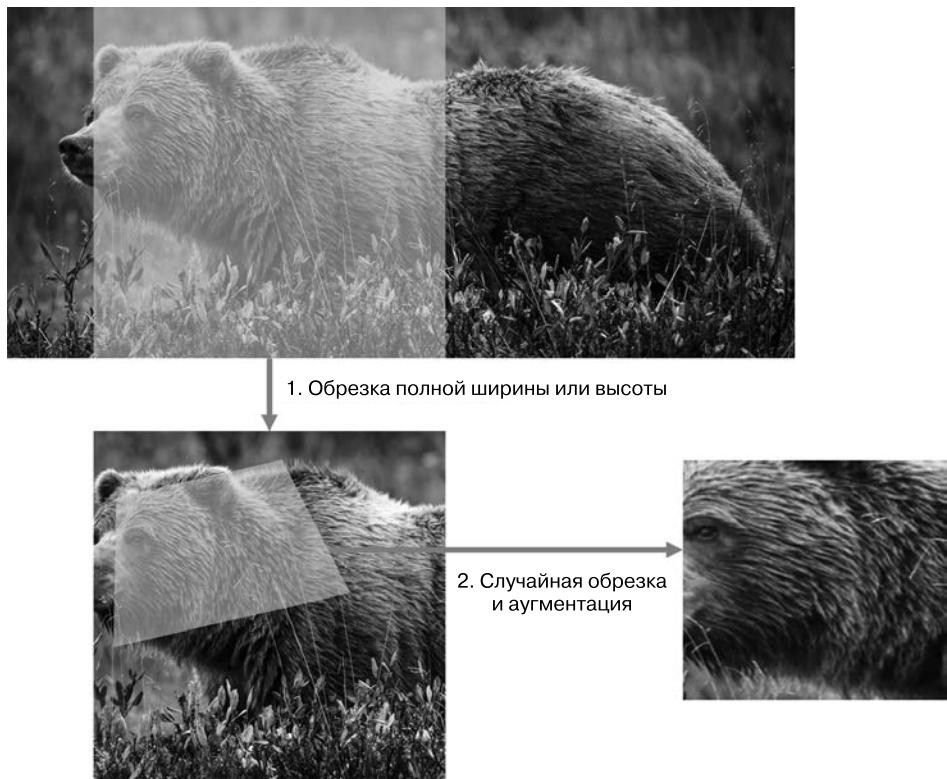


Рис. 5.1. Подготовка размера в обучающей выборке

На картинке показаны два шага.

1. *Обрезка полной ширины или высоты.* Это прописано в `item_tfms`, следовательно, применяется к каждомуциальному изображению до его копиро-

вания в GPU. Целью является обеспечить для всех изображений единый размер. В обучающей выборке область обрезки выбирается случайно. В контрольной же выборке всегда выбирается центральный квадрат изображения.

2. *Случайная обрезка и аугментация.* За это отвечает `batch_tfms`, а это означает применение этого процесса сразу ко всему пакету в GPU, что объясняет высокую скорость выполнения. Для контрольной выборки здесь выполняется только изменение размера до итогового варианта, необходимого модели. В обучающей же выборке сначала выполняется случайная обрезка и другие действия аугментации.

В fastai для реализации этого процесса вы используете `Resize` в качестве первичного преобразования (обрезки) элемента до крупного размера, а `RandomResizedCrop` — в качестве последующего преобразования пакета элементов меньшего размера. `RandomResizedCrop` будет добавлен автоматически, если вы включите в функцию `aug_transforms` параметр `min_scale`, как мы делали в вызове `DataBlock` в предыдущем разделе. В качестве альтернативы для изначального `Resize` вы можете использовать `pad` или `squish` вместо `crop` (по умолчанию).

На рис. 5.2 показано отличие между правым изображением, которое было увеличено, интерполировано, а затем снова интерполировано (такой подход используется другими библиотеками глубокого обучения), и левым, которое было увеличено и повернуто в одну операцию, а затем единожды интерполировано (подход fastai).



Рис. 5.2. Сравнение стратегии аугментации данных fastai (слева)
и традиционного подхода (справа)

Здесь видно, что изображение справа менее четкое, а в его левом нижнем углу присутствуют артефакты отражения. Кроме того, трава в верхней левой части полностью исчезла. Наблюдения показали, что на практике использование подготовки размера значительно улучшает точность моделей, а иногда и ускоряет их работу.

В fastai также есть простые способы проверить, как выглядят ваши данные, прежде чем начинать обучение модели, что является очень важным шагом, который мы и разберем следующим.

Проверка и отладка DataBlock

Никогда нельзя предполагать, что код просто так заработает идеально. Написание DataBlock можно сравнить с разработкой схемы. В случае обнаружения синтаксических ошибок в коде вы будете получать об этом сообщения, но нет никакой гарантии, что ваш шаблон будет работать с вашим источником данных, как вы планировали. Поэтому перед обучением модели всегда нужно проверять данные.

Это делается с помощью метода `show_batch`:

```
dls.show_batch(nrows=1, ncols=3)
```



Посмотрите на каждое изображение и убедитесь, что везде стоит правильная метка породы. Специалистам ML нередко приходится работать с данными, в которых они разбираются куда хуже экспертов соответствующей области. Например, я вот на самом деле не знаю породы многих домашних животных, так как экспертом в этой области не являюсь. Поэтому в данном случае я бы поискал изображения неизвестных мне пород в Google и сопоставил с тем, что вижу в текущем выводе.

Если при построении DataBlock вы допустите ошибку, то, скорее всего, не узнаете о ней до этого шага. Для исправления же подобных ошибок мы рекомендуем использовать метод `summary`. Он попробует создать из предоставленного вами источника пакет со множеством подробностей. Кроме того, если ему это не удастся, вы увидите, в какой именно точке произошла ошибка, а библиотека постарается в этом случае помочь. Например, распространенная ошибка — это забыть использовать преобразование `Resize`, в результате чего у вас получаются изображения разных размеров, которые невозможно собрать в пакет. Вот как в таком случае будет выглядеть `summary` (сводная информация) (обратите внимание, что с момента написания книги конкретный текст может измениться, но общий смысл будет прежним):

```
pets1 = DataBlock(blocks = (ImageBlock, CategoryBlock),
                  get_items=get_image_files,
                  splitter=RandomSplitter(seed=42),
                  get_y=using_attr(RegexLabeller(r'(.)_\d+.jpg$'), 'name'))
```

```
pets1.summary(path/"images")

Setting-up type transforms pipelines
Collecting items from /home/sgugger/.fastai/data/oxford-iiit-pet/images
Found 7390 items
2 datasets of sizes 5912,1478
Setting up Pipeline: PILBase.create
Setting up Pipeline: partial -> Categorize

Building one sample
Pipeline: PILBase.create
  starting from
    /home/sgugger/.fastai/data/oxford-iiit-pet/images/american_bulldog_83.jpg
  applying PILBase.create gives
    PILImage mode=RGB size=375x500

Pipeline: partial -> Categorize
  starting from
    /home/sgugger/.fastai/data/oxford-iiit-pet/images/american_bulldog_83.jpg
  applying partial gives
    american_bulldog
  applying Categorize gives
    TensorCategory(12)

Final sample: (PILImage mode=RGB size=375x500, TensorCategory(12))

Setting up after_item: Pipeline: ToTensor
Setting up before_batch: Pipeline:
Setting up after_batch: Pipeline: IntToFloatTensor

Building one batch
Applying item_tfms to the first sample:
Pipeline: ToTensor
  starting from
    (PILImage mode=RGB size=375x500, TensorCategory(12))
  applying ToTensor gives
    (TensorImage of size 3x500x375, TensorCategory(12))

Adding the next 3 samples

No before_batch transform to apply

Collating items in a batch
Error! It's not possible to collate your items in a batch
Could not collate the 0-th members of your tuples because got the following
shapes:
torch.Size([3, 500, 375]),torch.Size([3, 375, 500]),torch.Size([3, 333, 500]),
torch.Size([3, 375, 500])
```

Здесь видно, как именно мы собрали данные и разделили их, как перешли от имени файла к *образцу* (кортеж (изображение, категория)), какие после этого были применены преобразования и как произошел сбой объединения полученных образцов в пакет (из-за различия их форм).

Как только вы сочли, что данные выглядят как надо, попробуйте обучить простую модель. Мы замечаем, что люди слишком долго откладывают обучение реальной модели. В итоге они не могут разобраться, как выглядят их базовые результаты. Возможно, ваша задача не потребует больших инженерных усилий, связанных с конкретной предметной областью. А может, на основе используемых данных вообще не получится обучить модель. В этом нужно разобраться как можно раньше.

Для этого начального теста мы используем ту же простую модель, что использовали в главе 1:

```
learn = cnn_learner(dls, resnet34, metrics=error_rate)
learn.fine_tune(2)
```

epoch	train_loss	valid.loss	error_rate	time
0	1.491732	0.337355	0.108254	00:18

epoch	train_loss	valid.loss	error_rate	time
0	0.503154	0.293404	0.096076	00:23
1	0.314759	0.225316	0.066306	00:23

Как мы уже вкратце говорили, таблица, отображаемая в процессе подстройки модели, показывает нам результаты после каждой эпохи обучения. Помните, что эпоха — это один полный проход по всем изображениям набора данных. Столбцы таблицы показывают средние потери по элементам обучающей выборки, потери для контрольной выборки, а также любые запрошенные нами метрики — в этом случае частоту ошибок.

Помните, что *потери* — это та функция, которую мы решили использовать для оптимизации параметров модели. Но на самом деле мы не сообщали fastai, какую функцию потерь хотим использовать. Так что же она делает? Обычно fastai, как правило, старается выбрать подходящую функцию потерь на основе используемых типов данных и модели. В текущей ситуации мы обрабатываем данные изображений и делаем вывод по категориям, поэтому fastai по умолчанию использует функцию *потерь перекрестной энтропии*.

Перекрестная энтропия

Перекрестная энтропия — это функция потерь, аналогичная использованной в предыдущей главе, но (как мы увидим) она имеет два преимущества:

- работает даже при наличии в зависимой переменной более двух категорий;
- повышает скорость и надежность обучения.

Чтобы понять, как эта функция работает для зависимых переменных, имеющих более двух категорий, сначала нужно понять, как выглядят фактические данные и активации, которые ею рассматриваются.

Активации и метки

Посмотрим на активации нашей модели. Для получения пакета реальных данных из `DataLoaders` мы можем использовать метод `one_batch`:

```
x,y = dls.one_batch()
```

Как видите, это возвращает зависимые и независимые переменные в виде мини-пакета. Посмотрим, что содержится в зависимой:

```
y
```

```
TensorCategory([11, 0, 0, 5, 20, 4, 22, 31, 23, 10, 20, 2, 3, 27, 18, 23,
> 33, 5, 24, 7, 6, 12, 9, 11, 35, 14, 10, 15, 3, 3, 21, 5, 19, 14, 12,
> 15, 27, 1, 17, 10, 7, 6, 15, 23, 36, 1, 35, 6,
> 4, 29, 24, 32, 2, 14, 26, 25, 21, 0, 29, 31, 18, 7, 7, 17],
> device='cuda:5')
```

Размер пакета 64, значит, в этом тензоре 64 строки. Каждая строка — это одно целое число между 0 и 36, представляющее 37 возможных пород домашних животных. Мы можем просмотреть прогнозы (активации последнего слоя нейронной сети) с помощью `Learner.get_preds`. Эта функция получает либо индекс датасета (0 для обучающего и 1 для контрольного), либо итератор пакетов. Значит, для получения прогнозов мы можем передать ей простой список с пакетом. По умолчанию она возвращает прогнозы и цели, но так как цели у нас уже есть, то мы можем проигнорировать их, присвоив особой переменной `_`:

```
preds,_ = learn.get_preds(dl=[(x,y)])
preds[0]

tensor([7.9069e-04, 6.2350e-05, 3.7607e-05, 2.9260e-06, 1.3032e-05, 2.5760e-05,
6.2341e-08, 3.6400e-07, 4.1311e-06, 1.3310e-04, 2.3090e-03, 9.9281e-01, 4.6494e-05,
6.4266e-07, 1.9780e-06, 5.7005e-07,
3.3448e-06, 3.5691e-03, 3.4385e-06, 1.1578e-05, 1.5916e-06, 8.5567e-08,
5.0773e-08, 2.2978e-06, 1.4150e-06, 3.5459e-07, 1.4599e-04, 5.6198e-08, 3.4108e-07,
2.0813e-06, 8.0568e-07, 4.3381e-07,
1.0069e-05, 9.1020e-07, 4.8714e-06, 1.2734e-06, 2.4735e-06])
```

Фактические прогнозы — это 37 вероятностей между 0 и 1, которые вместе составляют 1:

```
len(preds[0]),preds[0].sum()
```

```
(37, tensor(1.0000))
```

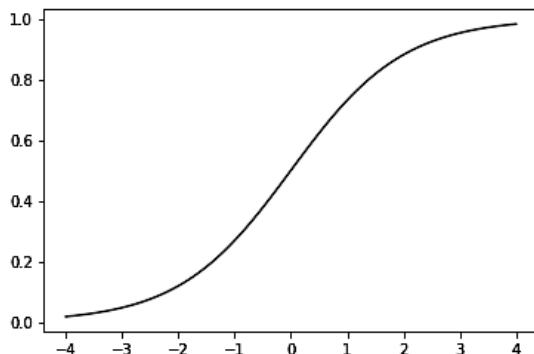
Для преобразования этих активаций модели в прогнозы, подобные этому, мы использовали функцию активации, называемую *softmax*.

Softmax

Используем функцию активации softmax в последнем слое нашей модели классификации, чтобы все активации входили в диапазон от 0 до 1 и вместе составляли 1.

Softmax является сигмоидной функцией, которую мы уже видели раньше, но на всякий случай напомним, что выглядит она так:

```
plot_function(torch.sigmoid, min=-4,max=4)
```



Можно применить эту функцию к одному столбцу активации нейронной сети и получить столбец чисел между 0 и 1 — это сделает ее очень полезной функцией активации для последнего слоя.

А теперь задумайтесь, что происходит, если требуется больше категорий для цели (как в нашем случае с 37 породами). Это означает, что уже нужен не один столбец, а гораздо больше активаций, то есть по активации для каждой *категории*. Мы можем создать, например, нейронную сеть, прогнозирующую тройки и семерки, которая возвращает по одной активации для каждого из этих двух классов, — это будет хорошим шагом на пути создания более общего подхода. Давайте для этого примера просто используем случайные числа со стандартным отклонением 2 (то есть мы умножаем `randn` на 2), предполагая, что у нас есть шесть изображений и две возможные категории (где первый столбец представляет тройки, а второй — семерки):

```
acts = torch.randn((6,2))*2
acts

tensor([[ 0.6734,  0.2576],
```

```
[ 0.4689,  0.4607],
[-2.2457, -0.3727],
[ 4.4164, -1.2760],
[ 0.9233,  0.5347],
[ 1.0698,  1.6187]])
```

Мы не можем просто непосредственно применить здесь сигмоиду, поскольку получаемые строки не складываются в 1 (нам нужно, чтобы вероятность тройки в сумме с вероятностью семерки составляла 1):

```
acts.sigmoid()

tensor([[0.6623, 0.5641],
       [0.6151, 0.6132],
       [0.0957, 0.4079],
       [0.9881, 0.2182],
       [0.7157, 0.6306],
       [0.7446, 0.8346]])
```

В главе 4 наша нейронная сеть создала для каждого изображения одну активацию, которую мы передавали через функцию `sigmoid`. Эта единственная активация представляла уверенность модели в том, что на вводе подана тройка. Проблемы двоичности — это особый случай в задачах по классификации, потому что цель может интерпретироваться как одно логическое значение, что мы и делали в `mnist_loss`. Но проблемы двоичности также можно рассматривать в контексте более общей группы классификаторов с любым числом категорий: в данном случае у нас их две. Как мы видели в модели классификации медведей, нейронная сеть возвращает одну активацию для каждой категории.

Так что же в этом двоичном случае реально отражают активации? Одна пара активаций просто отражает *относительную* уверенность в том, что на вводе подана тройка, а не семёрка. Результирующие значения, будь они оба высоки или низки, не имеют значения — значение имеет лишь то, какое из них выше и на сколько.

Можно предположить, что раз это просто другой способ выражения той же самой задачи, то мы можем применить `sigmoid` непосредственно к варианту нейронной сети с двумя активациями. Так оно и есть! Мы можем просто взять *разность* между активациями нейронной сети, поскольку она отражает степень уверенности, что на вводе подана тройка, а не семёрка, после чего вычислить для этой разности сигмоиду.

```
(acts[:,0]-acts[:,1]).sigmoid()

tensor([0.6025, 0.5021, 0.1332, 0.9966, 0.5959, 0.3661])
```

Тогда второй столбец (вероятность, что рассматривается семерка) будет просто этим значением, которое вычли из единицы. Теперь нам нужен такой способ

реализации, который также будет работать для более чем двух столбцов. Оказывается, функция `softmax` как раз им и является:

```
def softmax(x): return exp(x) / exp(x).sum(dim=1, keepdim=True)
```



ТЕРМИН: ЭКСПОНЕНТА (EXP)

Определяется как $c^{**}x$, где c — это особое число, приблизительно равное 2,718. Экспонента является функцией, обратной натуральному логарифму. Обратите внимание, что \exp всегда положительная и возрастает очень быстро!

Убедимся, что `softmax` возвращает для первого столбца те же значения, что и `sigmoid`, а для второго — результат вычитания этих значений из единицы:

```
sm_acts = torch.softmax(acts, dim=1)
sm_acts

tensor([[0.6025, 0.3975],
       [0.5021, 0.4979],
       [0.1332, 0.8668],
       [0.9966, 0.0034],
       [0.5959, 0.4041],
       [0.3661, 0.6339]])
```

`softmax` — это мультикатегориальный эквивалент `sigmoid`, который нам следует использовать во всех случаях, когда мы работаем с более чем двумя категориями, и вероятности этих категорий должны суммироваться в единицу. Мы также часто используем ее, когда есть всего две категории, просто чтобы добиться лучшей согласованности.

Можно создать и другие функции, имеющие свойства для всех активаций попадать в диапазон от 0 до 1 и суммироваться в единицу. Тем не менее ни одна из таких функций не будет иметь той же связи с сигмоидой, которая, как мы видели, отличается гладкостью и симметричностью. Кроме того, мы вскоре покажем вам, что функция `softmax` отлично работает в tandemе с функцией потерь, которую мы рассмотрим в следующем разделе.

Если у нас есть три выходные активации, как в известном нам классификаторе медведей, то вычисление `softmax` для одного изображения медведя будет выглядеть, как показано на рис. 5.3.

	output	exp	softmax
teddy		0.02 1.02	0.22
grizzly		-2.49 0.08	0.02
brown		1.25 3.49	0.76
		4.60	1.00

Рис. 5.3. Пример `softmax` для классификатора медведей

Что же эта функция делает на практике? Получение экспоненты гарантирует, что все числа будут положительными, а последующее деление на сумму гарантирует получение чисел, суммирующихся в единицу. У экспоненты, помимо этого, есть одно приятное свойство: если одно из чисел наших активаций x немного больше других, экспонента его увеличит (поскольку она возрастает экспоненциально), то есть в softmax число будет ближе к единице.

Очевидно, что функция softmax *действительно* как бы старается выбрать один класс среди других, поэтому идеально подходит для обучения классификатора, когда нам известно, что у каждой картинки есть определенная метка. (Обратите внимание, что результаты могут быть не такими идеальными, так как иногда необходимо, чтобы модель сообщала о том, что не может обнаружить классы, которые видела в процессе обучения, и не выбирала класс из-за его несколько завышенного показателя активации. В этом случае, возможно, будет лучше обучить модель, используя несколько двоичных столбцов вывода, каждый из которых будет использовать сигмоиду).

Softmax — это первая часть процесса определения потерь перекрестной энтропии, второй же частью является функция правдоподобия.

Логарифмическая функция правдоподобия

Когда мы вычисляли потери для примера MNIST в предыдущей главе, то использовали следующее:

```
def mnist_loss(inputs, targets):
    inputs = inputs.sigmoid()
    return torch.where(targets==1, 1-inputs, inputs).mean()
```

Теперь так же, как мы перешли от сигмоиды к softmax, нам нужно расширить функцию потерь, чтобы выйти за рамки двоичной классификации: она должна иметь возможность классифицировать любое число категорий (в данном случае у нас их тридцать семь). Наши активации, прошедшие softmax, представлены значениями между 0 и 1, суммируясь в единицу в каждой строке пакета прогнозов. Цели же представлены целыми числами от 0 до 36.

В двоичном случае для выбора между `inputs` и `1-inputs` мы использовали `torch.where`. Когда мы рассматриваем двоичную классификацию как общую задачу классификации с двумя категориями, то все становится еще легче, поскольку (как мы видели в предыдущем разделе) теперь у нас есть два столбца, содержащих эквивалент `inputs` и `1-inputs`. Поэтому все, что нам нужно сделать, — это произвести выбор из соответствующего столбца. Давайте попробуем реализовать это в PyTorch. Предположим, что для нашего синтетического примера троек и семерок мы используем следующие метки:

```
targ = tensor([0,1,0,1,1,0])
```

А вот активации softmax:

```
sm_acts
tensor([[0.6025, 0.3975],
       [0.5021, 0.4979],
       [0.1332, 0.8668],
       [0.9966, 0.0034],
       [0.5959, 0.4041],
       [0.3661, 0.6339]])
```

Тогда мы можем использовать все это для каждого элемента `targ`, чтобы выбрать подходящий столбец из `sm_acts`, использовав индексацию по тензору:

```
idx = range(6)
sm_acts[idx, targ]
tensor([0.6025, 0.4979, 0.1332, 0.0034, 0.4041, 0.3661])
```

Чтобы понять, что здесь происходит, объединим все столбцы в таблицу. Первые два столбца представляют наши активации, после них идут цели, индекс строки, и в конце показан результат приведенного выше кода:

3	7	targ	idx	loss
0.602469	0.397531	0	0	0.602469
0.502065	0.497935	1	1	0.497935
0.133188	0.866811	0	2	0.133188
0.99664	0.00336017	1	3	0.00336017
0.595949	0.404051	1	4	0.404051
0.366118	0.633882	0	5	0.366118

При рассмотрении этой таблицы становится очевидно, что последний столбец можно вычислить, взяв столбцы `targ` и `idx` в качестве индексов в матрице из двух столбцов, содержащей 3 и 7. Что и делает `sm_acts[idx, targ]`.

Самое интересное в том, что это также отлично работает и для большего числа столбцов. Убедиться в этом можно, рассмотрев, что произойдет, если мы добавим столбцы активации для каждой цифры (от 0 до 9) и `targ` будет содержать числа от 0 до 9. До тех пор пока эти столбцы активаций суммируются в единицу (так и будет, если использовать softmax), у нас будет получаться функция потерь, показывающая, насколько успешно мы прогнозируем каждую цифру.

Потери мы выбираем только из столбца, содержащего верную метку. Нам не нужно рассматривать другие столбцы, так как по определению softmax они складываются в единицу минус активация, соответствующая верной метке. Следовательно, увеличивая активацию верной метки до максимально возможного показателя, мы в то же время уменьшаем активации оставшихся столбцов.

PyTorch предоставляет функцию, которая делает то же, что и `sm_acts[range(n), targ]` (только получает отрицательные значения, так как при последующем применении логарифма мы получим отрицательные числа), и называется `nll_loss` (*NLL* означает *отрицательное логарифмическое правдоподобие*):

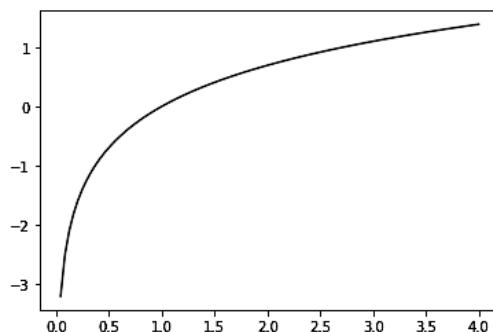
```
-sm_acts[idx, targ]
tensor([-0.6025, -0.4979, -0.1332, -0.0034, -0.4041, -0.3661])
F.nll_loss(sm_acts, targ, reduction='none')
tensor([-0.6025, -0.4979, -0.1332, -0.0034, -0.4041, -0.3661])
```

Несмотря на свое название, эта функция PyTorch не получает логарифм. В следующем разделе мы увидим причину этого, но сначала посмотрим, почему логарифм может оказаться полезным.

Применение логарифма

Функция, которую мы видели в предыдущем разделе, отлично работает в качестве функции потерь, но ее можно немного улучшить. Проблема в том, что мы используем вероятности, а вероятности не могут быть меньше 0 или больше 1. Это означает, что для нашей модели не будет разницы между прогнозом 0,99 и 0,999. Эти числа и впрямь очень близки, но с другой точки зрения, показатель 0,999 более уверененный, чем 0,99. Поэтому нам нужно преобразовать наши числа между 0 и 1 в такие, которые будут охватывать значения от минус бесконечности до 0. Для этой задачи подойдет особая математическая функция: *логарифм* (доступная в виде `torch.log`). Она не определяется для чисел меньше 0 и выглядит так:

```
plot_function(torch.log, min=0, max=4)
```



Припоминаете, что такое логарифм? Для логарифмической функции справедливо тождество:

```
y = b**a
a = log(y, b)
```

В этом случае мы предполагаем, что $\log(y, b)$ возвращает *log y по основанию b*. Тем не менее PyTorch не определяет \log так: в Python \log в качестве основания использует особое число e (2,718...).

Вполне возможно, что за последние лет двадцать вы о логарифмах даже и не вспоминали. Но они представляют собой математический принцип, который окажется необходим для многих процессов глубокого обучения, поэтому сейчас самое время освежить память. Главное, что нужно знать о логарифмах, — это следующее отношение:

$$\log(a*b) = \log(a)+\log(b)$$

Представленный в таком виде, он выглядит скучновато, но стоит подумать о том, что скрывается за этим значением. Смысл в том, что при экспоненциальном или мультипликативном увеличении основного сигнала логарифм увеличивается линейно. Этот принцип используется, к примеру, в шкале измерения силы землетрясений Рихтера или шкале децибелов, используемой для измерения уровня шума. Помимо этого, логарифмы используются для построения финансовых графиков, на которых требуется более отчетливо отразить совокупные темпы роста. Ученые в computer science любят использовать логарифмы, так как это означает, что преобразование, способное создать очень большие или очень малые числа, можно заменить сложением, которое уже не приведет к таким масштабам значений, которые компьютерам будет сложно обработать.

СЛОВО СИЛЬВЕЙНУ



Логарифмы любят не только ученые computer science! До появления компьютеров инженеры и ученые использовали особое аналоговое устройство, называемое логарифмической линейкой, с помощью которого выполняли множество математических операций, включая умножение чисел путем сложения их логарифмов. Логарифмы широко используются в физике для умножения очень больших и очень малых чисел, а также во многих других областях.

Получение среднего положительного или отрицательного логарифма наших вероятностей (в зависимости от того, верен класс или неверен) дает нам потери *отрицательной логарифмической вероятности*. В PyTorch `nll_loss` предполагает, что вы уже получили логарифм softmax, поэтому он не будет вычислять его за вас.

ОСТОРОЖНО, ОБМАНЧИВОЕ ИМЯ



`nll` в `nll_loss` означает «отрицательное логарифмическое правдоподобие», но фактически эта функция логарифм не получает, а предполагает, что вы уже его получили. В PyTorch есть функция под названием `log_softmax`, совмещающая `log` и `softmax` быстрым и точным образом. `nll_loss` же создана для использования после `log_softmax`.

Когда мы сначала вычисляем результат softmax, а затем его логарифмическое правдоподобие, то получаем комбинацию, которая и является *потерей перекрестной энтропии*. В PyTorch она выполняется как `nn.CrossEntropyLoss` (что фактически выполняет сначала `log_softmax`, а затем `nll_loss`):

```
loss_func = nn.CrossEntropyLoss()
```

Как вы видите, это класс, при инстанцировании которого вы получаете объект, действующий как функция:

```
loss_func(acts, targ)
tensor(1.8045)
```

Все функции потерь в PyTorch приводятся в двух формах: в только что показанной форме класса, а также в чисто функциональной форме, доступной в пространстве имен `F`:

```
F.cross_entropy(acts, targ)
tensor(1.8045)
```

Обе они работают прекрасно и могут использоваться в любой ситуации. Мы заметили, что большинство людей предпочитают версию класса, которая также чаще используется в официальной документации PyTorch и примерах, поэтому решили тоже придерживаться именно ее.

По умолчанию функции потерь PyTorch получают средний показатель потерь для всех элементов. Это можно отключить, определив `reduction='none'`:

```
nn.CrossEntropyLoss(reduction='none')(acts, targ)
tensor([0.5067, 0.6973, 2.0160, 5.6958, 0.9062, 1.0048])
```



СЛОВО СИЛЬВЕЙНУ

Интересный нюанс в отношении перекрестной энтропии возникает при рассмотрении ее градиента. Градиент `cross_entropy(a,b)` – это $\text{softmax}(a) - b$. Поскольку `softmax(a)` представляет финальную активацию модели, это означает, что градиент пропорционален разности между прогнозом и его целью. Это то же, что и среднеквадратичная ошибка в регрессии (если предположить, что нет функции финальной активации, такой как добавленная `y_range`), поскольку градиент $(a-b)^{**2}$ – это $2*(a-b)$. Поскольку градиент линеен, мы не увидим резких скачков или экспоненциальных возрастаний в градиентах, что обусловит более плавное обучение моделей.

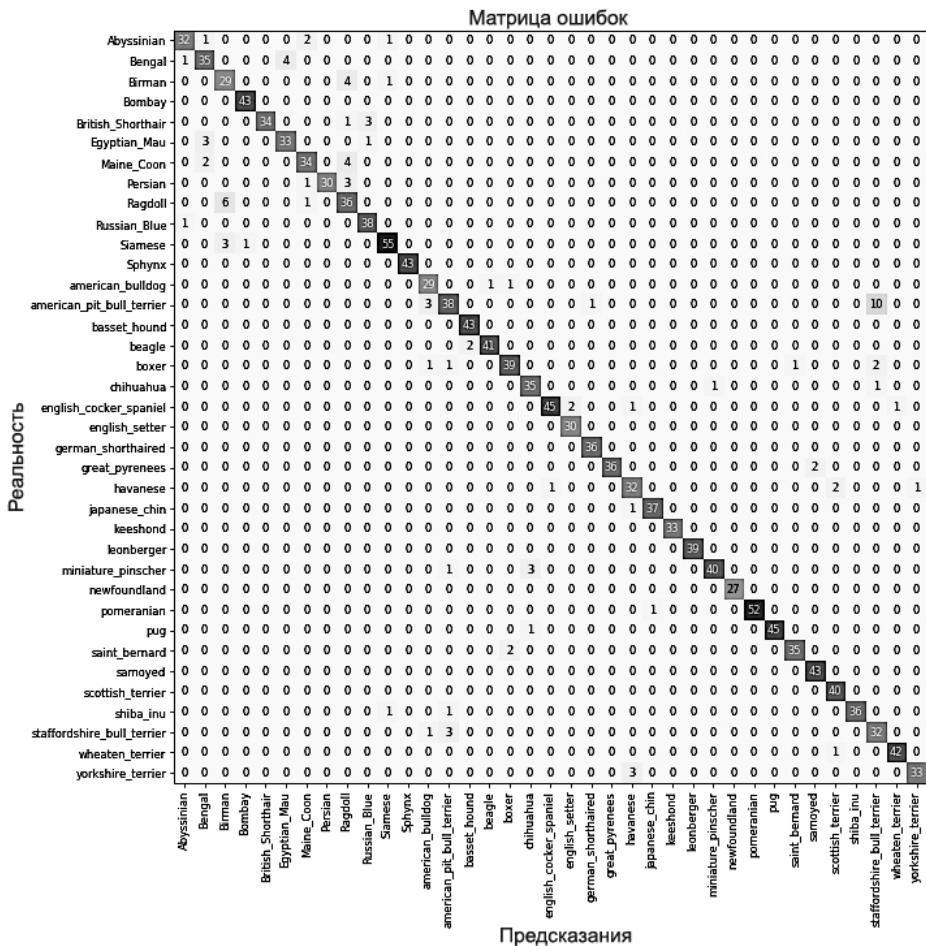
Вот мы и рассмотрели все стоящие за функцией потерь элементы. Но несмотря на то что с ее помощью мы получаем число, отражающее, как хорошо (или плохо) работает модель, это никак не помогает нам понять ее эффективность в целом. Давайте же теперь рассмотрим некоторые способы для интерпретирования прогнозов нашей модели.

Интерпретация модели

Очень трудно интерпретировать функции потерь непосредственно, потому что они созданы для дифференцирования и оптимизации компьютерами, а не для понимания людьми. Именно поэтому мы используем метрики. Они не используются в процессе оптимизации, но помогают людям понять, что вообще происходит. В этом случае точность уже выглядит достаточно неплохо. Так где же мы совершаем те немногие ошибки?

В главе 1 мы узнали, что можно использовать матрицу ошибок, чтобы увидеть, где модель справляется отлично, а где нет:

```
interp = ClassificationInterpretation.from_learner(learn)
interp.plot_confusion_matrix(figsize=(12,12), dpi=60)
```



Да уж, в этом случае матрицу читать сложновато. У нас есть 37 пород животных, а это значит, что мы получаем гигантскую матрицу из 37×37 записей. Вместо этого мы можем использовать метод `most_confused`, который просто покажет те ячейки матрицы, где прогнозы наименее верны (здесь со значением 5 и выше):

```
interp.most_confused(min_val=5)  
[('american_pit_bull_terrier', 'staffordshire_bull_terrier', 10),  
 ('Ragdoll', 'Birman', 6)]
```

Поскольку мы не эксперты по домашним породам, нам сложно понять, отражают ли эти ошибки в категориях фактические сложности в распознавании пород. И тут мы снова направляемся в Google, где недолгий поиск показывает, что обнаруженные ошибки относятся к таким различиям между породами, о которых спорят даже сами эксперты. Значит, можно быть спокойными — мы на верном пути.

Мы получили неплохую базовую модель, как же теперь ее улучшить?

Улучшение модели

Теперь мы рассмотрим ряд техник по улучшению обучения модели и ее совершенствованию. Параллельно с этим мы несколько подробнее расскажем о переносе обучения и о том, как тонко настраивать предварительно обученную модель наилучшим образом, не нарушая при этом предварительно обученных весов.

Первое, что нам понадобится установить при обучении модели, — это скорость обучения. В предыдущей главе мы видели, что от ее правильного выбора зависит итоговая эффективность всего процесса обучения. Так как же выбрать именно правильную? В fastai для этого предусмотрен специальный инструмент.

Поиск скорости обучения

Один из важнейших моментов, который мы должны учитывать при обучении модели, — это использование правильной скорости обучения. Если она будет слишком низкой, модели может потребоваться очень много эпох для обучения. Причем это приведет не только к дополнительной трате времени, но также создаст угрозу переобучения, потому что при каждом полном проходе по данным мы даем модели шанс запомнить их.

Может, тогда нам просто стоит установить скорость обучения как очень высокую? Давайте попробуем и посмотрим, что из этого выйдет:

```
learn = cnn_learner(dls, resnet34, metrics=error_rate)
learn.fine_tune(1, base_lr=0.1)
```

epoch	train_loss	valid_loss	error_loss	time
0	8.946717	47.954632	0.893775	00:20

epoch	train_loss	valid_loss	error_loss	time
0	7.231843	4.119265	0.954668	00:24

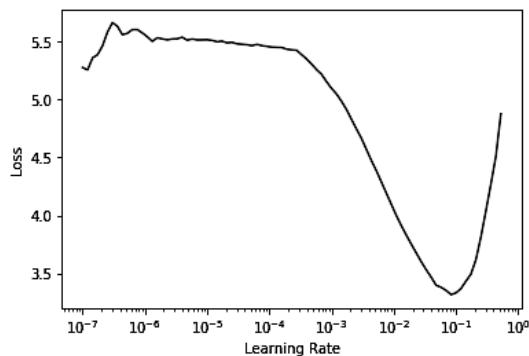
Выглядит не очень. Произошло же вот что. Оптимизатор совершил шаги в верном направлении, но перешагивал настолько далеко, что полностью перескочил область минимальных потерь. Многократное повторение этого процесса уводило его от этой области все дальше и дальше, вместо того чтобы приближать.

Что же нам делать для нахождения идеальной скорости обучения: не слишком высокой и не слишком низкой? В 2015 году исследователь Лесли Смит придумал замечательное решение, которое называл *искателем скорости обучения*. Его идея состояла в том, чтобы начать с минимально низкой скорости обучения, которая ни при каких условиях не выйдет из-под контроля. Мы применяем эту скорость к одному мини-пакету, после чего определяем потери, а затем увеличиваем ее на определенный процент (например, каждый раз удваиваем). После этого обрабатываем следующий мини-пакет, отслеживаем потери и снова удваиваем скорость обучения. Этот процесс мы повторяем до тех пор, пока уменьшение потерь не сменится их ростом. Дальнейшее увеличение скорости уже будет негативно сказываться на модели, поэтому мы выбираем ее значение, немного отступив от переломной точки назад. При этом мы советуем использовать один из двух вариантов.

- Выбрать значение на один порядок меньше, чем то, при котором были достигнуты минимальные потери (то есть минимум, деленный на 10).
- Выбрать последнюю точку, в которой потери отчетливо снижались.

В данном случае искатель скорости обучения помогает вам, вычисляя эти точки на кривой. Оба перечисленных правила обычно дают примерно одинаковое значение. В главе 1 мы не указывали скорость обучения, используя предустановленное значение из библиотеки fastai (а именно 1e-3):

```
learn = cnn_learner(dls, resnet34, metrics=error_rate)
lr_min, lr_stEEP = learn.lr_find()
```



```
print(f"Minimum/10: {lr_min:.2e}, steepest point: {lr_steepest:.2e}")

Minimum/10: 8.32e-03, steepest point: 6.31e-03
```

На этом графике видно, что в диапазоне от $1\text{e-}6$ до $1\text{e-}3$ фактически ничего не происходит, и модель не обучается. Затем потери начинают уменьшаться до момента достижения минимума, после чего начинают возрастать. Нам не нужна скорость обучения выше $1\text{e-}1$, поскольку она приведет к блужданию (можете попробовать сами), но даже $1\text{e-}1$ уже слишком высока, так как на этой стадии мы пропустили отрезок, где потери стабильно уменьшались.

Судя по графику, вырисовывается, что скорость обучения около $3\text{e-}3$ оказывается наиболее подходящей. Значит, ее мы и выберем:

```
learn = cnn_learner(dls, resnet34, metrics=error_rate)
learn.fine_tune(2, base_lr=3e-3)
```

epoch	train_loss	valid_loss	error_loss	time
0	1.071820	0.427476	0.133965	00:19
0	0.738273	0.541828	0.150880	00:24
1	0.401544	0.266623	0.081867	00:24



ЛОГАРИФМИЧЕСКИЙ МАСШТАБ

Определитель скорости обучения использует логарифмический масштаб, в связи с чем средняя точка между $1\text{e-}3$ и $1\text{e-}2$ находится между $3\text{e-}3$ и $4\text{e-}3$. Причина в том, что в первую очередь нас интересует порядок возрастания скорости обучения.

Интересно то, что искатель скорости обучения был разработан только в 2015 году, хотя разработка самих нейронных сетей уходит корнями в 1950-е. На протяжении всего этого времени нахождение оптимальной скорости обучения являлось, вероятно, важнейшей и сложнейшей задачей для практиков. Это решение не требует

использования продвинутых математических принципов, чудовищных вычислительных ресурсов, огромных датасетов или чего-либо еще, что сделало бы этот процесс недоступным для любого любопытного исследователя. Более того, Смит не был частью какой-то эксклюзивной лаборатории из Кремниевой долины, а работал простым военно-морским исследователем. Все это говорит о том, что прорывы и важнейшие работы в области глубокого обучения совершенно не требуют доступа к обширным ресурсам, связь с элитными командами или знания продвинутых математических концепций. Впереди предстоит сделать еще очень многое, для чего требуется только немного здравого смысла, креативности и упорства.

Теперь, когда мы подобрали хорошую скорость обучения, посмотрим, как можно точно настроить веса предварительно обученной модели.

Разморозка и перенос обучения

В главе 1 мы уже коротко говорили о принципе работы переноса обучения. Мы видели, что его основная идея в том, что предварительно обученная модель, которая потенциально обучалась на миллионах точек данных (например, ImageNet), может тонко настраиваться для другой задачи. Но что же это значит на самом деле?

Теперь нам известно, что сверточная нейронная сеть состоит из многих линейных слоев с функцией нелинейной активации между каждой их парой, сопровождаемых одним или несколькими заключительными слоями с функцией активации, такой как softmax, в конце. Заключительный линейный слой использует матрицу с таким количеством столбцов, чтобы размер вывода совпал с числом классов модели (предполагая, что мы выполняем классификацию).

Этот заключительный линейный слой вряд ли будет для нас полезен при выполнении тонкой настройки в условиях переноса обучения, потому что он спроектирован специально для классификации категорий в оригинальном датасете, использованном для предварительного обучения. Поэтому в процессе переноса обучения мы его удаляем и заменяем новым линейным слоем с нужным числом выводов, соответствующим нашей задаче (в данном случае будет 37 активаций).

Этот добавленный линейный слой будет иметь полностью случайные веса. Следовательно, модель до начала тонкой настройки будет иметь полностью случайные выводы. Но это вовсе не значит, что модель тоже будет абсолютно случайной. Все слои, предшествующие заключительному, были тщательно обучены для эффективного выполнения классификации изображений в целом. Как мы видели на изображениях из статьи Зейлера (Zeiler) и Фергуса (Fergus) (<https://oreil.ly/aTRwE>) в главе 1 (см. рис. 1.10–1.13), первые несколько слоев кодируют общие принципы, такие как градиенты и края, последующие же слои тоже кодируют полезные для нас принципы, например глаза и мех.

Нам нужно обучить модель таким образом, чтобы позволить ей запомнить все эти общие полезные идеи из предварительно обученной модели, использовать их

для решения конкретно нашей задачи (классификации домашних пород) и подстроить их ровно настолько, чтобы они соответствовали именно этой задаче.

Сложность в процессе тонкой настройки возникает при замене случайных весов добавленных слоев на такие, которые достигнут нужного результата в интересующей нас задаче (классификация домашних пород), не затрагивая тщательно обученные ранее веса и другие слои. Для этого мы используем простой трюк: сообщаем оптимизатору, что нужно обновить веса только в случайно добавленных заключительных слоях и совсем не менять веса в остальной части нейронной сети. Это называется *заморозкой* предварительно обученных слоев.

Когда мы создаем модель на основе предварительно обученной сети, `fastai` автоматически замораживает все предварительно обученные слои за нас. Когда мы вызываем метод `fine_tune`, `fastai` выполняет два действия.

- Обучает случайно добавленные слои в течение одной эпохи. При этом все остальные слои остаются замороженными.
- Размораживает все слои и обучает их указанное количество эпох.

Несмотря на то что это разумный предустановленный подход, вполне возможно, что для вашего датасета получится добиться лучших результатов, сделав все несколько иначе. В методе `fine_tune` есть параметры, с помощью которых можно менять его поведение, но для этого, вероятно, легче просто вызвать непосредственно внутренние методы, настроив тем самым нужное поведение. Помните, что для просмотра исходного кода метода можно использовать следующий синтаксис:

```
learn.fine_tune??
```

Попробуем сделать это сами вручную. Прежде всего, мы обучим случайно добавленные слои в течение трех эпох, используя для этого `fit_one_cycle`. Как уже говорилось в главе 1, `fit_one_cycle` предлагается использовать для обучения моделей без применения `fine_tune`. Несколько позже мы узнаем почему. Вкратце же скажем, что `fit_one_cycle` начинает обучение с низкой скоростью и постепенно повышает ее в процессе первого этапа обучения, после чего снова постепенно ее уменьшает на протяжении последнего этапа:

```
learn = cnn_learner(dls, resnet34, metrics=error_rate)
learn.fit_one_cycle(3, 3e-3)
```

epoch	train_loss	valid_loss	error_loss	time
0	1.188042	0.355024	0.102842	00:20
1	0.534234	0.302453	0.094723	00:20
2	0.325031	0.222268	0.074425	00:20

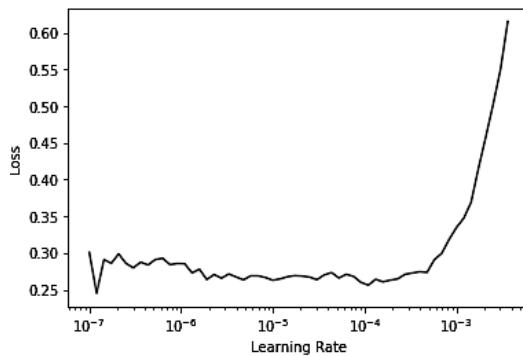
После этого мы разморозим модель:

```
learn.unfreeze()
```

А затем снова выполним `lr_find`, потому что появление дополнительных слоев для обучения и весов, которые уже обучались на протяжении трех эпох, означает, что ранее найденная скорость обучения больше не подходит:

```
learn.lr_find()
```

```
(1.0964782268274575e-05, 1.5848931980144698e-06)
```



Обратите внимание, что график здесь несколько отличается от того, где мы использовали случайные веса: здесь нет крутого спуска, указывающего, что модель обучается. Причина в том, что эта модель уже была обучена. Здесь мы наблюдаем в основном плоскую линию, после которой идет резкий подъем, и должны выбрать точку, достаточно удаленную от этого подъема: например, `le-5`. Точка с максимальным градиентом нас не интересует и должна игнорироваться.

Произведем обучение с подходящей скоростью:

```
learn.fit_one_cycle(6, lr_max=1e-5)
```

epoch	train_loss	valid_loss	error_loss	time
0	0.263579	0.217419	0.069012	00:24
1	0.253060	0.210346	0.062923	00:24
2	0.224340	0.207357	0.060217	00:24
3	0.200195	0.207244	0.061570	00:24
4	0.194269	0.200149	0.059540	00:25
5	0.173164	0.202301	0.059540	00:25

Модель немного улучшилась, но мы можем добиться большего. Самые глубокие слои предварительно обученной модели могут не нуждаться в такой же высокой скорости обучения, как первые, поэтому нам стоит использовать для них другие скорости, которые называются *дискриминативными* скоростями обучения.

Дискриминативные скорости обучения

Даже после разморозки нас все еще волнует качество предварительно обученных весов. Мы не ожидаем, что для них наилучшая скорость обучения будет настолько же высокой, как для добавленных случайных весов, даже после того, как мы обучили последние в течение нескольких эпох. Все потому, что предварительно обученные веса обучались сотни эпох и на миллионах изображений.

Кроме того, помните ли вы изображения из главы 1, показывающие, чему обучается каждый слой? Первый изучает очень простые основы, такие как края и детекторы градиентов. Это пригодится практически для любой задачи. Последующие слои изучают уже намного более сложные принципы, такие как «глаз» и «закат», которые могут быть и не нужны вашей задаче (например, если вы классифицируете модели автомобилей). Поэтому имеет смысл позволить поздним слоям выполнять тонкую подстройку быстрее, чем ранним.

В связи с этим fastai при переносе обучения по умолчанию использует дискриминативные скорости обучения. Эта техника изначально была разработана в подходе ULMFiT (Тонкая настройка универсальной языковой модели для классификации текста), примененном к переносу обучения NLP (обработке естественного языка), с которой мы познакомимся в главе 10. Как и у многих хороших идей в глубоком обучении, ее смысл достаточно прост: использовать более низкие скорости обучения для ранних слоев нейронной сети и более высокие — для поздних слоев (особенно для случайно добавленных). Эта идея основана на открытиях, сделанных Джейсоном Йосински и др. (Jason Yosinski), который в 2014 году показал, что при переносе обучения различные слои нейронной сети должны обучаться с разной скоростью (рис. 5.4).

Fastai позволяет передавать Python-объект `slice` везде, где ожидается параметр скорости обучения. Первое передаваемое значение будет скоростью обучения для ранних слоев, а второе — для заключительного. При этом величина скорости обучения для промежуточных слоев будет мультипликативно равноудаленной на протяжении всего их диапазона. Воспользуемся этим подходом для повтора предыдущего обучения, но на этот раз установим только *нижний* слой сети на скорость $1e-6$. Скорость остальных слоев при этом будет постепенно возрастать до $1e-4$. Проведем несколько эпох обучения и посмотрим, что это даст:

```
learn = cnn_learner(dls, resnet34, metrics=error_rate)
learn.fit_one_cycle(3, 3e-3)
learn.unfreeze()
learn.fit_one_cycle(12, lr_max=slice(1e-6,1e-4))
```

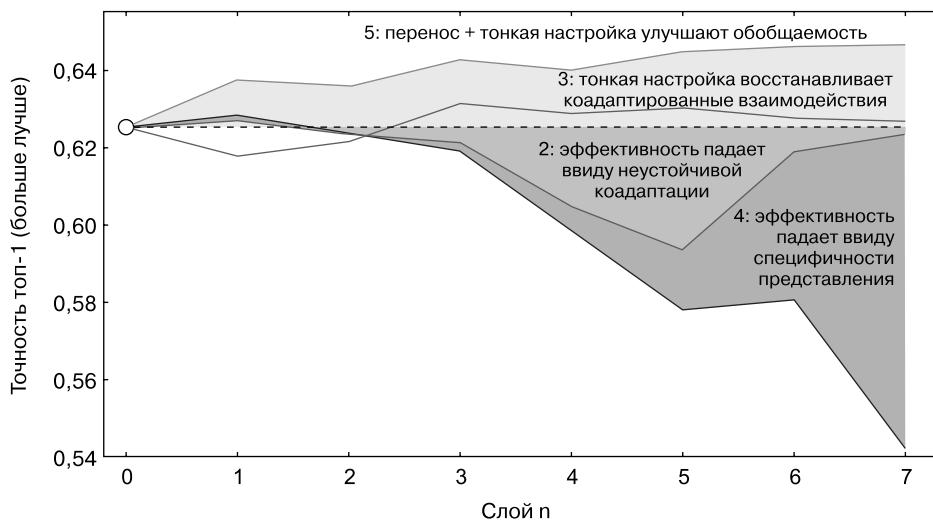


Рис. 5.4. Влияние разных слоев и методов обучения на перенос обучения
(график любезно предоставлен Джейсоном Йосински и др.)

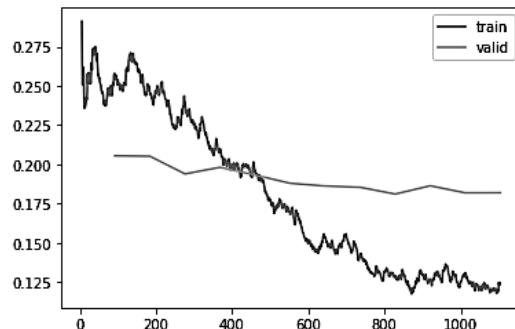
epoch	train_loss	valid_loss	error_loss	time
0	1.145300	0.345568	0.119756	00:20
1	0.533986	0.251944	0.077131	00:20
2	0.317696	0.208371	0.069012	00:20

epoch	train_loss	valid_loss	error_loss	time
0	0.257977	0.205400	0.067659	00:25
1	0.246763	0.205107	0.066306	00:25
2	0.240595	0.193848	0.062246	00:25
3	0.209988	0.198061	0.062923	00:25
4	0.194756	0.193130	0.064276	00:25
5	0.169985	0.187885	0.056157	00:25
6	0.153205	0.186145	0.058863	00:25
7	0.141480	0.185316	0.053451	00:25
8	0.128564	0.180999	0.051421	00:25
9	0.126941	0.186288	0.054127	00:25
10	0.130064	0.181764	0.054127	00:25
11	0.124281	0.181855	0.054127	00:25

Вот сейчас тонкая настройка работает отлично!

Теперь fastai может показать нам график изменения потерь для обучающей и контрольной выборки:

```
learn.recorder.plot_loss()
```



Здесь мы видим, что показатель потерь для обучающего набора продолжает снижаться. Но обратите внимание, что в конечном счете потери для контрольной выборки почти перестают уменьшаться, а иногда даже немного возрастают! Это именно та точка, в которой начинается переобучение модели. Говоря конкретнее, модель становится чрезмерно уверенной в своих прогнозах. Но это *не* означает, что она обязательно утрачивает точность. Взгляните на таблицу результатов обучения по каждой эпохе, в которой можно видеть, что точность продолжает улучшаться даже при ухудшении показателя потерь на контрольной выборке. В итоге нас интересует именно точность, или, говоря обобщенно, выбранные вами метрики, а не потери. Потери — это просто функция, которую мы представили компилятору, чтобы он помог с оптимизацией.

А теперь рассмотрим еще один немаловажный вопрос при обучении модели, а именно длительность обучения.

Выбор количества эпох

При выборе количества эпох для обучения вы часто будете сталкиваться скорее с нехваткой времени, чем с обобщением модели или ее точности. Поэтому начинать рекомендуется с обучения в течение такого количества эпох, какое вы готовы ждать. После этого рассмотрите графики потерь на обучающей и контрольной выборках, как мы это делали недавно, одновременно оценив выбранные вами метрики. Если обнаружится, что в заключительных эпохах показатели продолжают улучшаться, значит, с продолжительностью обучения модели вы еще не переборщили.

В обратном случае вы можете отчетливо увидеть, что к последним эпохам показатели метрик ухудшаются. Помните, что при оценке модели мы не только

смотрим на возможное повышение потерь, но также и на метрики. Потери на контрольной выборке в процессе обучения сначала будут расти из-за роста уверенности модели и лишь затем из-за неверного запоминания данных. Нас же фактически волнует только последнее. Как мы говорили, функция потерь используется только для оптимизатора, чтобы ему было что дифференцировать и оптимизировать, поэтому нас она на практике не интересует.

До начала использования политики обучения 1cycle модель сохраняли в конце каждой эпохи, после чего выбирали, какая из всех сохраненных моделей имела наилучшую точность. Это называется *обучением с остановкой*. Но такой способ вряд ли даст наилучший ответ, потому что во время прохождения средних эпох скорость обучения еще не снижается до уровня, при котором может быть найден лучший результат. Следовательно, если вы обнаружите, что столкнулись с переобучением, то нужно начать обучение модели с нуля, на этот раз выбрав общее число эпох, взяв за основу точку, в которой до этого были получены наилучшие результаты.

Если у вас недостаточно времени для обучения модели дополнительное число эпох, то вместо этого вы можете использовать это время для обучения большего числа параметров, то есть использовать углубленную архитектуру.

Углубленные архитектуры

Как правило, модель с большим числом параметров может моделировать данные более точно. (Конечно, тут множество подводных камней, и все будет зависеть от специфики конкретно используемой архитектуры, но сейчас этот принцип рассмотрим как универсальный.) Для большинства используемых в этой книге архитектур допустимо создавать их расширенные версии простым добавлением дополнительных слоев. Но поскольку мы хотим использовать предварительно обученные модели, то нужно убедиться, что число выбираемых слоев соответствует числу предварительно обученных.

Именно поэтому на практике архитектуры чаще представлены в небольшом количестве вариантов. Например, ResNet, которую мы используем в этой главе, предлагает варианты с 18, 34, 50, 101 и 152 слоями, предварительно обученными на ImageNet. Более глубокая (больше слоев и параметров; иногда описывается как *емкость* модели) версия ResNet всегда сможет обеспечить меньшие потери, но при этом будет склонна к переобучению, потому что для этого будет доступно больше параметров.

Обычно чем больше модель, тем выше ее способность фиксировать реальные фундаментальные связи в данных, но также фиксировать и запоминать конкретные детали отдельных изображений.

Но использование более глубокой модели потребует большего объема памяти GPU, поэтому может потребоваться уменьшить размер пакетов во избежание

ошибок нехватки памяти. Они возникают, когда вы стараетесь вместить в GPU слишком много данных, и выглядят так:

```
Cuda runtime error: out of memory
```

Подобные ошибки могут потребовать перезапуска блокнота. Решаются они простым уменьшением размера пакета, что означает передачу через модель в любой отдельно взятый момент времени меньших групп изображений. Вы можете передать в вызов нужный вам размер пакета, создав `DataLoaders` с `bs=`.

Еще один недостаток углубленных архитектур в том, что они требуют очень много времени на обучение. Одна из техник, которая может ускорить этот процесс, — *обучение со смешанной точностью*. Оно подразумевает использование в обучении менее точных чисел (*чисел половинной точности*, также называемых *fp16*) везде, где это возможно. Мы пишем книгу в начале 2020 года, когда практически все современные GPU от NVIDIA поддерживают соответствующую особую функцию под называнием *тензорные ядра*, которая может решительно ускорить обучение нейронной сети, а именно сократить его в 2–3 раза. К тому же эти числа требуют существенно меньше памяти GPU. Для активации такого варианта обучения в fastai просто добавьте `to_fp16()` после создания `Learner` (вам также потребуется импортировать модуль).

Как бы то ни было, вы не сможете узнать наперед, какой окажется наилучшая архитектура конкретно для вашей задачи, поэтому нужно просто попробовать для обучения те или иные варианты. Теперь же посмотрим, каких результатов удастся добиться от ResNet-50 при использовании смешанной точности:

```
from fastai2.callback.fp16 import *
learn = cnn_learner(dls, resnet50, metrics=error_rate).to_fp16()
learn.fine_tune(6, freeze_epochs=3)
```

epoch	train_loss	valid_loss	error_loss	time
0	1.427505	0.310554	0.098782	00:21
1	0.606785	0.302325	0.094723	00:22
2	0.409267	0.294803	0.091340	00:21

epoch	train_loss	valid_loss	error_loss	time
0	0.261121	0.274507	0.083897	00:26
1	0.296653	0.318649	0.084574	00:26
2	0.242356	0.253677	0.069012	00:26
3	0.150684	0.251438	0.065629	00:26
4	0.094997	0.239772	0.064276	00:26
5	0.061144	0.228082	0.054804	00:26

Здесь вы видите, что мы снова вернулись к использованию `fine_tune`, поскольку это очень удобно. Мы можем передать `freeze_epochs`, чтобы сообщить fastai, сколько эпох нужно проводить обучение с сохранением заморозки. Для большинства датасетов она будет автоматически менять скорость обучения.

В этом случае мы не видим отчетливого преимущества углубленной модели, и это очень важно запомнить — более глубокие модели не обязательно окажутся более эффективными для конкретного случая. Поэтому старайтесь начинать эксперименты с меньших вариантов и лишь потом переходите к масштабированию.

Резюме

В этой главе вы изучили ряд важных практических приемов, предназначенных как для подготовки данных к моделированию (подготовка размеров, получение сводной информации (`summary`) о блоке данных), так и для настройки модели (искатель скорости обучения, разморозка, дискриминативные скорости обучения, выбор количества эпох и использование углубленных архитектур). Все эти инструменты помогут вам создавать более точные модели для работы с изображениями в более сжатые сроки.

Помимо этого мы рассмотрели перекрестную энтропию. В целом на изучение и полноценное понимание этой главы стоит потратить побольше времени. Вам, конечно, вряд ли понадобится на практике самостоятельно реализовывать данную функцию с нуля, но при этом важно понимать ее вводы и выводы, потому что она (или, как мы увидим в следующей главе, ее разновидности) используется практически в каждой модели классификации. Поэтому когда вы сберегетесь заняться отладкой модели, отправить ее в работу или повысить точность, вам понадобится возможность видеть ее активации и потери, а также понимать, что вообще происходит и почему. А сделать все это должным образом у вас получится, только если вы будете хорошо понимать свою функцию потерь.

Если принцип работы перекрестной энтропии еще недостаточно для вас прояснился, то волноваться не стоит — вы его однозначно освоите. Для начала вернитесь к предыдущей главе и убедитесь, что хорошо поняли `mnist_loss`. Затем постепенно прорабатывайте ячейки блокнота для этой главы, в которых мы проходим по каждой части функции перекрестной энтропии. Увердитесь в понимании того, что реализует каждое вычисление и почему. Попробуйте сами создать несколько небольших тензоров и передать их в эти функции, рассмотрев возвращаемый результат.

Помните, что выборы, осуществляемые при реализации функции потерь перекрестной энтропии, не являются единственными возможными. Как и в случае рассмотрения регрессии, когда мы могли выбирать между среднеквадратичной ошибкой и средней абсолютной разностью (L_1), мы можем изменять детали

и здесь. Если у вас есть другие предположения относительно возможных функций, которые могли бы тут подойти, смело пробуйте их в блокноте этой главы. (Честное предупреждение: скорее всего, вы обнаружите, что модель станет медленнее обучаться и утратит точность. Причина в том, что градиент функции потери перекрестной энтропии пропорционален разности между активацией и целью, поэтому SGD всегда получает хорошо масштабированный шаг для весов.)

Вопросник

1. Почему мы сначала увеличиваем размер с помощью CPU, а затем уменьшаем его на GPU?
2. Если вы не знакомы с регулярными выражениями, найдите по ним руководство и несколько наборов задач для решения. На сайте книги предложено несколько вариантов.
3. Какими двумя способами данные чаще всего предоставляются для большинства датасетов глубокого обучения?
4. Изучите документацию к классу `L` и попробуйте использовать несколько из новых методов, которые он добавляет.
5. Ознакомьтесь с документацией для Python-модуля `pathlib` и попробуйте использовать несколько методов класса `Path`.
6. Приведите два примера того, как преобразования изображений могут ухудшить качество данных.
7. Какой метод в `fastai` позволяет просматривать данные в `DataLoaders`?
8. Какой метод в `fastai` позволяет вам отлаживать `DataBlock`?
9. Стоит ли вам отложить обучение модели до того, как данные будут тщательно очищены?
10. Из каких двух частей составляется функция перекрестной энтропии в PyTorch?
11. Какие два свойства активаций обеспечивает `softmax`? Почему это важно?
12. В каких случаях вам может потребоваться отсутствие у активаций этих двух свойств?
13. Вычислите столбцы `exp` и `softmax` на рис. 5.3 самостоятельно (то есть в электронной таблице, калькуляторе или блокноте).
14. Почему нельзя использовать `torch.where` при создании функции потерь для датасетов, в которых целевая метка может иметь более двух категорий?
15. Каково значение $\log(-2)$? Почему?

16. Каковы два основных рекомендуемых варианта выбора скорости обучения при использовании ее искателя?
17. Какие два шага выполняет метод `fine_tune`?
18. Как получить исходный код метода или функции в Jupyter Notebook?
19. Что такое дискриминативные скорости обучения?
20. Как Python-объект `slice` интерпретируется при передаче в качестве скорости обучения в `fastai`?
21. Почему обучение с остановкой считается плохим выбором при использовании подхода `1cycle`?
22. В чем разница между `resnet50` и `resnet101`?
23. Для чего используется `to_fp16`?

Дополнительные задания

1. Найдите и прочитайте работу Лесли Смита, в которой он вводит в оборот искатель скорости обучения.
2. Посмотрите, удастся ли вам повысить точность классификатора из этой главы. Какой наилучшей точности вы сможете достичь? Загляните на форумы сайта книги, чтобы сравнить свои результаты с результатами других студентов для этого же датасета.

ГЛАВА 6

Другие задачи компьютерного зрения

В предыдущей главе вы познакомились с несколькими важнейшими практическими техниками, являющимися частью обучения модели. Проработка таких моментов, как выбор скорости обучения и количества эпох, необходима для получения хороших результатов.

Здесь мы рассмотрим еще два типа задач компьютерного зрения: классификацию по нескольким меткам и регрессию. Первый вид относится к тем случаям, когда вам нужно прогнозировать для изображения более одной метки (иногда и вовсе ни одной), а второй — к тем, в которых метками является одно или несколько чисел, то есть количество вместо категории.

Мы также подробнее разберем выходные активации, цели и функции потерь.

Классификация по нескольким меткам

Классификация по нескольким меткам относится к задачам определения категорий объектов на изображениях, которые могут содержать более одного типа объекта. В таких случаях может быть как несколько соответствующих искомым классам объектов, так и ни одного.

Этот подход можно было бы очень удачно применить для нашего классификатора медведей. Одна из проблем классификатора, который мы развертывали в главе 2, заключалась в том, что при загрузке пользователем изображения, не относящегося ни к одному виду медведя, модель все равно сообщала, что это либо гризли, либо черный медведь, либо плюшевый, — у нее не было возможности предсказать, что «это вообще не медведь». На самом деле после чтения этой главы хорошо было бы вернуться к приложению для классификации изображений и попробовать повторно обучить модель, используя технику работы с несколькими метками. А после этого протестировать полученную

модель, передав ей изображение, не относящееся ни к одному из распознаваемых классов.

На практике нам редко попадались такие примеры, зато мы часто видели, что и пользователи, и разработчики жалуются на эту проблему. Оказывается, что это простое решение либо недостаточно широко понято, либо просто недооценено. Если учесть, что в работе мы чаще имеем дело с изображениями, на которых распознаваемый объект может как отсутствовать, так и присутствовать не в единственном числе, то напрашивается вывод, что классификаторы по некоторым меткам имеют более широкий спектр применения, чем их упрощенные версии с одной меткой.

Сначала посмотрим, как выглядит датасет с несколькими метками, после чего мы расскажем, как подготовить его для нашей модели. Вы увидите, что архитектура модели не изменится по сравнению с предыдущей главой и изменение коснется только функции потерь. Давайте начнем с данных.

Данные

Для нашего примера мы будем использовать датасет PASCAL, в котором на одном изображении может присутствовать более одного классифицируемого объекта.

Начнем мы, как обычно, с того, что скачаем и распакуем его:

```
from fastai.vision.all import *
path = untar_data(URLs.PASCAL_2007)
```

Этот датасет отличается от тех, что мы видели ранее, тем, что не структурирован по имени файлов или каталогам. Вместо этого он содержит CSV-файл, в котором прописано, какие метки и для каких изображений использовать. Мы можем изучить этот файл, прочитав его в Pandas DataFrame:

```
df = pd.read_csv(path/'train.csv')
df.head()
```

	fname	labels	ls_valid
0	000005.jpg	chair	True
1	000007.jpg	car	True
2	000009.jpg	horse person	True
3	000012.jpg	car	False
4	000016.jpg	bicycle	True

Как вы видите, перечень категорий каждого изображения представлен в виде строки с пробелами между этими категориями.

PANDAS И DATAFRAME

Конечно, речь не о панде. *Pandas* — это библиотека Python, используемая для управления табличными данными и данными временных рядов, а также для их анализа. Основным ее классом является *DataFrame*, который представляет собой таблицу, состоящую из строк и столбцов.

Можно получить *DataFrame* из CSV-файла, таблицы базы данных, словарей Python и многих других источников. В Jupyter *DataFrame* выводится в виде форматированной таблицы, как показано ниже.

Обращаться к его строкам и столбцам можно с помощью свойства *iloc*, как если бы это была матрица:

```
df.iloc[:,0]  
  
0      000005.jpg  
1      000007.jpg  
2      000009.jpg  
3      000012.jpg  
4      000016.jpg  
...  
5006    009954.jpg  
5007    009955.jpg  
5008    009958.jpg  
5009    009959.jpg  
5010    009961.jpg  
Name: fname, Length: 5011, dtype: object  
  
df.iloc[0,:]  
# Добавочное двоеточие нигде обязательным не считается  
# (в numpy, pytorch, pandas, и т. д.),  
# поэтому следующий вариант полностью равнозначен:  
df.iloc[0]  
  
fname    000005.jpg  
labels   chair  
is_valid True  
Name: 0, dtype: object
```

Можно также выбрать столбец по имени, обратившись напрямую к *DataFrame*:

```
df['fname']  
  
0      000005.jpg  
1      000007.jpg  
2      000009.jpg  
3      000012.jpg  
4      000016.jpg  
...  
5006    009954.jpg
```

```
5007    009955.jpg
5008    009958.jpg
5009    009959.jpg
5010    009961.jpg
Name: fname, Length: 5011, dtype: object
```

Можно создавать новые столбцы, а также производить вычисления с их помощью:

```
df1 = pd.DataFrame()
df1['a'] = [1,2,3,4]
df1
```

	a
0	1
1	2
2	3
3	4

```
df1['b'] = [10, 20, 30, 40]
df1['a'] + df1['b']

0    11
1    22
2    33
3    44
dtype: int64
```

Pandas — это быстрая и гибкая библиотека, которая есть в арсенале любого специалиста по данным, работающего в Python. К сожалению, его API несколько странноват и может запутать, поэтому нужно время, чтобы привыкнуть. Если раньше вы не использовали Pandas, мы предлагаем ознакомиться с одним из руководств. Лучше всего прочесть книгу *Python for Data Analysis* (<http://shop.oreilly.com/product/0636920050896.do>) (O'Reilly), написанную Уэсом Маккинни (Wes McKinney), создателем Pandas. В ней он рассказывает о других важных библиотеках — `matplotlib` и `NumPy`. Мы, конечно, постараемся вкратце описывать функциональность Pandas, которую будем использовать по ходу книги, но наше описание не сравнится с обстоятельным изложением темы в книге Маккинни.

Теперь, когда мы познакомились с нашими данными, займемся их подготовкой к обучению модели.

Построение DataBlock

Как преобразовать объект `DataFrame` в `DataLoaders`? Как правило, для создания объекта `DataLoaders` мы предлагаем везде по возможности использовать API блока данных, поскольку он обеспечивает отличное сочетание гибкости

и простоты. Далее на примере текущего датасета мы пошагово покажем, как построить `DataLoaders` с помощью этого API.

Как мы уже видели, в PyTorch и fastai есть два основных класса для представления обучающей и контрольной выборки, а также взаимодействия с ними.

- `Dataset` — коллекция, возвращающая кортеж с независимой и зависимой переменными для одного элемента.
- `DataLoader` — итератор, передающий поток мини-пакетов, каждый из которых состоит из пакета с независимыми переменными и пакета с зависимыми.

В добавок ко всему этому fastai предоставляет два класса для объединения обучающей и контрольной выборок.

- `Datasets` — итератор, содержащий обучающий `Dataset` и контрольный `Dataset`.
- `DataLoaders` — объект, содержащий обучающий `DataLoader` и контрольный `DataLoader`.

Поскольку `DataLoader` создается поверх `Dataset` и добавляет ему дополнительную функциональность (совмещение нескольких элементов в мини-пакет), то легче всего начать с создания и тестирования именно `Datasets`, а затем, когда он заработает, уже смотреть на `DataLoaders`.

Построение `DataBlock` мы делаем поступательно, шаг за шагом, и используем для попутной проверки данных блокнот. Это отличный способ придерживаться единого ритма в написании кода и уверенно отслеживать возможные проблемы. При этом тут же можно производить отладку, ведь возникновение любой ошибки будет вызвано только что введенным кодом.

Начнем с простейшего случая, а именно с блока данных без параметров:

```
dblock = DataBlock()
```

Мы можем создать на его основе объект `Datasets`. Единственное, что для этого нужно, — источник данных, которым в нашем случае является `DataFrame`:

```
dsets = dblock.datasets(df)
```

Он содержит датасет `train` и `valid`, к которым мы можем обратиться:

```
dsets.train[0]
```

```
(fname      008663.jpg
 labels     car person
 is_valid   False
```

```
Name: 4346, dtype: object,
fname      008663.jpg
labels     car person
is_valid   False
Name: 4346, dtype: object)
```

Как вы видите, эта команда просто возвращает строку из DataFrame дважды. Это происходит потому, что по умолчанию блок данных предполагает, что у нас есть два компонента: ввод и цель. Нам понадобится извлечь из DataFrame соответствующие поля, что можно сделать, передав блоку данных функции `get_x` и `get_y`:

```
dblock = DataBlock(get_x = lambda r: r['fname'], get_y = lambda r: r['labels'])
dsets = dblock.datasets(df)
dsets.train[0]

('005620.jpg', 'aeroplane')
```

Здесь вместо того, чтобы определять функцию привычным способом, используем ключевые слова `lambda`. Это просто сокращение, используемое для определения функции с последующим ее вызовом. Вот схожий, но более развернутый способ:

```
def get_x(r): return r['fname']
def get_y(r): return r['labels']
dblock = DataBlock(get_x = get_x, get_y = get_y)
dsets = dblock.datasets(df)
dsets.train[0]

('002549.jpg', 'tvmonitor')
```

Лямбда-функции отлично подходят для быстрого перебора, но они несовместимы с сериализацией, поэтому рекомендуем использовать более развернутый способ в случаях, когда нужно экспортить `Learner` после обучения (лямбды же отлично подойдут для экспериментов).

Здесь нам понадобится преобразовать независимую переменную в полный путь, чтобы иметь возможность открыть ее как изображение. Зависимую переменную нужно разделить пробелами (для чего в Python по умолчанию используется функция `split`), чтобы она стала списком:

```
def get_x(r): return path/'train'/r['fname']
def get_y(r): return r['labels'].split(' ')
dblock = DataBlock(get_x = get_x, get_y = get_y)
dsets = dblock.datasets(df)
dsets.train[0]

(Path('/home/sgugger/.fastai/data/pascal_2007/train/008663.jpg'),
 ['car', 'person'])
```

Чтобы открыть изображение и выполнить его превращение в тензоры, потребуется использовать преобразователи, которые нам предоставят типы блоков. Мы можем применить те же типы блоков, что и ранее, за одним исключением: `ImageBlock` будет работать отлично, так как путь, ведущий к реальному изображению, у нас есть, а вот `CategoryBlock` уже не сработает. Проблема в том, что этот блок возвращает одно целое число, а нам нужно иметь несколько меток для каждого элемента. Чтобы решить эту проблему, мы используем `MultiCategoryBlock`. Этот тип блока ожидает получения списка строк, как в нашем случае, так что протестируем его:

```
dblock = DataBlock(blocks=(ImageBlock, MultiCategoryBlock),
                   get_x = get_x, get_y = get_y)
dsets = dblock.datasets(df)
dsets.train[0]

(PILImage mode=RGB size=500x375,
 TensorMultiCategory([0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0.,
 0., 0., 0., 0., 0.]))
```

Здесь мы видим, что список категорий закодирован не так, как в случае с использованием `CategoryBlock`. В том случае у нас было одно целое число, отражающее представленную категорию на основе ее положения в словаре. В данном случае вместо него мы имеем список нулей с единицами в позициях, где представлена эта категория. Например, если во второй и четвертой позиции находится единица, это значит, что на изображении представлены второй и четвертый элементы словаря. Эта техника называется *быстрым кодированием*. Мы не можем просто взять и использовать список индексов категорий, потому что тогда каждый список будет разной длины, а PyTorch нужны тензоры, в которых все должно иметь одну длину.



ТЕРМИН: БЫСТРОЕ КОДИРОВАНИЕ

Кодирование списка целых чисел с помощью вектора из нулей с единицами в тех местах, где представлены данные.

Проверим, что представляют собой эти категории для данного примера (мы используем удобную функцию `torch.where`, которая сообщает нам все индексы, где наше условие `true` или `false`):

```
idxs = torch.where(dsets.train[0][1]==1.)[0]
dsets.train.vocab[idxs]

(#1) ['dog']
```

С помощью массивов NumPy, тензоров PyTorch и fastai-класса `L` мы можем выполнять индексацию напрямую, используя список или вектор, которые делают большую часть кода (как в данном примере) понятнее и лаконичнее.

Мы игнорировали столбец `is_valid` вплоть до текущего момента. Это означает, что `DataBlock` по умолчанию использовал случайное разделение. Чтобы выбрать

элементы для нашей контрольной выборки явно, нужно написать функцию и передать ее в `splitter` (в fastai для этого также можно использовать предопределенные функции или классы). Она получит эти элементы (в данном случае весь DataFrame) и должна будет вернуть два (или более) списка целых чисел:

```
def splitter(df):
    train = df.index[~df['is_valid']].tolist()
    valid = df.index[df['is_valid']].tolist()
    return train,valid

dblock = DataBlock(blocks=(ImageBlock, MultiCategoryBlock),
                   splitter=splitter,
                   get_x=get_x,
                   get_y=get_y)

dsets = dblock.datasets(df)
dsets.train[0]

(PILImage mode=RGB size=500x333,
 TensorMultiCategory([0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0.,
 0., 0., 0., 0., 0.]))
```

Как мы уже говорили, `DataLoader` собирает элементы из `Dataset` в мини-пакет. Это кортеж тензоров, в котором каждый тензор представляет стек элементов из этой (или соответствующей) позиции в элементе `Dataset`.

Теперь, когда мы знаем, что отдельные элементы выглядят хорошо, остается еще один шаг: нужно убедиться, что можно создать `DataLoaders`, то есть обеспечить, чтобы все элементы были одного размера. Для этого используем `RandomResizedCrop`:

```
dblock = DataBlock(blocks=(ImageBlock, MultiCategoryBlock),
                   splitter=splitter,
                   get_x=get_x,
                   get_y=get_y,
                   item_tfms = RandomResizedCrop(128, min_scale=0.35))
dls = dblock.dataloaders(df)
```

Теперь мы можем показать образец наших данных:

```
dls.show_batch(nrows=1, ncols=3)
```



Помните, если вдруг при создании `DataLoaders` из вашего `DataBlock` что-то пойдет не так или вы захотите посмотреть, что именно происходит с `DataBlock`, то используйте метод `summary`, о котором шла речь в прошлой главе.

Теперь данные готовы к обучению модели. Как мы увидим, при создании `Learner` ничего не изменится, но за кадром `fastai` выберет за нас другую функцию потерь: бинарную перекрестную энтропию.

Бинарная перекрестная энтропия

Займемся созданием `Learner`. В главе 4 мы видели, что объект `Learner` содержит четыре основных компонента: модель, объект `DataLoaders`, `Optimizer` и функцию потерь. У нас уже есть `DataLoaders`, мы можем задействовать `fastai` модели `resnet` (которые несколько позже мы научимся создавать с нуля) и знаем, как создавать оптимизатор `SGD`. Значит, остается убедиться, что мы используем подходящую функцию потерь. Для этого с помощью `cnn_learner` создадим `Learner`, чтобы можно было посмотреть его активации:

```
learn = cnn_learner(dls, resnet18)
```

Мы также видели, что модель в `Learner` обычно является объектом класса, наследующего от `nn.Module`, что ее можно вызвать с помощью круглых скобок и она вернет активации модели. Для этого нужно передать ей независимую переменную в виде мини-пакета, который мы можем взять из `DataLoader`:

```
x,y = dls.train.one_batch()
activs = learn.model(x)
activs.shape
```

Подумайте, почему у `activs` именно такая форма: размер нашего пакета 64, и требуется вычислить вероятность для каждой из 20 категорий. Вот как выглядит одна из их активаций:

```
activs[0]
tensor([ 2.0258, -1.3543,  1.4640,  1.7754, -1.2820, -5.8053,  3.6130,  0.7193,
-4.3683, -2.5001, -2.8373, -1.8037,  2.0122,  0.6189,  1.9729,  0.8999, -2.6769,
-0.3829,  1.2212,  1.6073],
device='cuda:0', grad_fn=<SelectBackward>)
```



ПОЛУЧЕНИЕ АКТИВАЦИЙ МОДЕЛИ

Умение вручную получать мини-пакет и передавать его в модель, чтобы посмотреть активации и потери, очень важно при отладке модели. Оно также весьма пригодится при обучении, потому что позволяет видеть, что именно происходит.

Активации пока не масштабированы до значений между 0 и 1, но в главе 4 мы научились это делать с помощью функции `sigmoid`. Мы также видели, как вычислять на основе этого потери, — это наша функция потерь из главы 4 с добавлением `log`, о чём мы говорили в прошлой главе:

```
def binary_cross_entropy(inputs, targets):
    inputs = inputs.sigmoid()
    return -torch.where(targets==1, inputs, 1-inputs).log().mean()
```

Обратите внимание, что так как наша зависимая переменная закодирована быстрым кодированием, мы не можем напрямую использовать `nll_loss` или `softmax` (а следовательно, и `cross_entropy`).

- `softmax`, как мы знаем, требует, чтобы все прогнозы суммировались в единицу, и склонна выталкивать одну активацию над остальными (из-за использования `exp`). Тем не менее у нас вполне может быть несколько объектов, в появлении которых на изображении мы уверены, поэтому ограничение максимальной суммы активаций до единицы не очень хорошая идея. По тем же соображениям нам может потребоваться, чтобы сумма была *меньше* единицы, если мы сочтем, что на изображении не появится *ни одна* из категорий.
- `nll_loss`, как нам известно, возвращает значение всего одной активации, соответствующей одной метке элемента. Это не имеет смысла при работе с несколькими метками.

С другой стороны, функция `binary_cross_entropy`, на деле являющаяся `mnist_loss`, совмещенной с `log`, дает то, что нужно, благодаря магии поэлементных операций PyTorch. Каждая активация будет сопоставлена с каждой целью для каждого столбца, поэтому эта функция вполне заработает для всех столбцов без нашего вмешательства.



СЛОВО ДЖЕРЕМИ

Что мне действительно нравится при работе с библиотеками вроде PyTorch, где используется уширение (broadcasting) и поэлементные операции, так это то, что можно написать код, одинаково работающий как для одного элемента, так и для их пакета. Функция `binary_cross_entropy` является отличным тому примером. При использовании этих операций исчезает необходимость в ручном прописывании циклов. Мы можем доверить этот процесс PyTorch, который выполняет их в соответствии с рангом тензора, с которым мы работаем.

PyTorch уже предоставляет эту функцию. Более того, он предоставляет ряд ее вариантов, причем с весьма неоднозначными именами.

`F.binary_cross_entropy` и ее модульный эквивалент `nn.BCELoss` вычисляют перекрестную энтропию для цели, закодированной быстро, но не включают началь-

ную `sigmoid`. Обычно для таких целей вам нужна `F.binary_cross_entropy_with_logits` (или `nn.BCEWithLogitsLoss`), которая выполняет и сигмоиду, и бинарную перекрестную энтропию, как в предыдущем примере.

Их эквивалентом для датасетов изображений с одиночными метками (типа MNIST или набора с домашними животными), где цель кодируется как одно целое число, будет `F.nll_loss` или `nn.NLLLoss` для версии без начальной softmax и `F.cross_entropy` или `nn.CrossEntropyLoss` для версии с ней.

Поскольку наша цель закодирована быстро, мы используем `BCEWithLogitsLoss`:

```
loss_func = nn.BCEWithLogitsLoss()
loss = loss_func(activs, y)
loss

tensor(1.0082, device='cuda:5', grad_fn=<BinaryCrossEntropyWithLogitsBackward>)
```

Нам не нужно просить fastai использовать эту функцию потерь (хотя при желании мы можем это сделать), потому что она будет выбрана автоматически. fastai знает, что в `DataLoaders` много категориальных меток, поэтому по умолчанию выбирает `nn.BCEWithLogitsLoss`.

Еще одно изменение — это используемая метрика. Поскольку задача подразумевает много меток, мы не можем использовать функцию точности. Почему? Что ж, точность сравнивала наши выводы с целями следующим образом:

```
def accuracy(inp, targ, axis=-1):
    "Compute accuracy with `targ` when `pred` is bs * n_classes"
    pred = inp.argmax(dim=axis)
    return (pred == targ).float().mean()
```

Был спрогнозирован класс с максимальной активацией (именно за это отвечает `argmax`). Здесь эта функция не сработает, потому что у нас может быть более одного прогноза для изображения. После применения сигмоиды к активациям (с целью приведения их к значению между 0 и 1) нам нужно решить, какие из них нули, а какие — единицы, выбрав порог. Каждое значение выше этого порога будет считаться единицей, а значения ниже будут относиться к нулю:

```
def accuracy_multi(inp, targ, thresh=0.5, sigmoid=True):
    "Compute accuracy when `inp` and `targ` are the same size."
    if sigmoid: inp = inp.sigmoid()
    return ((inp>thresh)==targ.bool()).float().mean()
```

Если мы передадим `accuracy_multi` в качестве метрики напрямую, то она будет использовать для `threshold` значение по умолчанию, то есть `0.5`. Нам же может понадобиться подстроить это значение и создать новую версию `accuracy_multi` с другим предустановленным значением. Здесь поможет Python-функция `partial`. Она позволяет связывать функцию с аргументами или именованными

аргументами, чтобы она при вызове всегда включала именно эти аргументы. Например, вот простая функция, получающая два аргумента:

```
def say_hello(name, say_what="Hello"): return f"{say_what} {name}."  
say_hello('Jeremy'), say_hello('Jeremy', 'Ahoy!')  
  
('Hello Jeremy.', 'Ahoy! Jeremy.')
```

Мы можем переключиться на французскую версию этой функции, используя `partial`:

```
f = partial(say_hello, say_what="Bonjour")  
f("Jeremy"), f("Sylvain")  
  
('Bonjour Jeremy.', 'Bonjour Sylvain.')
```

Теперь мы можем начать обучать модель. Давайте попробуем установить для нашей метрики порог точности 0.2:

```
learn = cnn_learner(dls, resnet50, metrics=partial(accuracy_multi, thresh=0.2))  
learn.fine_tune(3, base_lr=3e-3, freeze_epochs=4)
```

epoch	train_loss	valid_loss	accuracy_multi	time
0	0.903610	0.659728	0.263068	00:07
1	0.724266	0.346332	0.525458	00:07
2	0.415597	0.125662	0.937590	00:07
3	0.254987	0.116880	0.945418	00:07

epoch	train_loss	valid_loss	accuracy_multi	time
0	0.123872	0.132634	0.940179	00:08
1	0.112387	0.113758	0.949343	00:08
2	0.092151	0.104368	0.951195	00:08

Выбор порога очень важен. Если установить слишком низкий, то у вас зачастую не будет получаться выбирать верно размеченные объекты. Это можно увидеть, изменив метрику и затем вызвав `validate`, которая вернет контрольные потери и метрики:

```
learn.metrics = partial(accuracy_multi, thresh=0.1)  
learn.validate()  
  
(#2) [0.10436797887086868, 0.93057781457901]
```

Если же указать слишком высокий порог, то вы будете выбирать только те объекты, в которых модель очень уверена:

```
learn.metrics = partial(accuracy_multi, thresh=0.99)
learn.validate()

(#2) [0.10436797887086868, 0.9416930675506592]
```

Наилучшее значение порога можно найти, попробовав несколько вариантов и оценив результаты. Это делается существенно быстрее, если мы извлекаем прогнозы всего один раз:

```
preds,targs = learn.get_preds()
```

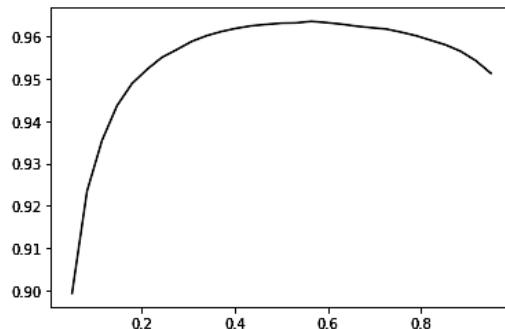
Затем можно вызвать метрику напрямую. Обратите внимание, что по умолчанию `get_preds` применяет выходную функцию активации (в этом случае сигмоиду) автоматически, поэтому нам нужно дать `accuracy_multi` команду не применять ее:

```
accuracy_multi(preds, targs, thresh=0.9, sigmoid=False)

TensorMultiCategory(0.9554)
```

Теперь можно использовать этот подход для определения наилучшего порогового значения:

```
xs = torch.linspace(0.05,0.95,29)
accs = [accuracy_multi(preds, targs, thresh=i, sigmoid=False) for i in xs]
plt.plot(xs,accs);
```



В данном случае для выбора гиперпараметра (порога) мы используем контрольный набор данных, для чего этот набор и предназначен. Иногда студентов беспокоит, что мы можем таким образом допустить *переобучение* в отношении контрольной выборки, поскольку поочередно перепробуем много значений в поиске наилучшего. Тем не менее, как вы сами видите на этом графике, измене-

нение порога в данном случае ведет к гладкой кривой, таким образом, мы точно не выбираем неподходящий выброс (*outlier*). Это хороший пример ситуации, в которой нужно внимательно различать разницу между теорией (не пробовать множество значений гиперпараметра, чтобы не переобучить модель в отношении контрольной выборки) и практикой (если отношение гладкое, то так делать можно).

На этом заканчивается часть главы, посвященная классификации по нескольким меткам. Дальше мы переходим к рассмотрению задачи регрессии.

Регрессия

Модели глубокого обучения очень удобно распределять по областям, например компьютерное зрение, NLP и т. д. Вообще, именно так и классифицирует свои приложения fastai, так как подобный подход привычнее для людей.

Но здесь есть кое-что еще. Модель определяется ее независимой и зависимой переменными наряду с функцией потерь. Это означает, что фактически есть гораздо более широкий спектр моделей, чем можно представить по областям их применения. Возможно, используемая независимая переменная является изображением, а зависимая — текстом (например, при генерации надписи из изображения). Или же наоборот, текстом представлена независимая переменная, а изображением — зависимая (например, при генерации изображения из надписи, что вполне по силам глубокому обучению). Также возможен вариант, что у нас есть изображения, тексты и табличные данные в виде независимых переменных, на основе которых мы стараемся прогнозировать приобретение товаров... Возможности действительно безграничны.

Чтобы выйти за рамки фиксированных вариантов применения и начать создавать собственные новаторские решения для передовых задач, важно хорошо разобраться в API блока данных (и, может быть, в API среднего уровня, с которым мы познакомимся позднее). В качестве примера рассмотрим задачу *регрессии изображений*. Она относится к обучению на датасете, в котором независимая переменная является изображением, а зависимая — одним или несколькими числами с плавающей запятой. Мы нередко видим, что люди используют регрессию изображений как совершенно отдельное приложение, но дальше вы поймете, что мы можем задействовать ее просто как еще одну CNN поверх API блока данных.

Мы собираемся сразу перейти к достаточно хитрому варианту регрессии изображений, потому что уверены: вы к этому готовы! Здесь мы реализуем модель определения ключевых точек. *Ключевая точка* означает конкретное место на изображении, в нашем случае мы будем использовать изображения людей и искать центр лиц. Это значит, что мы будем прогнозировать для каждого изображения *два* значения: строку и столбец, определяющие центр лица.

Сборка данных

Для этого раздела мы будем использовать датасет Biwi Kinect Head Pose (<https://oreil.ly/-4cO->). Начнем, как обычно, с его загрузки:

```
path = untar_data(URLs.BIWI_HEAD_POSE)
```

Посмотрим, что у нас тут.

```
path.ls()
```

```
(#50) [Path('13.obj'), Path('07.obj'), Path('06.obj'), Path('13'), Path('10'), Path('02'), Path('11'), Path('01'), Path('20.obj'), Path('17')...]
```

Здесь мы имеем 24 каталога с номерами от 01 до 24 (они соответствуют фотографиям разных людей) и связанный с каждым из них файл *.obj* (здесь они нам не понадобятся). Давайте заглянем в один из каталогов:

```
(path/'01').ls()
```

```
(#1000) [Path('01/frame_00281_pose.txt'), Path('01/frame_00078_pose.txt'), Path('01/frame_00349_rgb.jpg'), Path('01/frame_00304_pose.txt'), Path('01/frame_00207_pose.txt'), Path('01/frame_00116_rgb.jpg'), Path('01/frame_00084_rgb.jpg'), Path('01/frame_00070_rgb.jpg'), Path('01/frame_00125_pose.txt'), Path('01/frame_00324_rgb.jpg')...]
```

В подкаталогах расположены разные фреймы. Каждый из них содержит изображение (_rgb.jpg) и файл пространственного расположения (_pose.txt). Мы можем легко получить все файлы изображений рекурсивно с помощью `get_image_files`, а затем написать функцию, преобразующую имя файла изображения в ассоциированный с ним pose-файл:

```
img_files = get_image_files(path)
def img2pose(x): return Path(f'{str(x)[:-7]}pose.txt')
img2pose(img_files[0])
Path('13/frame_00349_pose.txt')
```

Взглянем на первое изображение:

```
im = PILImage.create(img_files[0])
im.shape
(480, 640)
im.to_thumb(160)
```



На сайте датасета Biwi (<https://oreil.ly/wHL28>) разъяснен формат текстового pose-файла, в котором указывается место расположения центра головы. Нас эти подробности не интересуют, поэтому мы просто покажем вам функцию для извлечения центральной точки головы:

```
cal = np.genfromtxt(path/'01'/'rgb.cal', skip_footer=6)
def get_ctr(f):
    ctr = np.genfromtxt(img2pose(f), skip_header=3)
    c1 = ctr[0] * cal[0][0]/ctr[2] + cal[0][2]
    c2 = ctr[1] * cal[1][1]/ctr[2] + cal[1][2]
    return tensor([c1,c2])
```

Данная функция возвращает координаты в виде тензора из двух элементов:

```
get_ctr(img_files[0])
tensor([384.6370, 259.4787])
```

Мы можем передать эту функцию в `DataBlock` как `get_y`, поскольку она отвечает за разметку каждого элемента. Чтобы немного ускорить обучение, мы также уменьшим изображения до половины их вводного размера.

Важно отметить, что в данном случае нам не следует использовать случайное разделение датасета. В нем на разных изображениях встречаются одни и те же люди, а нам нужно обеспечить, чтобы модель обобщалась до распознавания людей, которых она не видела. Каждый каталог датасета содержит изображения для одного человека. Следовательно, мы можем создать функцию разделения, возвращающую `true` только для одного человека, что приведет к формированию контрольной выборки, содержащей только его изображения.

Еще одно, последнее, отличие от предыдущих примеров блока данных в том, что вторым блоком является `PointBlock`. Таким образом мы сообщаем fastai, что метки представляют координаты. В этом случае она будет знать, что для этих координат нужно выполнить точно такую же аугментацию, что и для изображений:

```
biwi = DataBlock(
    blocks=(ImageBlock, PointBlock),
    get_items=get_image_files,
    get_y=get_ctr,
    splitter=FuncSplitter(lambda o: o.parent.name=='13'),
    batch_tfms=[*aug_transforms(size=(240,320)),
               Normalize.from_stats(*imagenet_stats)])
)
```

ТОЧКИ И АУГМЕНТАЦИЯ ДАННЫХ



Мы не знаем других библиотек (кроме fastai), которые бы автоматически и правильно применяли аугментацию данных к координатам. Поэтому если вы работаете с другой библиотекой, то для подобных задач вам может потребоваться отключить аугментацию.

Прежде чем приступить к моделированию, нам нужно взглянуть на наши данные и убедиться, что они в порядке:

```
dls = biwi.dataloaders(path)
dls.show_batch(max_n=9, figsize=(8,6))
```



Вполне неплохо! Помимо визуальной оценки пакета также полезно посмотреть на лежащие в его основе тензоры (это поможет лучше представить, как ваша модель видит данные):

```
xb,yb = dls.one_batch()
xb.shape,yb.shape

(torch.Size([64, 3, 240, 320]), torch.Size([64, 1, 2]))
```

Убедитесь, что понимаете, *почему* у наших мини-пакетов именно такие формы.

Вот пример одной строки из зависимой переменной:

```
yb[0]

tensor([[0.0111, 0.1810]], device='cuda:5')
```

Как вы видите, нам не пришлось использовать отдельное приложение *регрессии изображений*. Все, что нам понадобилось, — это разметить данные и сообщить fastai, какие их виды представляют независимая и зависимая переменные.

То же касается и создания Learner. Мы используем ту же функцию, что и ранее, только с одним новым параметром, после чего будем готовы переходить к обучению модели.

Обучение модели

Как обычно, мы можем создать `Learner` с помощью `cnn_learner`. Помните, как в главе 1 мы использовали `y_range`, чтобы сообщить `fastai` диапазон наших целей? То же самое мы сделаем здесь (координаты в `fastai` и PyTorch всегда масштабируются до диапазона между -1 и $+1$):

```
learn = cnn_learner(dls, resnet18, y_range=(-1,1))
```

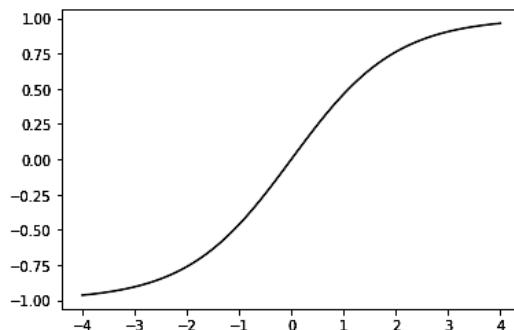
`y_range` реализуется в `fastai` с помощью `sigmoid_range`, которая определяется так:

```
def sigmoid_range(x, lo, hi): return torch.sigmoid(x) * (hi-lo) + lo
```

Если `y_range` определена, то она устанавливается как заключительный слой модели. Задумайтесь на минуту, что делает эта функция и почему она вынуждает модель выводить активации в диапазоне (lo, hi) .

Выглядит она так:

```
plot_function(partial(sigmoid_range,lo=-1,hi=1), min=-4, max=4)
```



Мы не определяли функцию потерь, в связи с чем получим ту, которую `fastai` выберет по умолчанию. Посмотрим, что она выбрала для этого случая:

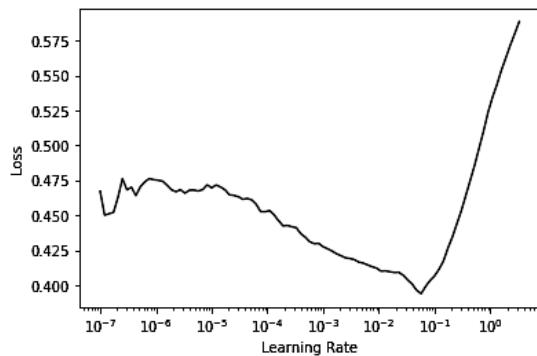
```
dls.loss_func
FlattenedLoss of MSELoss()
```

Выглядит разумно, поскольку при использовании координат в качестве зависимой переменной мы в большинстве случаев будем склонны прогнозировать наиболее близкие результаты. Именно это, по сути, и делает `MSELoss` (функция среднеквадратичной ошибки). Если вы хотите использовать другую функцию потерь, то можете передать ее в `cnn_learner`, используя параметр `loss_func`.

Обратите внимание, что мы не указывали метрики, потому что для подобной задачи MSE уже является полезной метрикой (хотя она все же станет более понятна, если мы извлечем из нее квадратный корень).

Мы можем выбрать хорошую скорость обучения с помощью уже знакомого нам метода поиска:

```
learn.lr_find()
```



Попробуем установить LR как 2e-2:

```
lr = 2e-2
learn.fit_one_cycle(5, lr)
```

epoch	train_loss	valid_loss	time
0	0.045840	0.012957	00:36
1	0.006369	0.001853	00:36
2	0.003000	0.000496	00:37
3	0.001963	0.000360	00:37
4	0.001584	0.000116	00:36

Обычно в данном случае потери составляют около 0,0001, что соответствует следующей средней ошибке прогноза координат:

```
math.sqrt(0.0001)
```

0.01

Выглядит очень точно! Но важно посмотреть на результаты с помощью `Learner.show_results`. В левой части — фактические (*истинные*) координаты, а в правой — прогнозы модели:

```
learn.show_results(ds_idx=1, max_n=3, figsize=(6,8))
```



Это удивительно, что спустя всего несколько минут вычислений мы создали настолько точную модель определения ключевых точек без использования каких-либо специализированных приложений. В этом мощь построения гибких API и использования переноса обучения. Особенно поражает, что мы смогли использовать перенос обучения настолько эффективно, что нам удалось тонко настроить модель классификации изображений на выполнение регрессии, задачи, сильно отличающейся от той, для которой она обучалась изначально.

Резюме

В задачах, которые на первый взгляд различаются (классификация по одной метке, классификация по нескольким меткам и регрессия), мы в результате используем ту же модель, просто с другим количеством выводов. Изменяется лишь функция потерь, в связи с чем важно дважды проверять, используете ли вы именно правильный ее вариант для конкретной задачи.

fastai автоматически старается выбирать функцию потерь на основе используемых вами данных, но если вы строите `DataLoaders` с помощью чистого PyTorch, то нужно хорошенько подумать над выбором подходящей функции и помнить, что, скорее всего, выбор будет следующим:

- `nn.CrossEntropyLoss` для классификации по одной метке;
- `nn.BCEWithLogitsLoss` для классификации по нескольким меткам;
- `nn.MSELoss` для регрессии.

Вопросник

1. Как можно использовать классификацию по нескольким меткам для повышения удобства и качества работы классификатора медведей?
2. Как мы кодируем зависимую переменную в задаче классификации по нескольким меткам?
3. Как обращаться к строкам и столбцам DataFrame, если рассматривать его как матрицу?
4. Как извлечь из DataFrame столбец по его имени?
5. В чем отличие `Dataset` от `DataLoader`?
6. Что обычно содержит объект `Datasets`?
7. Что обычно содержит объект `DataLoaders`?
8. Каково назначение `lambda` в Python?
9. С помощью каких методов настраивается создание независимых и зависимых переменных с помощью API блока данных?
10. Почему softmax не подходит в качестве функции активации при использовании цели, закодированной унитарно?
11. Почему `nn.l1_loss` не подходит в качестве функции потерь при использовании цели, закодированной унитарно?
12. В чем разница между `nn.BCELoss` и `nn.BCEWithLogitsLoss`?
13. Почему мы не можем использовать обычную метрику точности в задаче классификации по нескольким меткам?
14. Когда допускается настройка параметра с помощью контрольной выборки?
15. Как в fastai реализуется `y_range`? (Попробуйте реализовать ее самостоятельно и протестировать не подглядывая.)
16. Что подразумевает задача по регрессии? Какую функцию потерь следует для нее использовать?
17. Что нужно сделать, чтобы fastai применила одинаковую аугментацию данных к вводным изображениям и координатам целевой точки?

Дополнительные задания

1. Прочтите руководство по Pandas DataFrames и поэкспериментируйте с несколькими методами, которые покажутся вам наиболее интересными. Рекомендованные учебные материалы доступны на сайте книги.
2. Обучите повторно классификатор медведей с помощью классификации по нескольким меткам. Попробуйте добиться от него эффективной работы с изображениями, где нет медведей, включая показ соответствующей информации в веб-приложении. Попробуйте изображение с двумя видами медведей. Проверьте, изменилась ли точность из-за использования классификации по нескольким меткам для датасета с одиночными метками.

ГЛАВА 7

Обучение современной модели

Эта глава познакомит с более продвинутыми техниками обучения модели классификации изображений и получения самых качественных результатов. Но пока ее можно пропустить, если вы хотите сначала побольше узнать о других применениях глубокого обучения, — знание материалов этой главы не потребуется для понимания последующих.

Здесь же мы познакомимся с нормализацией, освоим мощную технику аугментации данных под называнием Mixup, изучим прогрессивный подход к изменению размеров и узнаем, что такое аугментация времени тестирования. Чтобы все это показать, мы будем обучать модель с нуля (без использования переноса обучения) с помощью подмножества ImageNet под названием Imagenette (<https://oreil.ly/1uj3x>). Оно содержит десять сильно отличающихся категорий из оригинального датасета ImageNet, что позволяет ускорить обучение, когда требуется поэкспериментировать.

В данном случае будет намного сложнее добиться отличных результатов, чем с предыдущими датасетами, потому что мы будем использовать полноцветные и полноразмерные изображения, представленные фотографиями объектов разных размеров, в разных положениях, с разным освещением и т. д. Здесь мы познакомим вас с важными техниками для получения максимальной отдачи от наших датасетов, особенно при обучении с нуля или при использовании переноса обучения для подготовки модели к датасету,циальному от использованного в предварительном обучении.

Imagenette

Когда компания fast.ai только появилась, для построения и тестирования моделей компьютерного зрения использовались три основных датасета:

- *ImageNet* — 1,3 миллиона изображений разных размеров с диагональю около 500 пикселей, разделенных на 1000 категорий. На обучение с таким датасетом уходило до нескольких дней.

- *MNIST* – 50 000 рукописных цифр в оттенках серого размером 28×28 пикселей.
- *CIFAR10* – 60 000 цветных изображений размером 32×32 пикселя, разделенных на десять классов.

Проблема была в том, что датасеты меньших размеров не могли эффективно обобщаться на огромный датасет ImageNet. Решения, которые бы отлично работали на ImageNet, должны были разрабатываться и обучаться именно на ImageNet. В результате многие люди стали верить, что полноценно участвовать в разработке алгоритмов классификации изображений могут только исследователи с огромными вычислительными ресурсами.

Мы же решили, что это совсем не так. Нам не встречалось исследование, которое доказывало бы, что ImageNet имеет «золотой» размер и что нельзя разработать другие датасеты, позволяющие получать полезные наработки. Поэтому мы решили создать новый датасет, на котором исследователи могли бы быстро и экономично тестировать свои алгоритмы и который также позволял бы в дальнейшем применять полученные наработки на полноценном датасете ImageNet.

Примерно спустя три часа был создан Imagenette. Мы выбрали десять классов из ImageNet, которые сильно отличались друг от друга. Мы надеялись, что сможем быстро и экономно создать классификатор, способный распознавать эти классы. После этого мы попробовали несколько алгоритмических приемов, оценив их влияние на Imagenette. Выяснилось, что некоторые из них работали очень даже хорошо, и мы протестировали их также на ImageNet, в результате чего убедились в их работоспособности и для этого датасета!

Здесь подразумевается важный нюанс: выбранный или полученный датасет не обязательно будет таким, какой вам нужен. Если говорить точнее, то такой датасет вряд ли подойдет для разработки и прототипирования. Стремитесь к скорости итерации, не превышающей пары минут, — то есть при желании опробовать появившуюся идею у вас должна быть возможность обучить модель и оценить ее результат в течение пары минут. Если подобный эксперимент занимает больше времени, подумайте о том, как можно сократить датасет или упростить модель, чтобы добиться увеличения скорости эксперимента. Чем больше экспериментов, тем лучше!

Начнем работать с нашим новым датасетом:

```
from fastai.vision.all import *
path = untar_data(URLs.IMGNETTE)
```

Сначала мы определим его в объект `DataLoaders`, используя метод *подготовки размера*, с которым мы познакомились в главе 5:

```
dblock = DataBlock(blocks=(ImageBlock(), CategoryBlock()),
                   get_items=get_image_files,
                   get_y=parent_label,
                   item_tfms=Resize(460),
                   batch_tfms=aug_transforms(size=224, min_scale=0.75))
dls = dblock.dataloaders(path, bs=64)
```

Затем произведем первичное обучение, результаты которого будем считать базовыми:

```
model = xresnet50()
learn = Learner(dls, model, loss_func=CrossEntropyLossFlat(), metrics=accuracy)
learn.fit_one_cycle(5, 3e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	1.583403	2.064317	0.401792	01:03
1	1.208877	1.260106	0.601568	01:02
2	0.925265	1.036154	0.664302	01:03
3	0.730190	0.700906	0.777819	01:03
4	0.585707	0.541810	0.825243	01:03

Получились неплохие базовые показатели, поскольку мы не используем предварительно обученную модель, но мы можем и лучше. При работе с моделями, обуляемыми с нуля или тонко настраиваемыми после совершенно другого датасета предварительного обучения, важно использовать некоторые дополнительные техники. Оставшаяся часть главы будет посвящена ряду ключевых подходов, которые вам однозначно пригодятся. Первый из них — *нормализация* данных.

Нормализация

При обучении модели очень удобно, если входные данные *нормализованы*, то есть имеют среднее 0 и стандартное отклонение 1. Но многие изображения и библиотеки компьютерного зрения используют для пикселей значения от 0 до 255 или между 0 и 1. В любом из этих случаев среднее ваших данных не будет иметь значение ноль, а стандартное отклонение не будет равным единице.

Возьмем пакет данных и рассмотрим их значения, взяв среднее по всем осям, кроме оси канала, то есть оси 1:

```
x,y = dls.one_batch()
x.mean(dim=[0,2,3]),x.std(dim=[0,2,3])

(TensorImage([0.4842, 0.4711, 0.4511], device='cuda:5'),
 TensorImage([0.2873, 0.2893, 0.3110], device='cuda:5'))
```

Как и ожидалось, среднее и стандартное отклонение далеки от нужных нам значений. К счастью, в fastai для нормализации данных есть удобная трансформация `Normalize`. Она применяется ко всему мини-пакету сразу, поэтому вы можете добавить ее в раздел `batch_tfms` блока данных. Этой трансформации нужно передать среднее и стандартное отклонение, которые вы хотите использовать. В fastai изначально определено стандартное среднее и стандартное отклонение датасета ImageNet. (Если вы не передадите в `Normalize` никакие статистики, то fastai автоматически вычислит их на основе одного пакета ваших данных.)

Добавим эту трансформацию (используя `imagenet_stats`, так как `Imagenette` является подмножеством ImageNet) и посмотрим на один из пакетов:

```
def get_dls(bs, size):
    dblock = DataBlock(blocks=(ImageBlock, CategoryBlock),
                      get_items=get_image_files,
                      get_y=parent_label,
                      item_tfms=Resize(460),
                      batch_tfms=[*aug_transforms(size=size, min_scale=0.75),
                                  Normalize.from_stats(*imagenet_stats)])
    return dblock.dataloaders(path, bs=bs)

dls = get_dls(64, 224)

x,y = dls.one_batch()
x.mean(dim=[0,2,3]),x.std(dim=[0,2,3])

(TensorImage([-0.0787, 0.0525, 0.2136], device='cuda:5'),
 TensorImage([1.2330, 1.2112, 1.3031], device='cuda:5'))
```

Проверим, какой эффект это оказало на обучение модели:

```
model = xresnet50()
learn = Learner(dls, model, loss_func=CrossEntropyLossFlat(), metrics=accuracy)
learn.fit_one_cycle(5, 3e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	1.632865	2.250024	0.391337	01:02
1	1.294041	1.579932	0.517177	01:02
2	0.960535	1.069164	0.657207	01:04
3	0.730220	0.767433	0.771845	01:05
4	0.577889	0.550673	0.824496	01:06

Несмотря на то что здесь она помогла не сильно, нормализация становится особенно важной при использовании предварительно обученных моделей. Такие

модели знают, как работать только с теми данными, на которых обучались. Если среднее значение пикселя в этих данных было ноль, а у ваших данных его минимальное возможное значение ноль, тогда модель увидит совсем не то, что от нее ожидается.

Это означает, что при дистрибуции модели вам нужно также прилагать статистики, использованные для нормализации, поскольку тем, кто будет использовать ее для переноса обучения или вывода, понадобится задействовать те же самые статистики. Точно так же, если вы используете модель, которую обучали до вас, обязательно выясните, какие статистики нормализации применялись, и сопоставьте их.

В предыдущих главах нам не требовалось обрабатывать нормализацию, потому что при использовании предварительно обученной модели через `cnn_learner` библиотека fastai автоматически добавляет подходящую трансформацию `Normalize`. Эта модель предварительно обучалась при использовании в `Normalize` конкретных статистик (обычно они берутся из датасета ImageNet), поэтому библиотека может подставить их за вас. Обратите внимание, что это применимо только к предварительно обученным моделям, в связи с чем здесь, при обучении модели с нуля, нам нужно добавить эту информацию вручную.

Все виды обучения вплоть до текущего момента выполнялись нами при размере 224 пикселей. На деле же мы можем начать с меньшего размера, а затем постепенно его увеличивать. Эта техника называется *прогрессивным изменением размера*.

Прогрессивное изменение размера

Когда fast.ai совместно со своей командой студентов победили на соревновании DAWN Bench в 2018-м (<https://oreil.ly/16tar>), они, помимо прочего, использовали одно важное, но при этом очень простое нововведение: начать обучение с маленьких изображений, а закончить большими. Выполнение большинства эпох обучением на небольших изображениях существенно ускоряет весь процесс. Завершение же при использовании крупных изображений существенно повышает итоговую точность. Этот подход называется *прогрессивным изменением размера* (*progressive resizing*).



ТЕРМИН: ПРОГРЕССИВНОЕ ИЗМЕНЕНИЕ РАЗМЕРА

Постепенное увеличение размера изображений, используемых в процессе обучения.

Как мы уже видели, признаки, изучаемые сверточными нейронными сетями, ни в каком смысле не зависят от размера изображения: ранние слои находят

такие основы, как края и градиенты, а более поздние слои могут находить уже, к примеру, носы и закаты. Поэтому когда мы изменяем размер изображения в середине обучения, это не означает, что нам нужно находить для нашей модели полностью другие параметры.

Но при этом также очевидно, что между маленькими и большими изображениями есть отличия, поэтому нам не следует ожидать, что наша модель продолжит так же отлично работать, если мы не внесем никаких изменений. Ничего не напоминает? Когда мы разработали эту идею, сразу возникли ассоциации с переносом обучения. Мы стараемся добиться, чтобы наша модель научилась делать что-то не совсем так, как была обучена ранее. Из этого следует, что у нас должна быть возможность использовать после изменения размеров изображений метод `fine_tune`.

У описываемой техники есть и еще одно преимущество: она представляет собой дополнительный вид аугментации. Это значит, что при обучении с изменением размеров изображений в результате можно добиться лучшей генерализации моделей.

Для реализации данной методики удобнее всего сначала создать функцию `get_dls`, получающую размер изображения и размер пакета, как мы делали в предыдущем разделе, и возвращающую `DataLoaders`.

Теперь можно создать `DataLoaders` с маленьким размером и использовать `fit_one_cycle` обычным способом, выполнив обучение на протяжении меньшего числа эпох, чем в стандартном случае.

```
dls = get_dls(128, 128)
learn = Learner(dls, xresnet50(), loss_func=CrossEntropyLossFlat(),
                 metrics=accuracy)
learn.fit_one_cycle(4, 3e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	1.902943	2.447006	0.401419	00:30
1	1.315203	1.572992	0.525765	00:30
2	1.001199	0.767886	0.759149	00:30
3	0.765864	0.665562	0.797984	00:30

А затем заменить `DataLoaders` в `Learner` и произвести тонкую настройку:

```
learn.dls = get_dls(64, 224)
learn.fine_tune(5, 1e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	0.985213	1.654063	0.565721	01:06

epoch	train_loss	valid_loss	accuracy	time
0	0.706869	0.689622	0.784541	01:07
1	0.739217	0.928541	0.712472	01:07
2	0.629462	0.788906	0.764003	01:07
3	0.491912	0.502622	0.836445	01:06
4	0.414880	0.431332	0.863331	01:06

Как вы видите, показатели эффективности модели существенно выросли, а начальное обучение на маленьких изображениях существенно сократило продолжительность каждой эпохи.

Вы можете повторять этот процесс, увеличивая размер и обучая модель на протяжении дополнительного числа эпох на изображениях любого размера, но естественно, вы не получите никакой пользы, используя изображения большего размера, чем их исходный вариант на диске.

Обратите внимание, что при переносе обучения прогрессивное изменение размера может и навредить качеству модели. Это, скорее всего, произойдет, если используемая предварительно обученная модель очень близка к вашей целевой задаче и датасету плюс обучалась на изображениях аналогичного размера, в связи с чем веса должны меняться очень незначительно. В таком случае обучение на уменьшенных изображениях может изменить эти предварительно обученные веса.

С другой стороны, если задача, для которой используется перенос обучения, будет использовать изображения других размеров, форм или стилей, то прогрессивное изменение размера, скорее всего, поможет повысить итоговую эффективность. В любом случае, чтобы узнать, насколько в этом есть смысл, нужно просто проверить.

Помимо перечисленных приемов, мы также можем применить аугментацию данных к контрольной выборке. До сих пор мы применяли ее только к обучающей, и контрольная всегда получала неизмененные изображения. Но может, нам стоит попробовать сделать прогнозы для нескольких аугментированных вариантов контрольной выборки и усреднить их? Этот подход мы и рассмотрим следующим.

Аугментация во время тестирования

Мы использовали случайную обрезку изображений как способ полезной аугментации данных, позволяющий повысить обобщаемость модели и снизить необходимый минимум обучающих данных. Когда мы используем такой способ обрезки, fastai автоматически применяет для изображений контрольной выборки

обрезку по центру, то есть выбирает максимальный квадрат в середине каждого изображения, не выходя за его края.

Тем не менее такой подход может оказаться проблематичным. Например, в датасетах с множественными метками на изображениях иногда присутствуют миниатюрные объекты, расположенные близко к краям. При обрезке по центру такие объекты могут быть полностью утрачены. Даже для таких задач, как наш пример с классификацией пород домашних животных, может случиться, что будет обрезан важный признак, необходимый для верного определения породы, например цвет носа.

Одно из решений этой проблемы — полностью избегать случайной обрезки. Вместо этого можно просто сжимать или растягивать прямоугольные изображения, заполняя ими квадрат. Но тогда мы лишимся очень полезной аугментации данных, усложнив при этом распознавание изображений для модели, потому что она будет вынуждена обучаться узнавать их не в естественных пропорциях, а в сжатой и растянутой форме.

Еще одно решение — не обрезать центральную часть изображений для контрольной выборки, а выбирать ряд областей, вырезаемых из оригинального прямоугольного изображения, пропускать каждое из них через модель и получать максимальное или среднее значение прогнозов. Фактически мы могли бы так сделать не только для разных вырезаемых участков изображений, но также и для разных значений всех параметров аугментации времени тестирования. Этот метод называется *аугментацией во время тестирования*.



ТЕРМИН: АУГМЕНТАЦИЯ ВО ВРЕМЯ ТЕСТИРОВАНИЯ (TTA)

Создание множества версий каждого изображения путем аугментации в процессе вывода или контроля модели, сопровождаемое получением среднего или максимального значения прогнозов для каждой аугментированной версии изображения.

В зависимости от датасета, аугментация во время тестирования может приводить к существенному повышению точности. Она не влияет на продолжительность обучения, но увеличивает время, необходимое для контроля модели или ее вывода, что связано с запрошенным количеством аугментированных изображений. По умолчанию fastai будет применять неаугментированное обрезание изображений по центру плюс четыре случайно аугментированных изображения.

Вы можете передать в fastai-метод `tta` любой `DataLoader`. По умолчанию же библиотека использует контрольную выборку:

```
preds,targs = learn.tta()  
accuracy(preds, targs).item()
```

0.8737863898277283

Мы видим, что применение TTA дает существенный прирост к показателю эффективности, не требуя при этом дополнительного обучения. Тем не менее процесс вывода замедляется: если выбрать для TTA получение среднего значения по пяти изображениям, то скорость вывода замедлится в пять раз.

Мы уже видели несколько примеров того, как аугментация данных помогает добиться лучших результатов от обучения модели. Теперь уделим внимание новой технике аугментации под названием *Mixup*.

Mixup

Метод аугментации Mixup был введен в работе Хунги Джан и др. (Hongyi Zhang) *mixup: Beyond Empirical Risk Minimization* (<https://oreil.ly/UvIkN>) («Mixup: за пределами минимизации рисков»), написанной в 2017 году. Это мощная техника, способная обеспечить очень высокий прирост точности, особенно в условиях дефицита данных и отсутствия модели, предварительно обученной на близких к вашему датасету данных. В работе дается следующее пояснение: «Несмотря на то что аугментация данных последовательно ведет к повышению обобщаемости модели, эта процедура зависит от датасета, в связи с чем требует от нас экспериментального уровня знаний». Например, общепринятый прием аугментации — это переворачивание изображения, но стоит ли переворачивать его только горизонтально? Ответ будет зависеть от вашего датасета. Кроме того, если переворачивание (например) не обеспечит достаточной степени аугментации, то вы уже не сможете «перевернуть больше». Поэтому полезно иметь в арсенале такие техники аугментации, которые «увеличивают» или «уменьшают» объем изменений, позволяя выбрать оптимальный вариант.

Для каждого изображения mixup работает так:

1. Случайно выбирает другое изображение из датасета.
2. Случайно выбирает вес.
3. Получает средневзвешенное (используя вес из шага 2) выбранного изображения с начальным изображением, определяя, таким образом, независимую переменную.
4. Получает средневзвешенное значение (взяв тот же вес) меток выбранного изображения с начальным, определяя зависимую переменную.

В псевдокоде мы делаем следующее (t является весом для средневзвешенного значения):

```
image2,target2 = dataset[randint(0,len(dataset))]
t = random_float(0.5,1.0)
new_image = t * image1 + (1-t) * image2
new_target = t * target1 + (1-t) * target2
```

Чтобы это сработало, цели должны быть закодированы унитарно. В работе это описывается с использованием уравнений, приведенных на рис. 7.1 (где λ представляет то же, что и t в нашем псевдокоде).

Contribution Motivated by these issues, we introduce a simple and data-agnostic data augmentation routine, termed *mixup* (Section 2). In a nutshell, *mixup* constructs virtual training examples

$$\begin{aligned}\hat{x} &= \lambda x_i + (1 - \lambda)x_j, && \text{where } x_i, x_j \text{ are raw input vectors} \\ \hat{y} &= \lambda y_i + (1 - \lambda)y_j, && \text{where } y_i, y_j \text{ are one-hot label encodings}\end{aligned}$$

Рис. 7.1. Отрывок из работы, раскрывающей принцип Mixup

ИССЛЕДОВАТЕЛЬСКИЕ РАБОТЫ И МАТЕМАТИКА

В книге мы будем все чаще обращаться к различным научным работам. Теперь, когда вы уже владеете базовой терминологией, вы удивитесь, как много из этих работ понятны при минимуме практики. Тем не менее одной из сложностей в процессе знакомства окажутся греческие буквы, например λ , которые встречаются во многих публикациях. Полезно выучить их названия, иначе чтение и запоминание соответствующей информации вызовет затруднения (сложности могут возникнуть и при чтении кода, так как в нем часто используются слова, обозначающие греческие буквы, например `lambda`).

Еще одна сложность научных трудов в том, что в них для объяснения происходящего используется математика. Если ваш опыт в этой области невелик, то первое время это будет несколько пугать и запутывать. При этом нужно помнить, что все показанное на математическом языке будет так или иначе реализовано в коде. Это просто другой способ говорить об одном и том же. Прочитав буквально несколько работ, вы будете понимать все больше и больше обозначений. Если вы не узнаете какой-то символ, попробуйте обратиться к таблице математических символов Википедии или нарисовать его на ресурсе Detexify (<https://oreil.ly/92u4d>), который (с помощью машинного обучения) найдет название написанного от руки символа. По этому названию вы уже сможете поискать в интернете и выяснить, что он вообще обозначает и для чего используется.

На рис. 7.2 показано, как выглядит линейное *совмещение изображений*, производимое аугментацией Mixup.

Третье изображение создано сложением первого и второго в соотношении 0,3 к 0,7 соответственно. Как вы думаете, в этом примере модель должна спрогнозировать «церковь» или «заправку»? Правильный ответ — на 30 % церковь и на 70 % заправку, поскольку именно это мы и получим, если рассмотрим линейное совмещение этих унитарно закодированных целей. Предположим, что у нас есть 10 классов и «церковь» представлена индексом 2, а «заправка» индексом 7. Унитарная кодировка в таком случае будет следующей:

[0, 0, 1, 0, 0, 0, 0, 0, 0] и [0, 0, 0, 0, 0, 0, 1, 0, 0]



Рис. 7.2. Совмещение церкви и заправки

А вот наша итоговая цель:

[0, 0, 0.3, 0, 0, 0, 0, 0.7, 0, 0]

Все это выполняется за нас внутри fastai путем добавления *обратного вызова* (Callback) в Learner. Такие обратные вызовы используются в fastai для внедрения пользовательского поведения в обучающий цикл (наподобие расписания скорости обучения или обучения со смешанной точностью). В главе 16 вы узнаете об обратных вызовах все, в том числе как создавать собственные. Сейчас же вам нужно знать лишь то, что для передачи обратных вызовов в Learner используется параметр `cbs`.

Вот вариант обучения модели с применением Mixup:

```
model = xresnet50()
learn = Learner(dls, model, loss_func=CrossEntropyLossFlat(),
                metrics=accuracy, cbs=Mixup)
learn.fit_one_cycle(5, 3e-3)
```

Что происходит, когда мы обучаем модель с помощью данных, «смешанных» таким образом? Очевидно, что учиться ей станет сложнее, потому что сложнее понять, что именно представлено на изображении. При этом модели придется прогнозировать не по одной, а по две метки для каждого из них, а также определять вес каждой. Тем не менее в таком случае риск переобучения существенно снижается, потому что мы показываем не одно и то же изображение в каждой эпохе, а произвольную комбинацию двух изображений.

Если сравнивать Mixup с другими подходами аугментации, которые мы видели, то при его использовании для достижения лучших показателей точности требуется выполнять обучение на протяжении гораздо большего числа эпох. Вы можете попробовать обучить Imagenette с Mixup и без, используя скрипт `examples/train_imagenette.py`, размещенный в репозитории fastai (<https://oreil.ly/lrGXE>). На момент написания книги таблица рекордов репозитория Imagenette показывает, что Mixup использовался для получения всех лидиру-

ющих результатов с обучением больше 80 эпох и что для меньшего числа эпох он не использовался вообще. Мы в своих экспериментах с использованием Mixup пришли к такому же ограничению.

Одна из прелестей Mixup в том, что его можно применять не только к фотографиям. На практике некоторые исследователи добивались хороших результатов, применяя его к активациям *внутри* их моделей, а не только к входным данным, — это уже позволяет использовать Mixup для NLP и других типов данных.

У Mixup есть и еще одна тонкая особенность, которая работает нам на руку. Дело в том, что при использовании всех моделей, которые мы видели до этого момента, фактически невозможно добиться идеальных показателей потерь. Все потому, что наши метки — это 1 и 0, а выходные значения softmax и сигмоиды никогда не могут быть равны 0. Это означает, что при обучении наша модель все сильнее сдвигает активации к этим значениям таким образом, что чем дольше мы обучаем, тем более экстремальными эти значения становятся.

При использовании Mixup эта проблема нас больше не тревожит, потому что наши метки будут конкретно 1 или 0, только если мы вдруг произведем «смешивание» с изображением того же класса. В других случаях метки будут линейной комбинацией, такой как 0,7 и 0,3, которую мы получили в примере с церковью и заправкой.

Как бы то ни было, но здесь все же присутствует одна сложность, которая заключается в том, что Mixup «случайно» делает метки больше 0 или меньше 1. Это значит, что мы *неявно* говорим модели, что хотим изменить метки таким образом. Поэтому если нам нужно, чтобы изменение приближало/удаляло метки к 0 и 1, то нужно изменить объем Mixup, что также изменит объем аугментации данных, а нам это может помешать. В таком случае есть способ решить эту проблему более непосредственно, и называется он *сглаживанием меток*.

Сглаживание меток

При теоретическом выражении потерь в задачах классификации наши цели кодируются унитарно (на практике мы склонны избегать этого для экономии памяти, но вычисляем мы те же потери, как если бы использовали унитарное кодирование). Это означает, что модель обучается возвращать 0 для всех категорий, кроме одной, для которой она обучается возвращать 1. Даже 0,999 недостаточно. Модель будет получать градиенты и учиться прогнозировать активации с еще большей уверенностью. Это вызывает переобучение и дает на выводе такую модель, которая не сможет выдавать никаких значимых вероятностей: она всегда будет называть единицу для прогнозируемой категории, даже если не будет особо уверена, просто потому, что так была обучена.

Это может существенно навредить, если данные размечены не идеально. В классификаторе медведей, изученном нами в главе 2, мы видели, что некоторые изображения были размечены с ошибками или представляли два разных вида медведей. Как правило, использовать вы будете не идеальные данные. Даже если метки создавались людьми вручную, при этом могли быть допущены ошибки или возникнуть разногласия в мнениях относительно изображений, которые сложно определить точно.

Вместо этого мы можем заменить все единицы на число чуть меньше 1, а все нули на число чуть больше 0 и потом уже начать обучение. Эта техника называется *сглаживанием меток*. Склоняя модель к меньшей уверенности, сглаживание меток делает процесс обучения более надежным, даже если в датасете присутствуют неверно размеченные данные. В результате на выводе получается более обобщающаяся модель.

На практике этот метод работает так: мы начинаем с закодированных унитарно меток, затем заменяем все 0 на $\frac{\epsilon}{N}$ (это греческая буква *эпсилон*, которая использовалась в работе, посвященной методу сглаживания меток (<https://oreil.ly/L3yrf>). Она также встречается в коде fastai), где N — это число классов, а ϵ — это параметр (обычно 0,1, что означает, что мы на 10 % не уверены в метках). Поскольку нам нужно, чтобы метки складывались в единицу, мы также заменяем все единицы на $1 - \epsilon + \frac{\epsilon}{N}$. Таким образом мы не склоняем модель прогнозировать что-то с излишней уверенностью. В нашем примере с Imagenette, где используется 10 классов, цели принимают следующий вид (приведен пример для цели, соответствующей индексу 3):

```
[0.01, 0.01, 0.01, 0.91, 0.01, 0.01, 0.01, 0.01, 0.01]
```

На практике мы не хотим кодировать метки унитарно, и к счастью, нам это не понадобится (унитарное кодирование просто удобно для объяснения сглаживания меток и наглядного представления этого процесса).

Чтобы использовать это на практике, нам нужно просто изменить функцию потерь в вызове к Learner:

```
model = xresnet50()
learn = Learner(dls, model, loss_func=LabelSmoothingCrossEntropy(),
                metrics=accuracy)
learn.fit_one_cycle(5, 3e-3)
```

Как и в случае с Mixup, вы, как правило, не увидите существенных улучшений от сглаживания меток, пока не проведете обучение на протяжении большего числа эпох. Проверьте сами: сколько эпох нужно обучать модель, чтобы сглаживание меток дало ощутимые улучшения?

СГЛАЖИВАНИЕ МЕТОК, НАУЧНАЯ РАБОТА

Вот какой ход рассуждений стоял за разработкой сглаживания меток в работе Кристиана Сегеди и др. (Christian Szegedy): «Этот максимум недостижим для конечного Z_k , но к нему стремится $Z_y >> Z_k$ при всех $k \neq y$ — это верно, если логит, соответствующий истинной метке, намного [больше] остальных логитов. Тем не менее это может вызвать две проблемы. Во-первых, привести к переобучению: если модель обучается присваивать полную вероятность истинной метке для каждого обучающего примера, то она не обязательно будет обобщаться. Во-вторых, это вызывает увеличение разности между самым большим логитом и всеми остальными, что совместно с привязанным градиентом $\frac{dl}{dz_k}$ снижает способность модели адаптироваться. Очевидно, что это происходит из-за того, что модель становится слишком уверенной в своих прогнозах».

Попрактикуем наш навык чтения научных работ и попробуем понять написанное. «Этот максимум» относится к предыдущей части абзаца, где говорилось о том, что 1 — это значение метки для утвердительного класса. Поэтому никакое значение (кроме бесконечности) не может привести к 1 после сигмоиды или softmax. В научных работах вы, как правило, не встретите выражений «любое значение». Вместо этого будет приведен символ, которым в данном случае является Z_k . Такое сокращение помогает в работах, потому что на него можно сослаться позднее, и читатель поймет, о чем идет речь.

Далее мы видим «если $Z_y >> Z_k$ для всех $k \neq y$ ». В этом случае в работе сразу идет математическое изложение с описанием на английском, что очень удобно, потому что можно просто его прочитать. В математике *у*значает цель (*у* определяется ранее в этой работе; иногда сложно найти, где именно определены *те* или *иные* символы, но почти во всех работах все употребляемые символы где-то да определяются), а Z_y — это соответствующая цели активация. Поэтому чтобы приблизиться к 1, этой активации нужно быть намного выше всех остальных в этом прогнозе.

Далее рассмотрим фразу «если модель обучается присваивать полную вероятность истинной метке для каждого примера обучения, она не обязательно будет обобщаться». Здесь говорится, что сильное увеличение Z_y означает, что нам во всей модели понадобятся большие веса и высокие активации. Большие веса ведут к «неровным» функциям, в которых небольшое изменение ввода очень сильно влияет на прогнозы. Это очень плохо для обобщаемости, так как означает, что всего одно небольшое изменение пикселя может полностью изменить прогноз.

Заканчивается же цитата так: «Это вызывает увеличение разности между самым большим логитом и всеми остальными, что совместно с ограниченным градиентом $\frac{dl}{dz_k}$ снижает способность модели адаптироваться. Вспомните, что градиент перекрестной энтропии — это, по сути, *output* — *target*. И *output*, и *target* находятся между 0 и 1, поэтому разность находится между -1 и 1, почему в работе и сказано, что градиент «ограничен» (он не может быть бесконечным), в результате чего шаги SGD тоже ограничены. «Снижает способность модели адаптироваться» означает, что ей сложно обновляться в условиях переноса обучения. Это следует из того, что разница в потерях, вызванная неверными прогнозами, не ограничена, а мы можем каждый раз совершать только ограниченный шаг.

Резюме

Теперь вы увидели все, что требуется для обучения эталонной модели компьютерного зрения, будь то обучение с нуля или при его переносе. Теперь от вас требуется лишь экспериментировать на своих задачах. Проверьте, избегает ли более длительное обучение с применением Mixup и/или сглаживания меток феномена переобучения и дает ли лучшие результаты? Попробуйте также постепенное увеличение размера и аугментацию во время тестирования.

При этом самое важное — запомнить, что если у вас большой датасет, то нет смысла делать на нем полное прототипирование. Выделите небольшое подмножество, которое будет представлять датасет целиком, как мы сделали с Imagenette, и экспериментируйте с ним.

В следующих трех главах мы рассмотрим другие способы применения глубокого обучения, непосредственно поддерживаемые fastai, а именно совместную фильтрацию, табличное моделирование и работу с текстом. Позднее мы еще вернемся к компьютерному зрению, а в главе 13 подробно разберем сверточные нейронные сети.

Вопросник

1. В чем отличие между ImageNet и Imagenette? С какой базой данных и в каких ситуациях следует экспериментировать?
2. Что такое нормализация?
3. Почему нас не волновала нормализация, когда мы использовали предварительно обученную модель?
4. Что такое постепенное увеличение размера?
5. Реализуйте постепенное увеличение размера в собственном проекте. Помогло ли оно?
6. Что такое аугментация во время тестирования? Как она используется в fastai?
7. Использование ТТА в процессе вывода модели замедляет или ускоряет этот процесс?
8. Что такое Mixup? Как этот метод используется в fastai?
9. Почему Mixup не позволяет модели быть слишком уверенной?
10. Почему обучение с Mixup в течение пяти эпох приводит к худшему показателю эффективности, чем обучение без него?
11. Какая идея стоит за сглаживанием меток?

12. Какие проблемы с данными может решить сглаживание меток?
13. Какая цель будет ассоциироваться с индексом 1 при использовании сглаживания меток по пяти категориям?
14. С чего следует начинать подготовку к проведению быстрых экспериментов по прототипированию с новым датасетом?

Дополнительные задания

1. Используйте документацию fastai для построения функции, обрезающей изображение до квадрата в каждом из четырех углов. Затем реализуйте метод TTA, усредняющий прогнозы обрезки по центру и тех четырех частей. Улучшился ли результат? Лучше ли этот метод, чем TTA в fastai?
2. Найдите и прочитайте работу по исследованию Mixup на arXiv. Ознакомьтесь с одной или двумя недавними статьями, описывающими вариации Mixup, и попробуйте реализовать эти вариации в своем проекте.
3. Найдите скрипт обучения Imagenette с использованием Mixup и примените его в качестве примера для построения скрипта длительного обучения в своем проекте. Выполните его и посмотрите, улучшил ли это результат.
4. Прочтите сноску «Сглаживание меток, научная работа», затем откройте связанный с ней раздел оригинала работы и попробуйте его понять. Не стесняйтесь обращаться за помощью!

ГЛАВА 8

Коллаборативная фильтрация

Одна из наиболее востребованных задач в области машинного обучения — это генерация полезных для пользователей рекомендаций относительно имеющихся в предложении товаров. В этом плане возможны разные ситуации: например, какое кино рекомендовать (как на Netflix), что выделять для пользователя на главной странице, какие истории показывать в ленте соцсетей и т. д. Для этой задачи есть общее решение — *коллаборативная фильтрация, или, как еще говорят, совместная фильтрация*. Принцип ее действия такой: просмотреть товары, которые пользователь отметил как понравившиеся или использовал/купил, найти других пользователей, которые тоже использовали или отметили аналогичные товары, и порекомендовать ему другие товары, отмеченные или использованные этими пользователями.

На Netflix вы могли просмотреть много фильмов из области научной фантастики в жанре экшен и снятых в 1970-е. Netflix могут не быть известны эти конкретные свойства просмотренных вами фильмов, но платформа сможет увидеть, что некоторые пользователи, смотревшие эти же фильмы, также смотрели и другие научно-фантастические ленты из 1970-х с выраженным уровнем экшена. Другими словами, для использования этого подхода не обязательно знать подробности фильмов, достаточно лишь знать, кто любит их смотреть.

Есть и более обобщенный класс задач, которые этот подход может решать и которые не обязательно включают пользователей и товары. Фактически при реализации техники совместной фильтрации мы обычно обращаемся к *элементам*, а не к *товарам*. Элементами могут быть ссылки, по которым переходят люди, диагнозы, выбираемые для пациентов, и т. д.

Основная идея заключается в *скрытых факторах*. В примере с Netflix мы начали с предположения, что вам интересны старые научно-фантастические фильмы с элементами экшена. Но вы никогда не говорили Netflix, что такие фильмы вам нравятся. При этом данной платформе не требуется добавлять в таблицы фильма столбцы с указанием типов этих фильмов. Тем не менее в основе должна лежать некая общая концепция научной фантастики, уровня экшена

и возрастной категории фильма, которой руководствуются все пользователи (или хотя бы некоторые из них) при выборе этих фильмов.

В этой главе мы как раз проработаем задачу рекомендации фильмов. Начнем мы с получения подходящих для модели совместной фильтрации данных.

Первый взгляд на данные

У нас нет доступа ко всему датасету истории просмотров Netflix, но мы можем использовать другой отличный датасет — MovieLens (<https://oreil.ly/gP3Q5>). Он содержит десятки миллионов оценок фильмов (состоящих из ID фильма, ID пользователя и самой оценки), хотя мы для нашего примера задействуем только подмножество из 100 000. Если вам интересно, то в качестве тренировки будет здорово попробовать повторить этот подход для всего рекомендательного датасета из 25 миллионов единиц данных, который доступен на сайте MovieLens.

Текущий датасет можно получить с помощью обычной функции fastai:

```
from fastai.collab import *
from fastai.tabular.all import *
path = untar_data(URLs.ML_100k)
```

Согласно README, основная таблица находится в файле *u.data*. Она разбита с помощью табуляции и содержит столбцы user, movie, rating и timestamp. Поскольку эти имена не закодированы, то при чтении файла с помощью Pandas нам нужно их указать. Вот как можно открыть и просмотреть данную таблицу:

```
ratings = pd.read_csv(path/'u.data', delimiter='\t', header=None,
                      names=['user', 'movie', 'rating', 'timestamp'])
ratings.head()
```

	user	movie	rating	timestamp
0	196	242	3	881250949
1	186	302	3	891717742
2	22	377	1	878887116
3	244	51	2	880606923
4	166	346	1	886397596

Несмотря на то что здесь есть вся необходимая нам информация, этот способ просмотра данных не особо удобен для восприятия. На рис. 8.1 показаны те же данные, но уже в более удобной перекрестной таблице.

	movield															
	27	49	57	72	79	89	92	99	143	179	180	197	402	417	505	
userId	14	3.0	5.0	1.0	3.0	4.0	4.0	5.0	2.0	5.0	5.0	4.0	5.0	5.0	2.0	5.0
29	5.0	5.0	5.0	4.0	5.0	4.0	4.0	5.0	4.0	4.0	5.0	5.0	5.0	3.0	4.0	5.0
72	4.0	5.0	5.0	4.0	5.0	3.0	4.5	5.0	4.5	5.0	5.0	5.0	4.5	5.0	4.0	
211	5.0	4.0	4.0	3.0	5.0	3.0	4.0	4.5	4.0		3.0	3.0	5.0	3.0		
212	2.5		2.0	5.0		4.0	2.5		5.0	5.0	3.0	3.0	4.0	3.0	2.0	
293	3.0		4.0	4.0	4.0	3.0		3.0	4.0	4.0	4.5	4.0	4.5	4.0		
310	3.0	3.0	5.0	4.5	5.0	4.5	2.0	4.5	4.0	3.0	4.5	4.5	4.0	3.0	4.0	
379	5.0	5.0	5.0	4.0		4.0	5.0	4.0	4.0	4.0		3.0	5.0	4.0	4.0	
451	4.0	5.0	4.0	5.0	4.0	4.0	5.0	5.0	4.0	4.0	4.0	4.0	2.0	3.5	5.0	
467	3.0	3.5	3.0	2.5			3.0	3.5	3.5	3.0	3.5	3.0	3.0	4.0	4.0	
508	5.0	5.0	4.0	3.0	5.0	2.0	4.0	4.0	5.0	5.0	5.0	3.0	4.5	3.0	4.5	
546		5.0	2.0	3.0	5.0		5.0	5.0		2.5	2.0	3.5	3.5	3.5	5.0	
563	1.0	5.0	3.0	5.0	4.0	5.0	5.0		2.0	5.0	5.0	3.0	3.0	4.0	5.0	
579	4.5	4.5	3.5	3.0	4.0	4.5	4.0	4.0	4.0	4.0	3.5	3.0	4.5	4.0	4.5	
623		5.0	3.0	3.0		3.0	5.0		5.0	5.0	5.0	5.0	2.0	5.0	4.0	

Рис. 8.1. Перекрестная таблица фильмов и пользователей

Для этой таблицы мы выбрали только некоторые из наиболее популярных фильмов совместно с наиболее активными пользователями. Пустые ячейки представляют сегменты, которые наша модель должна научиться заполнять. Фактически эти сегменты означают, что пользователь этот фильм еще не оценивал, вероятно, потому что еще не смотрел. Наша задача состоит в том, чтобы выяснить наиболее предпочтительный для каждого пользователя фильм из представленных.

Если мы будем знать о каждом пользователе, в какой степени ему нравится та или иная категория, например жанр, год выхода фильма, режиссеры, актеры и т. д., и будем иметь ту же информацию по каждому фильму, то проще всего заполнить таблицу, перемножив эти данные по каждому фильму и использовав их итоговую комбинацию. Например, если предположить, что эти факторы находятся в диапазоне от -1 до $+1$, где положительные числа означают большее соответствие, а категориями являются «научная фантастика», «экшен» и «старое кино», то можно представить фильм *The Rise of Skywalker* («Скайуокер. Восход») так:

```
rise_skywalker = np.array([0.98, 0.9, -0.9])
```

Здесь мы, к примеру, оцениваем фильм как *очень научно-фантастический* значением 0,98 и *очень нестарый* значением -0.9 . Пользователя, который любит современные научно-фантастические экшены, мы можем представить так:

```
user1 = np.array([0.9, 0.8, -0.6])
```

Теперь можно вычислить соответствие между этой комбинацией:

```
(user1*rise_skywalker).sum()
```

```
2.1420000000000003
```

Когда мы умножаем элементы двух векторов и складываем результаты, то получаем скалярное произведение. Оно широко используется в машинном обучении и лежит в основе матричного умножения. Мы будем подробно рассматривать матричное умножение и скалярное произведение в главе 17.



ТЕРМИН: СКАЛЯРНОЕ ПРОИЗВЕДЕНИЕ

Математическая операция умножения элементов двух векторов с последующим суммированием результатов.

С другой стороны, мы можем представить фильм *Casablanca* («Касабланка») так:

```
casablanca = np.array([-0.99, -0.3, 0.8])
```

Соответствие между этой комбинацией получается следующим:

```
(user1*casarblanca).sum()
```

```
-1.611
```

Поскольку нам неизвестны латентные (скрытые) факторы и мы не знаем, как оценивать их для каждого пользователя и фильма, то их нужно изучить.

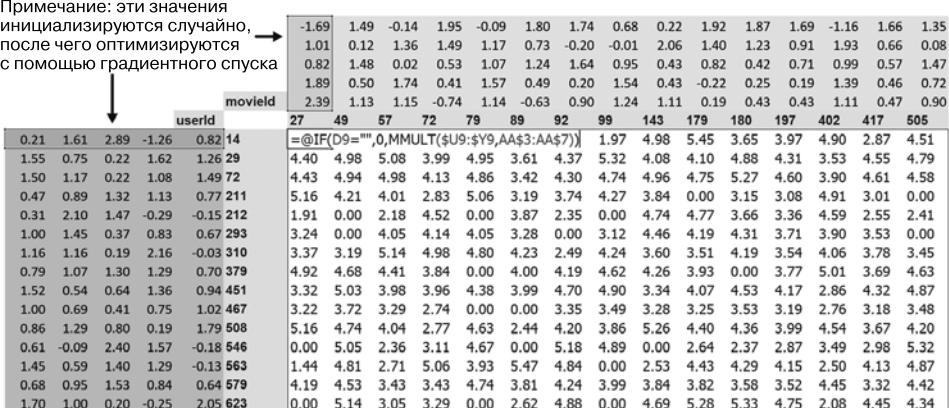
Обучение скрытых факторов

Между определением структуры модели, которое мы делали в предыдущем разделе, и ее обучением на удивление мало отличий, поскольку можно просто использовать общий подход градиентного спуска.

Шаг 1 этого подхода подразумевает случайную инициализацию ряда параметров. Эти параметры будут набором скрытых факторов для каждого пользователя и фильма. При этом нам потребуется решить, сколько именно их нужно задействовать. Мы еще обсудим, как определять это количество, но для наглядности давайте пока что возьмем 5. Так как набор этих факторов будет у каждого пользователя и у каждого фильма, мы можем показывать их инициализированные случайные значения в таблице рядом с пользователями и фильмами и впоследствии вставлять между ними скалярные произведения для каждой из комбинаций. На рис. 8.2 показано, как это будет выглядеть в Microsoft Excel, где в верхней левой ячейке таблицы показан пример формулы.

Шагом 2 будет вычисление прогнозов. Как мы уже говорили, для этого нужно просто взять скалярное произведение каждого фильма и пользователя. Если, к примеру, первый скрытый фактор пользователя представляет, насколько пользователь любит экшены, а первый скрытый фактор фильма выражает степень, в которой он является экшеном, то скалярное произведение будет большим, если пользователь любит такие фильмы и этот фильм в существенной степени к ним относится либо если пользователь не любит экшены и фильм в себе элементов экшена не содержит. С другой стороны, если есть несоответствие (пользователь любит экшены, но в фильме экшена нет), то произведение получится очень низким.

Примечание: эти значения инициализируются случайно, →
после чего оптимизируются
с помощью градиентного спуска



	movielid															
userId	27	49	57	72	79	89	92	99	143	179	180	197	402	417	505	
0.21	1.61	2.89	-1.26	0.82	14											
1.55	0.75	0.22	1.62	1.26	29											
1.50	1.17	0.22	1.08	1.49	72											
0.47	0.89	1.32	1.13	0.77	211											
0.31	2.10	1.47	-0.29	-0.15	212											
1.00	1.45	0.37	0.83	0.67	293											
1.16	1.16	0.19	2.16	-0.03	310											
0.79	1.07	1.30	1.29	0.70	379											
1.52	0.54	0.64	1.36	0.94	451											
1.00	0.69	0.41	0.75	1.02	467											
0.86	1.29	0.80	0.19	1.79	508											
0.61	-0.09	2.40	1.57	-0.18	546											
1.45	0.59	1.40	1.29	-0.13	563											
0.68	0.95	1.53	0.84	0.64	579											
1.70	1.00	0.20	-0.25	2.05	623											
						=@IF(D9="",0,MMULT(\$U9:\$Y9,AA\$3:AA\$7))			1.97	4.98	5.45	3.65	3.97	4.90	2.87	4.51
						4.40	4.98	5.08	3.99	4.95	3.61	4.37	5.32	4.08	4.10	4.88
						4.43	4.94	4.98	4.13	4.86	3.42	4.30	4.74	4.96	4.75	5.27
						5.16	4.21	4.01	2.83	5.06	3.19	3.74	4.27	3.84	0.00	3.15
						1.91	0.00	2.18	4.52	0.00	3.87	2.35	0.00	4.74	4.77	3.66
						3.24	0.00	4.05	4.14	4.05	3.28	0.00	3.12	4.46	4.19	4.31
						3.37	3.19	5.14	4.98	4.80	4.23	2.49	4.24	3.60	3.51	4.19
						4.92	4.68	4.41	3.84	0.00	4.00	4.19	4.62	4.26	3.93	0.00
						3.32	5.03	3.98	3.96	4.38	3.99	4.70	4.90	3.34	4.07	4.53
						3.22	3.72	3.29	2.74	0.00	0.00	3.35	3.49	3.28	3.25	3.53
						5.16	4.74	4.04	2.77	4.63	2.44	4.20	3.86	5.26	4.40	4.36
						0.00	5.05	2.36	3.11	4.67	0.00	5.18	4.89	0.00	2.64	2.37
						1.44	4.81	2.71	5.06	3.93	5.47	4.84	0.00	2.53	4.43	4.29
						4.19	4.53	3.43	3.43	4.74	3.81	4.24	3.99	3.84	3.82	3.58
						0.00	5.14	3.05	3.29	0.00	2.62	4.88	0.00	4.69	5.28	5.33

Рис. 8.2. Скрытые факторы в перекрестной таблице

Шагом 3 будет вычисление потерь, для чего подойдет любая предпочтительная функция. Давайте возьмем среднеквадратичную ошибку, поскольку это один из наиболее подходящих способов представить точность прогноза.

Вот и все, что нам нужно. Теперь мы можем оптимизировать параметры (латентные факторы) с помощью стохастического градиентного спуска, как и в случае минимизации потерь. На каждом этапе оптимизатор будет вычислять соответствие между каждым фильмом и пользователем, используя скалярное произведение. Он будет сравнивать полученный результат с фактической оценкой, которую пользователь каждому фильму присвоил. После этого он будет вычислять производную этого значения и смещать веса, умножая ее на скорость обучения. После многократного повторения этого процесса потери станут улучшаться, а вместе с этим повысится и качество прогнозов.

Чтобы использовать обычную функцию `Learner.fit`, нам понадобится поместить данные в `DataLoaders`, чем мы сейчас и займемся.

Создание DataLoaders

При показе данных будет удобнее отображать названия фильмов, а не их ID. Таблица `u.item` содержит соотнесение ID с названиями:

```
movies = pd.read_csv(path/'u.item', delimiter='|', encoding='latin-1',
                      usecols=(0,1), names=('movie','title'), header=None)
movies.head()
```

	movie	title
0	1	Toy Story (1995)
1	2	GoldenEye (1995)
2	3	Four Rooms (1995)
3	4	Get Shorty (1995)
4	5	Copycat (1995)

Мы можем совместить ее с таблицей `ratings`, сгруппировав пользовательские оценки по названию фильма:

```
ratings = ratings.merge(movies)
ratings.head()
```

	user	movie	rating	timestamp	title
0	196	242	3	881250949	Kolya (1996)
1	63	242	3	875747190	Kolya (1996)
2	226	242	5	883888671	Kolya (1996)
3	154	242	3	879138235	Kolya (1996)
4	306	242	5	876503793	Kolya (1996)

Затем можно создать из этой таблицы объект `DataLoaders`. По умолчанию первый столбец отводится пользователям, второй — элементам (здесь это фильмы), а третий — оценкам. В нашем случае нужно изменить значение `item_name`, чтобы использовать вместо ID названия:

```
dls = CollabDataLoaders.from_df(ratings, item_name='title', bs=64)
dls.show_batch()
```

	user	title	rating
0	207	Four Weddings and a Funeral (1994)	3
1	565	Remains of the Day, The (1993)	5
2	506	Kids (1995)	1
3	845	Chasing Amy (1997)	3
4	798	Being Human (1993)	2
5	500	Down by Law (1986)	4
6	409	Much Ado About Nothing (1993)	3
7	721	Braveheart (1995)	5
8	316	Psycho (1960)	2
9	883	Judgment Night (1993)	5

Для представления совместной фильтрации в PyTorch можно просто задействовать сам формат перекрестной таблицы, особенно если нам нужно, чтобы она подходила под наш фреймворк глубокого обучения. Таблицы латентных факторов пользователей и фильмов представляются в виде простых матриц:

```
n_users = len(dls.classes['user'])
n_movies = len(dls.classes['title'])
n_factors = 5

user_factors = torch.randn(n_users, n_factors)
movie_factors = torch.randn(n_movies, n_factors)
```

Чтобы вычислить результат для конкретной комбинации фильма и пользователя, нужно найти индекс этого фильма в матрице скрытых факторов фильмов и индекс пользователя в матрице скрытых факторов пользователей. После этого можно получить скалярное произведение между двумя векторами этих латентных факторов. Но модели глубокого обучения не знают, как выполнять операцию *поиска в индексе*, зато знают, как выполнять матричное произведение и функции активации.

К счастью, оказывается, что мы можем представить *поиск в индексе* в виде матричного произведения. Хитрость в том, чтобы заменить индексы на унитарно закодированные векторы. Вот пример того, что произойдет, если мы умножим вектор на унитарно закодированный вектор, представляющий индекс 3:

```
one_hot_3 = one_hot(3, n_users).float()
user_factors.t() @ one_hot_3

tensor([-0.4586, -0.9915, -0.4052, -0.3621, -0.5908])
```

В итоге мы получаем тот же вектор, что и для индекса 3, в матрице:

```
user_factors[3]
tensor([-0.4586, -0.9915, -0.4052, -0.3621, -0.5908])
```

Если мы сделаем это сразу для нескольких индексов, то получим матрицу унитарно закодированных векторов, а сама операция будет называться *матричным умножением*. Построение моделей на основе такой архитектуры оказалось бы идеальным решением, но потребовало бы намного больше ресурсов памяти и времени. Нам известно, что нет никакой фундаментальной причины хранить вектор в унитарной кодировке или выполнять по нему поиск числа 1, — у нас должна быть возможность напрямую обращаться к индексу массива, используя целое число. В связи с этим большинство библиотек глубокого обучения, включая PyTorch, содержат особый слой, который отвечает только за это. Данный слой обращается к индексу вектора, используя целое число, но производную при этом вычисляет таким образом, что она получается идентичной той, какая получилась бы при выполнении матричного умножения с использованием вектора в унитарной кодировке. Это называется *вложением*.



ТЕРМИН: ВЛОЖЕНИЕ

Умножение на быстро закодированную матрицу с использованием вычислительного сокращения, которое реализуется обращением к индексу напрямую. Это несколько вычурное слово для очень простого принципа. То, на что вы умножаете матрицу в быстрой кодировке (или с помощью вычислительного сокращения обращаетесь к индексу напрямую), называется *матрицей вложений*.

В компьютерном зрении есть очень простой способ получения всей информации о пикселе через его значения RGB: каждый пиксель цветного изображения представлен тремя числами. Эти три числа определяют степень красного, зеленого и синего, чего нашей модели достаточно для дальнейшей работы.

Для текущей задачи у нас нет настолько же легкого способа охарактеризовать пользователя или фильм. Вероятно, есть зависимость от жанров: если пользователь любит мелодрамы, то наверняка будет высоко оценивать именно такие фильмы. К другим факторам можно причислить то, относится фильм к экшну или, наоборот, содержит много диалогов, снимался ли в нем конкретный интересный пользователю актер и т. д.

Как же мы определяем числа для описания всего этого? Ответ прост: мы не определяем. Мы позволяем модели *выучить* их. Анализируя существующие связи между пользователями и фильмами, наша модель сама может выяснить признаки, которые являются важными, и наоборот.

Так и работают вложения. Мы будем приписывать каждому пользователю и фильму случайный вектор конкретной длины (здесь им будет `n_factor=5`)

и создавать тем самым обучаемые параметры. Это означает, что на каждом шаге при вычислении потерь путем сравнивания прогнозов с целями мы будем вычислять градиенты потерь в отношении этих векторов вложений и обновлять их согласно правилам SGD (или другого оптимизатора).

Вначале эти числа ничего не будут значить, поскольку выбираются случайно, но к концу обучения их значение станет решающим. Изучая связи между пользователями и фильмами на существующих данных и не имея никакой дополнительной информации, они все равно будут находить важные признаки и отделять блокбастеры от детективов, экшены от мелодрам и т. д.

Вот теперь мы готовы к созданию всей модели с нуля.

Коллаборативная фильтрация с нуля

Прежде чем писать модель в PyTorch, нужно изучить основы ООП и Python. Если вы еще не занимались объектно-ориентированным программированием, то мы кратко о нем расскажем, но при этом советуем самостоятельно поискать обучающий материал и попрактиковать.

Ключевой идеей ООП является *класс*. Мы использовали классы на протяжении всей книги: ими были `DataLoader`, `String` и `Learner`. При этом Python также облегчает для нас создание новых классов. Вот простой пример:

```
class Example:  
    def __init__(self, a): self.a = a  
    def say(self,x): return f'Hello {self.a}, {x}.'
```

Самой важной частью здесь является особый метод под названием `__init__` (произносится как *дандер init*). В Python любой метод, окруженный нижними подчеркиваниями, рассматривается как особый. Он указывает, что с ним связано некое дополнительное поведение. Если рассматривать метод `__init__`, то он вызывается при создании объекта. Это значит, что в нем можно установить нужное состояние, которое должно инициализироваться при создании объекта. Любые параметры, включаемые пользователем при построении экземпляра вашего класса, будут переданы в качестве параметров методу `__init__`. Обратите внимание, что первый параметр любого определяемого в классе метода — это `self`, который можно использовать для установки и получения нужных вам атрибутов:

```
ex = Example('Sylvain')  
ex.say('nice to meet you')  
  
'Hello Sylvain, nice to meet you.'
```

Заметьте также, что создание нового модуля PyTorch требует наследования от `Module`. *Наследование* — это важный принцип ООП, который подробно мы по-

яснять не будем. Говоря коротко, он означает, что мы можем добавлять в существующий класс дополнительное поведение. В PyTorch уже присутствует класс **Module**, представляющий базовые основы, с которых мы и начнем построение. Итак, мы добавим имя этого *суперкласса* после имени класса, который определяем, как показано в следующих примерах.

Последнее, что вам нужно знать для создания нового модуля PyTorch, — это то, что при вызове вашего модуля PyTorch будет вызывать метод класса, называемый **forward**, и передавать ему любые параметры, включенные в начальный вызов. Ниже вы видите класс, определяющий модель скалярного произведения:

```
class DotProduct(Module):
    def __init__(self, n_users, n_movies, n_factors):
        self.user_factors = Embedding(n_users, n_factors)
        self.movie_factors = Embedding(n_movies, n_factors)

    def forward(self, x):
        users = self.user_factors(x[:,0])
        movies = self.movie_factors(x[:,1])
        return (users * movies).sum(dim=1)
```

Если вам не доводилось встречаться с ООП, не беспокойтесь — в этой книге его особо использовать не придется. Мы просто упоминаем здесь этот подход, потому что большинство онлайн-уроков и документация содержат объектно-ориентированный синтаксис.

Обратите внимание, что ввод модели — это тензор формы **batch_size** × 2, где первый столбец (**x[:, 0]**) содержит ID пользователей, а второй (**x[:, 1]**) — ID фильмов. Как уже говорилось, мы используем слои *вложений* для представления наших матриц латентных факторов пользователей и фильмов:

```
x,y = dls.one_batch()
x.shape

torch.Size([64, 2])
```

Теперь, когда мы определили архитектуру и создали матрицы параметров, нужно создать **Learner** для оптимизации модели. В предыдущих примерах мы использовали специальные функции под названием **cnn_learner**, которые всё настраивали за нас для конкретного приложения. Но раз уж здесь мы все делаем с нуля, то будем использовать простой класс **Learner**:

```
model = DotProduct(n_users, n_movies, 50)
learn = Learner(dls, model, loss_func=MSELossFlat())
```

Теперь займемся настройкой модели:

```
learn.fit_one_cycle(5, 5e-3)
```

epoch	train_loss	valid_loss	time
0	1.326261	1.295701	00:12
1	1.091352	1.091475	00:11
2	0.961574	0.977690	00:11
3	0.829995	0.893122	00:11
4	0.781661	0.876511	00:12

Первое, что можно сделать для небольшого улучшения модели, — это ограничить возможные прогнозы до значений между 0 и 5. Для этого нужно просто использовать `sigmoid_range`, как в главе 6. На основе проб и ошибок мы выяснили, что лучше, если диапазон будет несколько выходить за 5, поэтому здесь используем $(0, 5.5)$:

```
class DotProduct(Module):
    def __init__(self, n_users, n_movies, n_factors, y_range=(0,5.5)):
        self.user_factors = Embedding(n_users, n_factors)
        self.movie_factors = Embedding(n_movies, n_factors)
        self.y_range = y_range

    def forward(self, x):
        users = self.user_factors(x[:,0])
        movies = self.movie_factors(x[:,1])
        return sigmoid_range((users * movies).sum(dim=1), *self.y_range)

model = DotProduct(n_users, n_movies, 50)
learn = Learner(dls, model, loss_func=MSELossFlat())
learn.fit_one_cycle(5, 5e-3)
```

epoch	train_loss	valid_loss	time
0	0.976380	1.001455	00:12
1	0.875964	0.919960	00:12
2	0.685377	0.870664	00:12
3	0.483701	0.874071	00:12
4	0.385249	0.878055	00:12

Неплохое начало, но можно добиться лучшего. Один очевидно упущеный фактор в том, что некоторые пользователи дают намного более позитивные или негативные оценки, и некоторые фильмы на порядок либо лучше, либо хуже других. Но в нашем представлении скалярного произведения у нас нет способа закодировать ни то ни другое. По правде говоря, если вы можете

сказать о фильме, например, что он научно-фантастический, содержит много экшена и нестарый, то у вас нет способа выразить, нравится ли он большинству людей.

Причина в том, что в этот момент мы имеем только веса, но не имеем смещений. Если у нас есть число для каждого пользователя, которое можно добавить к оценкам, и то же самое для каждого фильма, то это вполне компенсирует упущеный фактор. Поэтому для начала подкорректируем архитектуру модели:

```
class DotProductBias(Module):
    def __init__(self, n_users, n_movies, n_factors, y_range=(0,5.5)):
        self.user_factors = Embedding(n_users, n_factors)
        self.user_bias = Embedding(n_users, 1)
        self.movie_factors = Embedding(n_movies, n_factors)
        self.movie_bias = Embedding(n_movies, 1)
        self.y_range = y_range

    def forward(self, x):
        users = self.user_factors(x[:,0])
        movies = self.movie_factors(x[:,1])
        res = (users * movies).sum(dim=1, keepdim=True)
        res += self.user_bias(x[:,0]) + self.movie_bias(x[:,1])
        return sigmoid_range(res, *self.y_range)
```

Попробуем ее обучить и оценим результат:

```
model = DotProductBias(n_users, n_movies, 50)
learn = Learner(dls, model, loss_func=MSELossFlat())
learn.fit_one_cycle(5, 5e-3)
```

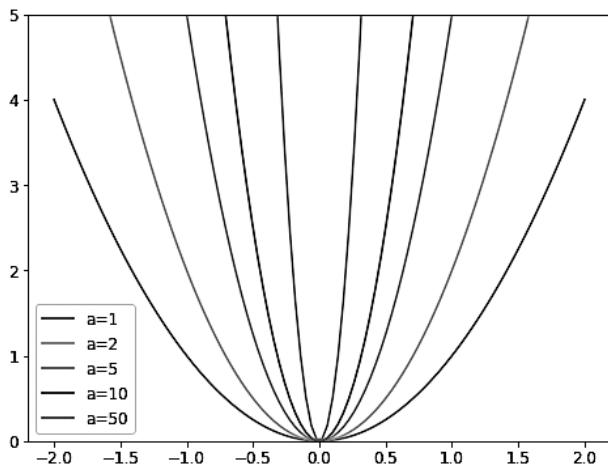
epoch	train_loss	valid_loss	time
0	0.929161	0.936303	00:13
1	0.820444	0.861306	00:13
2	0.621612	0.865306	00:14
3	0.404648	0.886448	00:13
4	0.292948	0.892580	00:13

Вместо улучшения показатели потерь ухудшаются (по крайней мере в конце обучения). Почему так происходит? Если внимательно посмотреть на оба этапа обучения, то мы увидим, что показатель для контрольной выборки перестал улучшаться в середине, откуда его значение начало, наоборот, ухудшаться. Как мы уже знаем, это очевидный признак переобучения. В данном случае нет возможности использовать аугментацию данных, поэтому мы вынуждены использовать другую технику регуляризации, в качестве которой нам подойдет подход *сокращения весов*.

Сокращение весов

Сокращение весов, иначе называемое *регуляризацией L2*, представляет собой добавление к функции потерь, возведенной в квадрат, суммы всех весов. Зачем нам это делать? Потому что при вычислении градиентов это будет добавлять им вклад, способствующий максимальному сокращению весов.

Почему это должно предотвратить переобучение? Суть в том, что чем больше коэффициенты, тем круче спуски, получаемые в функции потерь. Если взять простой пример параболы $y = a * (x^2)$, то чем больше a , тем более *узкой* получается парабола:



Поэтому если позволить модели обучать высокие параметры, то она подстроит все точки данных обучающей выборки с помощью сверхсложной функции, имеющей резкие изменения, что приведет к переобучению.

Препятствуя чрезмерному увеличению весов, мы замедлим обучение модели, но она окажется в состоянии, способствующем лучшей обобщаемости. Если кратко вернуться к теории, то сокращение весов (или просто wd) — это параметр, контролирующий сумму квадратов, добавляемых нами к потерям (предполагая, что `parameters` является тензором всех параметров):

```
loss_with_wd = loss + wd * (parameters**2).sum()
```

Тем не менее на практике будет очень неэффективно (а может, и численно нестабильно) вычислять настолько большую сумму и добавлять ее к потерям. Вы можете припомнить из курса высшей математики, что производная r^{**2} по отношению к r равна $2*r$. Поэтому добавление такой большой суммы к потерям эквивалентно следующему:

```
parameters.grad += wd * 2 * parameters
```

В реальности, так как `wd` является выбираемым нами параметром, мы можем его удвоить, исключив из данного уравнения *2. Для использования сокращения весов в `fastai` нужно передать `wd` в вызов `fit` или `fit_one_cycle` (можно передать в оба):

```
model = DotProductBias(n_users, n_movies, 50)
learn = Learner(dls, model, loss_func=MSELossFlat())
learn.fit_one_cycle(5, 5e-3, wd=0.1)
```

epoch	train_loss	valid_loss	time
0	0.972090	0.962366	00:13
1	0.875591	0.885106	00:13
2	0.723798	0.839880	00:13
3	0.586002	0.823225	00:13
4	0.490980	0.823060	00:13

Намного лучше!

Создание собственного модуля вложений

До сих пор мы использовали `Embedding`, не обращая внимания на то, как он фактически работает. Давайте воссоздадим `DotProductBias` без применения этого класса. Нам понадобится матрица случайным образом инициализированных весов для каждого из вложений. При этом нужно быть осторожными. В главе 4 мы писали, что оптимизаторы требуют возможности получения всех параметров модуля из его метода `parameters`. Тем не менее это происходит не полностью автоматически. Если мы просто добавим тензор к `Module` в качестве атрибута, то в `parameters` он включен не будет:

```
class T(Module):
    def __init__(self): self.a = torch.ones(3)

L(T().parameters())
(#0) []
```

Чтобы сообщить `Module` о своем желании рассматривать тензор как параметр, нужно обернуть его в класс `nn.Parameter`. Этот класс не привносит никакую функциональность (за исключением автоматического вызова `requires_grad`). Он просто используется как «маркер», показывающий, что нужно включить в `parameters`:

```
class T(Module):
    def __init__(self): self.a = nn.Parameter(torch.ones(3))
```

```
L(T().parameters())
(#1) [Parameter containing:
tensor([1., 1., 1.], requires_grad=True)]
```

Все модули PyTorch задействуют `nn.Parameter` для всех обучаемых параметров, почему нам и не требовалось специально использовать эту обертку до текущего момента:

```
class T(Module):
    def __init__(self): self.a = nn.Linear(1, 3, bias=False)

t = T()
L(t.parameters())

(#1) [Parameter containing:
tensor([-0.9595,
        [-0.8490],
        [ 0.8159]], requires_grad=True)]

type(t.a.weight)
torch.nn.parameter.Parameter
```

Мы можем создать тензор в качестве параметра, используя случайную инициализацию:

```
def create_params(size):
    return nn.Parameter(torch.zeros(*size).normal_(0, 0.01))
```

Используем это для повторного создания `DotProductBias`, но без `Embedding`:

```
class DotProductBias(Module):
    def __init__(self, n_users, n_movies, n_factors, y_range=(0,5.5)):
        self.user_factors = create_params([n_users, n_factors])
        self.user_bias = create_params([n_users])
        self.movie_factors = create_params([n_movies, n_factors])
        self.movie_bias = create_params([n_movies])
        self.y_range = y_range

    def forward(self, x):
        users = self.user_factors[x[:,0]]
        movies = self.movie_factors[x[:,1]]
        res = (users*movies).sum(dim=1)
        res += self.user_bias[x[:,0]] + self.movie_bias[x[:,1]]
        return sigmoid_range(res, *self.y_range)
```

Еще раз проведем обучение, чтобы проверить, получим ли мы такой же результат, что и в предыдущем разделе:

```
model = DotProductBias(n_users, n_movies, 50)
learn = Learner(dls, model, loss_func=MSELossFlat())
learn.fit_one_cycle(5, 5e-3, wd=0.1)
```

epoch	train_loss	valid_loss	time
0	0.962146	0.936952	00:14
1	0.858084	0.884951	00:14
2	0.740883	0.838549	00:14
3	0.592497	0.823599	00:14
4	0.473570	0.824263	00:14

А теперь посмотрим, чему наша модель научилась.

Интерпретация вложений и смещений

Наша модель уже достаточно эффективна в своей способности предоставлять рекомендации для наших пользователей, но при этом очень интересно узнать, какие же параметры она обнаружила. Пуще всего интерпретировать смещения. Вот фильмы с наименьшими значениями векторе смещений:

```
movie_bias = learn.model.movie_bias.squeeze()
idxs = movie_bias.argsort()[:5]
[dls.classes['title'][i] for i in idxs]

['Children of the Corn: The Gathering (1996)',
 'Lawnmower Man 2: Beyond Cyberspace (1996)',
 'Beautician and the Beast, The (1997)',
 'Crow: City of Angels, The (1996)',
 'Home Alone 3 (1997)']
```

Подумайте о том, что это значит. Здесь говорится, что, несмотря на высокое соответствие пользователя латентным факторам приведенных фильмов (которые, как мы вскоре увидим, представляют уровень экшена, возрастную категорию фильма и т. д.), сам фильм, как правило, пользователю не нравится. Мы могли бы просто отсортировать фильмы непосредственно по их средней оценке, но, глядя на обученное смещение, становится понятно кое-что более интересное. Это говорит нам не только о том, что фильм относится к категории тех, которые пользователи не любят смотреть, но и о том, что пользователи не хотели бы смотреть такой фильм, даже если он относится к жанру, который их интересует. А вот аналогичный результат, но уже для фильмов с самым высоким смещением:

```
idxs = movie_bias.argsort(descending=True)[:5]
[dls.classes['title'][i] for i in idxs]

['L.A. Confidential (1997)',
 'Titanic (1997)',
 'Silence of the Lambs, The (1991)',
 'Shawshank Redemption, The (1994)',
 'Star Wars (1977)']
```

То есть, например, даже если обычно вы не цените детективные фильмы, *LA Confidential* может вам понравиться.

Напрямую интерпретировать матрицы вложений не так-то просто. Здесь существует слишком много факторов для просмотра. Но при этом есть техника, которая способна выделить в такой матрице наиболее важные определяющие направления. Называется эта техника *методом главных компонентов* (principal component analysis – PCA).

Подробно рассматривать эту технику мы не станем, потому что ее понимание не особо важно для успешной практики глубокого обучения. Но если вам интересно, то предлагаем ознакомиться с курсом fast.ai под названием «Вычислительная линейная алгебра для программистов» (<https://oreil.ly/NLj2R>). На рис. 8.3 показано, как выглядят фильмы на основе двух самых сильных компонентов PCA.

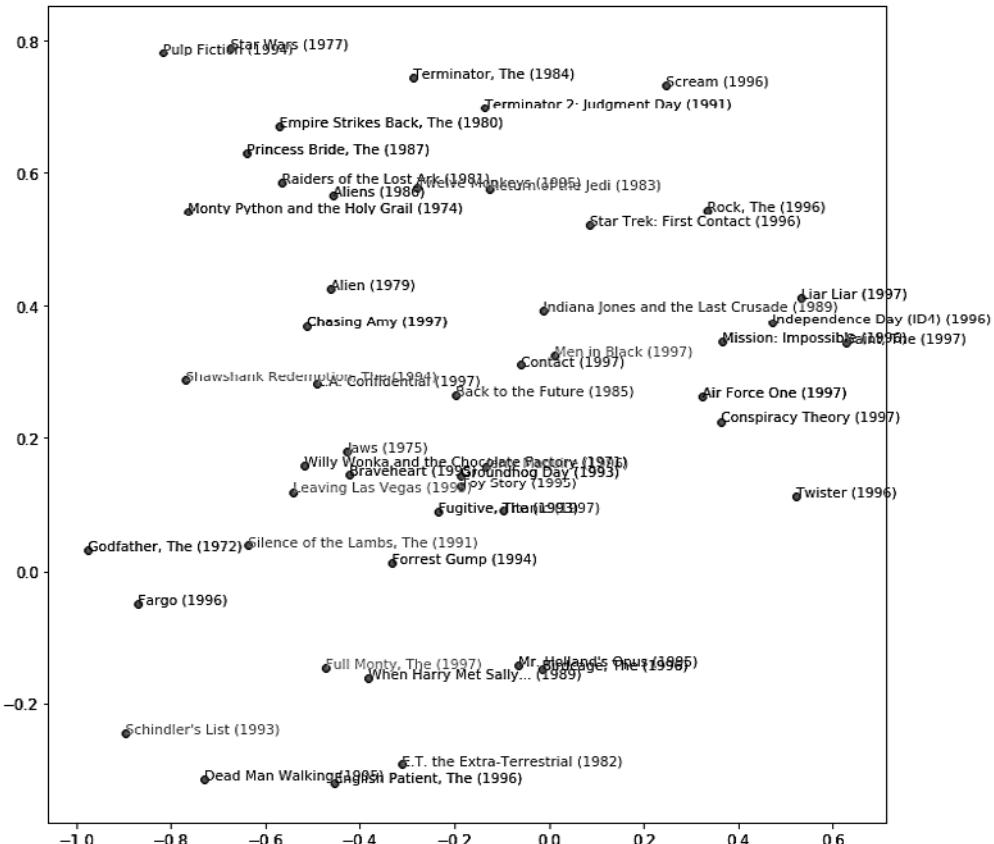


Рис. 8.3. Представление фильмов на основе двух самых сильных компонентов PCA

Мы видим, что модель обнаружила принцип *классического кино* в противовес фильмам о поп-культуре или, возможно, здесь представлены фильмы, получившие *признание критиков*.



СЛОВО ДЖЕРЕМИ

Сколько бы моделей я ни обучал, меня всегда поражает то, как этим случайно инициализируемым набором чисел, которые обучаются простыми механиками, удается самостоятельно узнавать подробности о моих данных. Это будто жульничество: я могу создать выполняющий полезные действия код, не объясняя ему, как именно эти действия нужно выполнять.

Мы создали модель с нуля с целью показать вам, что кроется у нее внутри, но вы можете использовать для ее создания непосредственно библиотеку fastai. Как раз далее мы и рассмотрим этот процесс.

Использование fastai.collab

Мы можем создать и обучить модель совместной фильтрации с помощью той же показанной ранее структуры, используя `collab_learner`:

```
learn = collab_learner(dls, n_factors=50, y_range=(0, 5.5))
learn.fit_one_cycle(5, 5e-3, wd=0.1)
```

epoch	train_loss	valid_loss	time
0	0.931751	0.953806	00:13
1	0.851826	0.878119	00:13
2	0.715254	0.834711	00:13
3	0.583173	0.821470	00:13
4	0.496625	0.821688	00:13

Просмотреть названия слоев можно так:

```
learn.model

EmbeddingDotBias(
    (u_weight): Embedding(944, 50)
    (i_weight): Embedding(1635, 50)
    (u_bias): Embedding(944, 1)
    (i_bias): Embedding(1635, 1)
)
```

Их можно использовать для воспроизведения любых видов анализов, которые мы делали в предыдущем разделе, например:

```
movie_bias = learn.model.i_bias.weight.squeeze()
idxs = movie_bias.argsort(descending=True)[:5]
[dls.classes['title'][i] for i in idxs]

['Titanic (1997)',
 "Schindler's List (1993)",
 'Shawshank Redemption, The (1994)',
 'L.A. Confidential (1997)',
 'Silence of the Lambs, The (1991)']
```

Кроме того, мы можем просмотреть еще один интересный показатель, а именно *расстояние* между этими обученными вложениями.

Расстояние между вложениями

На двумерной карте мы можем вычислить расстояние между двумя точками координат, используя формулу Пифагора: $\sqrt{x^2 + y^2}$ (предполагая, что x и y – это расстояния между координатами на каждой оси). Для 50-мерного вложения можно сделать то же самое, только сложить квадраты всех 50 расстояний между координатами.

Если мы возьмем два практически идентичных фильма, то их векторы вложений тоже будут почти идентичными, потому что пользователи, которым бы они понравились, также будут очень похожи. Здесь открывается и более общий принцип: сходство фильмов можно определить по сходству пользователей, которым эти фильмы нравятся. Это, в свою очередь, означает, что данное сходство можно определить по расстоянию между векторами вложений двух фильмов. На основе этого можно найти фильмы, наиболее схожие с *Silence of the Lambs*:

```
movie_factors = learn.model.i_weight.weight
idx = dls.classes['title'].o2i['Silence of the Lambs, The (1991)']
distances = nn.CosineSimilarity(dim=1)(movie_factors, movie_factors[idx][None])
idx = distances.argsort(descending=True)[1]
dls.classes['title'][idx]

'Dial M for Murder (1954)'
```

Мы успешно обучили модель, и пора посмотреть, как решить ситуацию, когда нет данных о пользователях. Как вообще генерировать рекомендации для новых пользователей?

Бутстрэппинг модели коллаборативной фильтрации

Самую большую сложность при использовании моделей совместной фильтрации представляет *бутстрэппинг*. В наиболее выраженной форме эта про-

блема проявляется отсутствием пользователей и, соответственно, истории, на которой можно обучать модель. Какие товары вы рекомендуете самому первому пользователю?

Но даже если вы относитесь к хорошо организованной компании с большой историей пользовательских записей, то все равно остается вопрос: что вы делаете при регистрации нового пользователя? А что делаете, когда добавляете в портфолио компании новый продукт? Для этой проблемы нет волшебного решения, да и вообще все предлагаемые нами решения являются лишь вариациями на тему «*используйте здравый смысл*». Вы могли бы присваивать новым пользователям среднее значение всех векторов вложений всех других пользователей, но здесь проблема в том, что конкретная комбинация скрытых факторов может вовсе не оказаться общей (например, среднее значение для фактора научной фантастики может быть велико, а среднее для фактора экшена мало, но найти людей, которые любят научную фантастику без экшена, достаточно сложно). Наверняка будет лучше выбрать конкретного пользователя для представления *усредненного вкуса*.

А еще лучше задействовать табличную модель, основанную на метаданных пользователя, чтобы построить начальный вектор вложений. Подумайте, какие вопросы можно задать пользователю при регистрации, чтобы лучше понять его вкусы. Тогда вы можете создать модель, в которой зависимая переменная является вектором вложений пользователя, а независимые переменные — ответами на заданные вопросы, а также метаданными из регистрационной формы пользователя. В следующем разделе вы увидите, как создавать такие табличные модели. (Вы могли заметить, что при регистрации на таких сервисах, как Pandora или Netflix, вам задают несколько вопросов относительно предпочтительных жанров кино/музыки: именно так они создают начальные рекомендации на основе совместной фильтрации).

Помимо этого, нужно следить за тем, чтобы небольшое число особо увлеченных пользователей не явились причиной генерации рекомендаций для всей базы данных пользователей. Эта проблема весьма актуальна, например, в рекомендательной системе кино. Люди, которые смотрят аниме, склонны смотреть его очень много, в основном игнорируя другие жанры и при этом активно выставляя оценки на сайтах. В результате аниме начинает очень часто фигурировать в списках *лучших фильмов всех времен*. В этом конкретном случае наличие смещения представления может быть вполне очевидно, но если смещение проявляется в скрытых факторах, то очевидным оно может и не быть.

Такая проблема способна полностью изменить внешний вид базы пользователей и поведение системы. Это особенно актуально из-за положительных петель обратной связи. Если небольшое число пользователей определяет направление рекомендаций системы, то они естественным образом при-

влекают пользователей с похожими интересами. Это, конечно же, приведет к еще большему смещению представления. Такой тип смещения имеет тенденцию к экспоненциальному росту. Вы могли встречать примеры случаев, когда руководители компаний выражают удивление по поводу того, как их онлайн-платформы резко деградировали и начинали представлять ценности, противоречащие ценностям, заложенным их основателями. При наличии подобных петель обратной связи легко видеть, что такое отклонение может произойти не только очень быстро, но и оставаться скрытым до тех пор, пока не станет слишком поздно.

В таких саморазвивающихся системах стоит воспринимать петли обратной связи скорее как норму, а не как исключение. Следует быть готовым к встрече с ними и соответствующим образом спланировать действия по их устраниению. Подумайте, как могут проявиться петли обратной связи в вашей системе и как вы можете обнаружить их в своих данных. В итоге это возвращает нас к изначальному совету о том, как избежать катастрофы при внедрении любой системы машинного обучения. В этом совете все сводится к следующему: участию в процессе людей, наличию постоянного мониторинга, а также продуманному поэтапному внедрению самой системы.

Наша модель скалярного произведения работает достаточно хорошо, являясь основой многих успешных рекомендательных систем в реальном мире. Этот подход совместной фильтрации известен как *вероятностное разложение матриц* (PMF). Еще одним подходом, который при получении таких же данных обычно работает так же хорошо, является глубокое обучение.

Глубокое обучение для коллаборативной фильтрации

Для превращения нашей архитектуры в модель глубокого обучения в первую очередь нужно взять результаты поиска по вложениям и конкатенировать эти активации. Так мы получим матрицу, которую затем сможем передать через линейные слои и нелинейные функции обычным способом.

Поскольку мы будем конкатенировать вложения, а не получать их скалярное произведение, то два таких вложения могут иметь разные размеры (разное число скрытых факторов). В fastai есть функция `get_emb_sz`, которая возвращает рекомендованные размеры вложений для ваших данных на основе эвристики, которая, по наблюдениям fast.ai, лучше всего работает на практике:

```
embs = get_emb_sz(dls)
embs
```

```
[(944, 74), (1635, 101)]
```

Реализуем этот класс:

```
class CollabNN(Module):
    def __init__(self, user_sz, item_sz, y_range=(0,5.5), n_act=100):
        self.user_factors = Embedding(*user_sz)
        self.item_factors = Embedding(*item_sz)
        self.layers = nn.Sequential(
            nn.Linear(user_sz[1]+item_sz[1], n_act),
            nn.ReLU(),
            nn.Linear(n_act, 1)) self.y_range = y_range

    def forward(self, x):
        embs = self.user_factors(x[:,0]),self.item_factors(x[:,1])
        x = self.layers(torch.cat(embs, dim=1))
        return sigmoid_range(x, *self.y_range)
```

И используем его для создания модели:

```
model = CollabNN(*embs)
```

`CollabNN` создает слои `Embedding` тем же способом, что и предыдущие классы этой главы, только теперь мы используем размеры `embs`. `self.layers` идентично мини-нейронной сети, которую мы создали в главе 4 для `MNIST`. Далее в `forward` мы применяем вложения, конкатенируем результаты и передаем итог через мини-нейронную сеть. В завершение мы применяем `sigmoid_range` так же, как и в предыдущих моделях.

Проверим ее в обучении:

```
learn = Learner(dls, model, loss_func=MSELossFlat())
learn.fit_one_cycle(5, 5e-3, wd=0.01)
```

epoch	train_loss	valid_loss	time
0	0.940104	0.959786	00:15
1	0.893943	0.905222	00:14
2	0.865591	0.875238	00:14
3	0.800177	0.867468	00:14
4	0.760255	0.867455	00:14

`fastai` предоставляет эту модель в `fastai.collab`, если вы передадите в вызов к `collab_learner` параметр `use_nn=True` (а также автоматически вызывает `get_emb_sz`), позволяя вам легко создать больше слоев. Например, здесь мы создаем два скрытых слоя с размерами 100 и 50 соответственно:

```
learn = collab_learner(dls, use_nn=True, y_range=(0, 5.5), layers=[100,50])
learn.fit_one_cycle(5, 5e-3, wd=0.1)
```

epoch	train_loss	valid_loss	time
0	1.002747	0.972392	00:16
1	0.926903	0.922348	00:16
2	0.877160	0.893401	00:16
3	0.838334	0.865040	00:16
4	0.781666	0.864936	00:16

`learn.model` является объектом типа `EmbeddingNN`. Взглянем на код `fastai` для этого класса:

```
@delegates(TabularModel)
class EmbeddingNN(TabularModel):
    def __init__(self, emb_szs, layers, **kwargs):
        super().__init__(emb_szs, layers=layers, n_cont=0, out_sz=1, **kwargs)
```

А кода-то совсем немного! Этот класс *наследует* от `TabularModel`, который и снабжает его всей функциональностью. В `__init__` он вызывает тот же метод в `TabularModel`, передавая `n_cont=0` и `out_sz=1`. В остальном он передает только те аргументы, которые получил.

KWARGS И DELEGATES

`EmbeddingNN` включает `**kwargs` в качестве параметра для `_init_`. В Python присутствие `**kwargs` в списке параметров означает «поместить любые дополнительные именованные аргументы в словарь `kwargs`». А присутствие `**kwargs` в списке аргументов означает «вставить сюда все пары ключ/значение словаря `kwargs` как именованные аргументы». Этот подход используется во многих популярных библиотеках, таких как `matplotlib`, где главная функция `plot` просто имеет сигнатуру `plot(*args, **kwargs)`. В документации к `plot` (<https://oreil.ly/P9A8T>) говорится «`kwargs` — это свойства `Line2D`», после чего эти свойства перечисляются.

Мы используем `**kwargs` в `EmbeddingNN`, чтобы избежать необходимости прописывать все аргументы в `TabularModel` второй раз и сохранить их согласованность. Тем не менее это несколько усложняет работу с API, потому что теперь Jupyter Notebook не знает, какие параметры доступны. В итоге перестанут работать такие функции, как завершение имен параметров с помощью табуляции и всплывающие списки сигнатур.

`fastai` решает эту проблему, предоставляемый особый декоратор `@delegates`, который автоматически изменяет сигнатуру класса или функции (`EmbeddingNN` в данном случае), чтобы вставить в нее все именованные аргументы.

Несмотря на то что результаты `EmbeddingNN` несколько хуже, чем при использовании подхода со скалярным произведением (что показывает преимущество

щество внимательного построения архитектуры для предметной области), он позволяет нам делать кое-что очень важное: теперь мы можем напрямую внедрять информацию о других пользователях и фильмах, датах и времени, а также любую другую, которая окажется подходящей для данной рекомендации. Именно это и делает `TabularModel`. В действительности теперь мы видели, что `EmbeddingNN` — это просто `TabularModel` с `n_cont=0` и `out_sz=1`. Это значит, нам стоит уделить время изучению `TabularModel` и тому, как его использовать для получения наилучших результатов. Этим мы и займемся в следующей главе.

Резюме

Изучая нашу первую задачу, не связанную с компьютерным зрением, мы рассмотрели рекомендательную систему и увидели, как с помощью градиентного спуска модель может изучать характерные факторы или искажения элементов из истории оценок, в результате чего они получают способность предоставлять нам информацию о самих данных.

Кроме этого, мы создали нашу первую модель в PyTorch. Мы еще не раз этим займемся на протяжении оставшейся части книги, но сначала закончим знакомство с другими основными сферами применения глубокого обучения, перейдя к табличным данным.

Вопросник

1. Какую задачу решает совместная фильтрация?
2. Как она ее решает?
3. Почему модель прогнозирования на основе совместной фильтрации может не оказаться эффективной в качестве рекомендательной системы?
4. Как выглядит представление данных для совместной фильтрации в форме перекрестной таблицы?
5. Напишите код для создания представления перекрестной таблицы данных из MovieLens (возможно, придется погуглить).
6. Что такое латентный фактор? Почему он «скрытый»?
7. Что такое скалярное произведение? Вычислите его вручную, используя чистый Python со списками.
8. Что делает `pandas.DataFrame.merge`?
9. Что такое матрица вложений?

10. Какая связь между вложением и матрицей векторов, закодированных быстро?
11. Зачем нам нужен `Embedding`, если для того же самого можно использовать векторы в быстрой кодировке?
12. Что содержит вложение до того, как мы начинаем обучение (предполагая, что мы не используем предварительно обученную модель)?
13. Создайте класс (по возможности не подглядывая) и используйте его.
14. Что возвращает `x[:, 0]`?
15. Перепишите класс `DotProduct` (по возможности не подглядывая) и обучите модель с его помощью.
16. Какая функция потерь хорошо подойдет для MovieLens и почему?
17. Что произойдет, если мы используем с MovieLens перекрестную энтропию? Как нам придется изменить модель?
18. Что, если мы задействуем в модели скалярного произведения смещение?
19. Как иначе называется сокращение весов?
20. Напишите уравнение для сокращения весов (не подглядывая!).
21. Напишите уравнение для градиента сокращения весов. Почему оно помогает уменьшать веса?
22. Почему сокращение весов ведет к лучшей обобщаемости?
23. Что в PyTorch выполняет `argsort`?
24. Дает ли сортировка смещений в фильмах тот же результат, что и усреднение общих оценок по фильму? Почему?
25. Как вывести названия слоев модели и их подробности?
26. Что означает «проблема бутстрэппинга» при использовании совместной фильтрации?
27. Как решить проблему бутстрэппинга для новых пользователей? А для новых фильмов?
28. Как могут петли обратной связи повлиять на системы, работающие по принципу совместной фильтрации?
29. Почему при использовании в совместной фильтрации нейронной сети у нас могут быть разные факторы для пользователей и фильмов?
30. Зачем в модели `CollabNN` нужен `nn.Sequential`?
31. Какой вид модели следует использовать, когда в модель совместной фильтрации нужно добавлять метаданные о пользователях и элементах или, например, дату и время?

Дополнительные задания

1. Рассмотрите отличия между Embedding-версией `DotProductBias` и версией `create_params` и постараитесь понять, почему потребовалось каждое изменение. Если не уверены, попробуйте отменить эти изменения поочереди и посмотреть, что произойдет. (Обратите внимание даже на изменение типа скобок в `forward`.)
2. Найдите три области, где используется совместная фильтрация, оценив на их примере достоинства и недостатки ее применения.
3. Заполните этот блокнот, используя полноценный датасет MovieLens, и сравните результаты с онлайн-бенчмарками. Посмотрите, удастся ли вам повысить точность. В поиске идей можете обратиться к сайту книги и форумам fast.ai. Обратите внимание, что в датасете есть и другие столбцы, — попробуйте использовать их тоже (возможные варианты можно будет почерпнуть из следующей главы).
4. Создайте модель для MovieLens, работающую с функцией перекрестной энтропии, и сравните ее с моделью из этой главы.

ГЛАВА 9

Табличное моделирование

Табличные модели принимают данные в форме таблицы (например, электронной таблицы или CSV). Задача — спрогнозировать значение в столбце на основе значений в других столбцах. В текущей главе мы рассмотрим не только глубокое обучение, но также и более общие техники ML, такие как случайные леса, поскольку в определенной области задач они могут давать лучшие результаты.

Мы узнаем, как следует предварительно обрабатывать и чистить данные, а также как интерпретировать результат модели после ее обучения. Но сначала посмотрим, как с помощью вложений передавать модели, ожидающей числа, столбцы, содержащие категории.

Категориальные вложения

В таблицах некоторые столбцы могут содержать численные данные, такие как «age» (возраст), в то время как в остальных будут строчные значения, например «sex» (пол). Численные данные можно передавать модели непосредственно (при желании можно их предварительно обработать), но столбцы в другом формате необходимо сначала преобразовать в числа. Поскольку значения в этих столбцах соответствуют разным категориям, мы зачастую называем такой тип переменных *категориальными*. Переменные первого типа называются *непрерывными*.



ТЕРМИН: НЕПРЕРЫВНЫЕ И КАТЕГОРИАЛЬНЫЕ ПЕРЕМЕННЫЕ

Непрерывные переменные — это численные данные, такие как «age», которые допускается передавать напрямую в модель, так как их можно складывать и умножать непосредственно. Категориальные переменные содержат ряд дискретных уровней, таких как movie ID, для которых сложение и умножение бессмысленно (даже если они сохранены в виде чисел).

В конце 2015 года на Kaggle компанией Rossmann проводилось соревнование по прогнозированию продаж (https://oreil.ly/U85_1). Участникам соревнования был

предоставлен обширный набор данных по различным магазинам в Германии. Задачей же было попытаться спрогнозировать продажи в течение определенного количества дней. Основная цель — помочь компании организовать рациональное управление запасами, удовлетворяя потребительский спрос без лишней загрузки складских резервов. В официальной обучающей выборке предоставлялся большой объем информации по магазинам. Помимо прочего, участникам разрешалось использовать дополнительные данные при условии, что эти данные будут открытыми и доступными для всех участников.

Один из золотых медалистов в одном из первых известных примеров эталонной табличной модели применил глубокое обучение. Его метод задействовал намного меньше инженерии признаков, основанной на знании предметной области, чем методы других золотых медалистов. Этот подход описан в работе *Entity Embeddings of Categorical Variables* (<https://oreil.ly/VmgoU>) («Вложения категориальных переменных в виде сущностей»). В онлайн-главе на сайте книги (<https://book.fast.ai/>) мы показываем, как повторить его с нуля и достичь такой же точности, какая приводится в работе. В аннотации к этому труду его авторы Чен Го (Cheng Guo) и Феликс Бекхан (Felix Bekhahn) пишут:

Вложение сущностей не только уменьшает потребление памяти и ускоряет нейронные сети в сравнении с быстрым кодированием, но, что более важно, путем отображения схожих значений, близких друг к другу в пространстве вложений, оно раскрывает внутренние свойства категориальных переменных... [Это] особенно полезно для датасетов с большим количеством признаков с высокой кардинальностью, где другие методы склонны приводить к переобучению... Поскольку вложение сущностей определяет меру расстояния для категориальных переменных, его можно использовать для визуализации категориальных данных и для кластеризации данных.

Мы уже отмечали все эти моменты, когда создавали модель коллаборативной фильтрации. Теперь же становится очевидно, что эти открытия выходят далеко за ее рамки.

В данной работе также указывается, что (как мы говорили в предыдущей главе) слой вложений в точности равнозначен размещению стандартного линейного слоя после каждого вводного слоя, закодированного быстро. Авторы для демонстрации этой равнозначности использовали схему, приведенную на рис. 9.1. Имейте в виду, что термин *dense layer* означает то же, что «линейный слой», а *one-hot encoding layer* (слои в быстрой кодировке) представляют собой вводы.

Это открытие важно, потому что мы уже знаем, как обучать линейные слои, а значит, это показывает, что с точки зрения архитектуры и алгоритма обучения слой вложений — это просто другой слой. Мы также видели это на практике в предыдущей главе, когда создавали нейронную сеть коллаборативной фильтрации, которая выглядит в точности как эта схема.

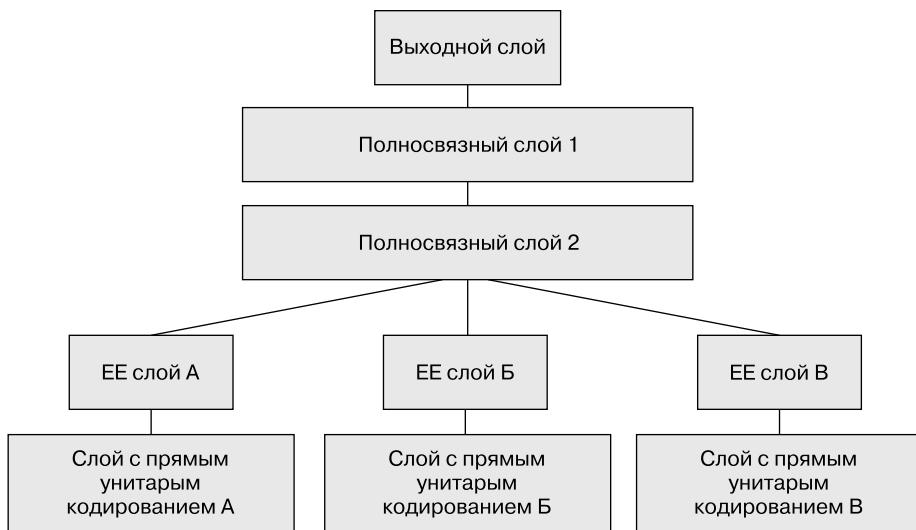


Рис. 9.1. Вложения сущностей в нейронной сети (схема любезно предоставлена Ченом Го и Феликсом Берханом)

Точно так же, как мы анализировали веса вложений для оценок фильмов, авторы работы проанализировали веса вложений для их модели прогнозирования продаж. При этом они обнаружили весьма удивительный факт, который отражает их второе ключевое открытие: вложение преобразует категориальные переменные во вводы, которые являются непрерывными и в то же время несут смысл.

Изображения на рис. 9.2 отражают эти идеи. Они основаны на подходах, использованных в этой работе, а также на некоторой добавленной нами аналитике.

Слева показан график матрицы вложений для возможных значений категории *State*. Возможные значения категориальной переменной мы называем «уровнями» (или «категориями», или «классами»), поэтому здесь один уровень — это «Berlin», следующий — «Hamburg» и т. д. Справа приведена карта Германии. Фактические местоположения германских земель не являлись частью предоставленных данных, тем не менее модель сама выучила, где они должны находиться, взяв за основу одни только данные о продажах в магазинах.

Помните ли вы наш разговор о *расстоянии* между вложениями? Авторы рассматриваемой работы нарисовали график расстояний между вложениями магазинов, сопоставив их с фактическими географическими расстояниями между этими магазинами (рис. 9.3), и обнаружили таким образом очень близкое соответствие!

Мы даже пробовали отобразить на графике вложения для дней недели и месяцев, выяснив, что те дни и месяцы, которые находятся близко друг к другу в календаре, также оказались близки и во вложениях. Это показано на рис. 9.4.

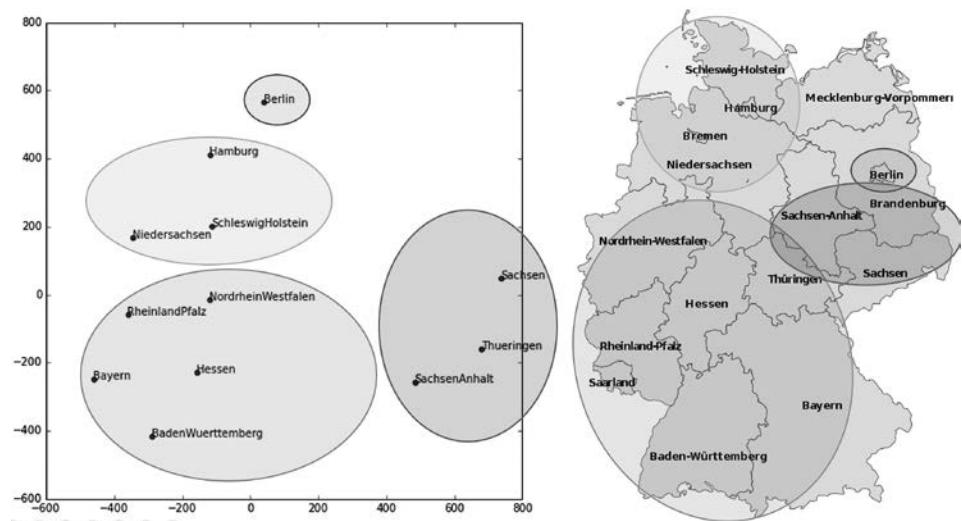


Рис. 9.2. Вложения для категории «Земли Германии» и карта (любезно предоставлены Ченом Го и Феликсом Бекханом)

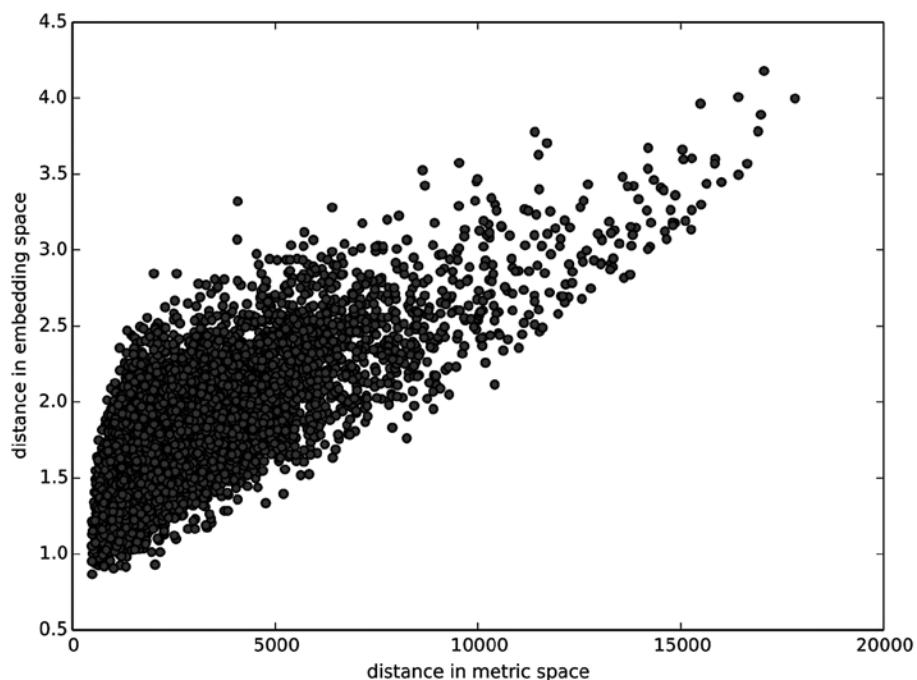


Рис. 9.3. Расстояния между магазинами (график любезно предоставлен Ченом Го и Феликсом Бекханом)

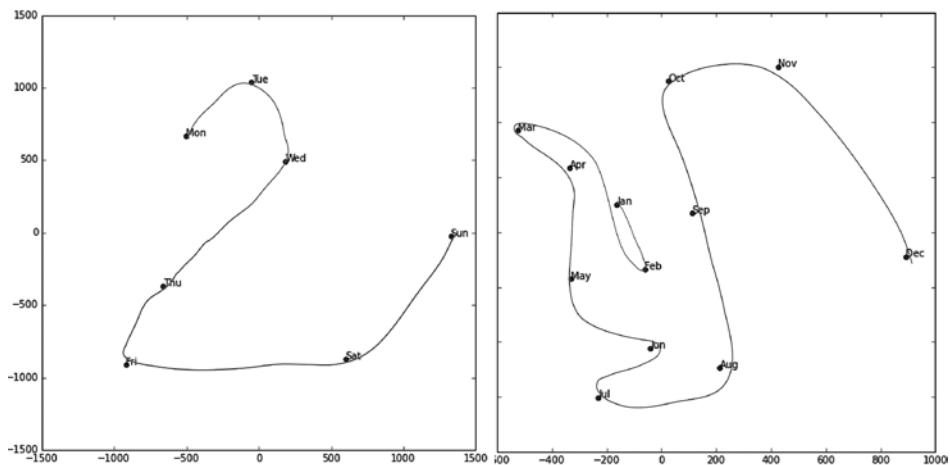


Рис. 9.4. Вложения дат (график любезно предоставлен Ченом Го и Феликсом Бекханом)

В обоих этих примерах выделяется то, что мы передаем модели конкретно категориальные данные о дискретных сущностях (например, земли Германии или дни недели), после чего модель изучает вложение для этих сущностей, которое определяет непрерывное обозначение расстояния между ними. Поскольку расстояние между вложениями было выучено на основе реальных паттернов в данных, оно склонно соответствовать нашему представлению.

Помимо этого, непрерывность вложений ценна сама по себе, потому что модели лучше понимают именно непрерывные переменные. Это неудивительно, учитывая, что модели построены из множества непрерывных весов параметров и непрерывных значений активаций, которые обновляются посредством градиентного спуска (алгоритм обучения для нахождения минимумов непрерывных функций).

Еще одно преимущество в том, что мы можем легко совмещать непрерывные значения вложений с истинными непрерывными входными данными: мы просто конкатенируем эти переменные и передаем результат в первый линейный слой. Другими словами, сырье категориальные данные преобразуются слоем вложений перед тем, как он взаимодействует с сырьем непрерывными входными данными. Именно так fastai, а также Го и Бекхан обрабатывают табличные модели, содержащие непрерывные и категориальные переменные.

К примерам использования этого подхода с конкатенацией можно отнести то, как Google формирует свои рекомендации в Google Play, что объясняется в работе *Wide & Deep Learning for Recommender Systems* (<https://oreil.ly/wsnvQ>) («Широкое и глубокое обучение для рекомендательных систем»). Эта система показана на рис. 9.5.

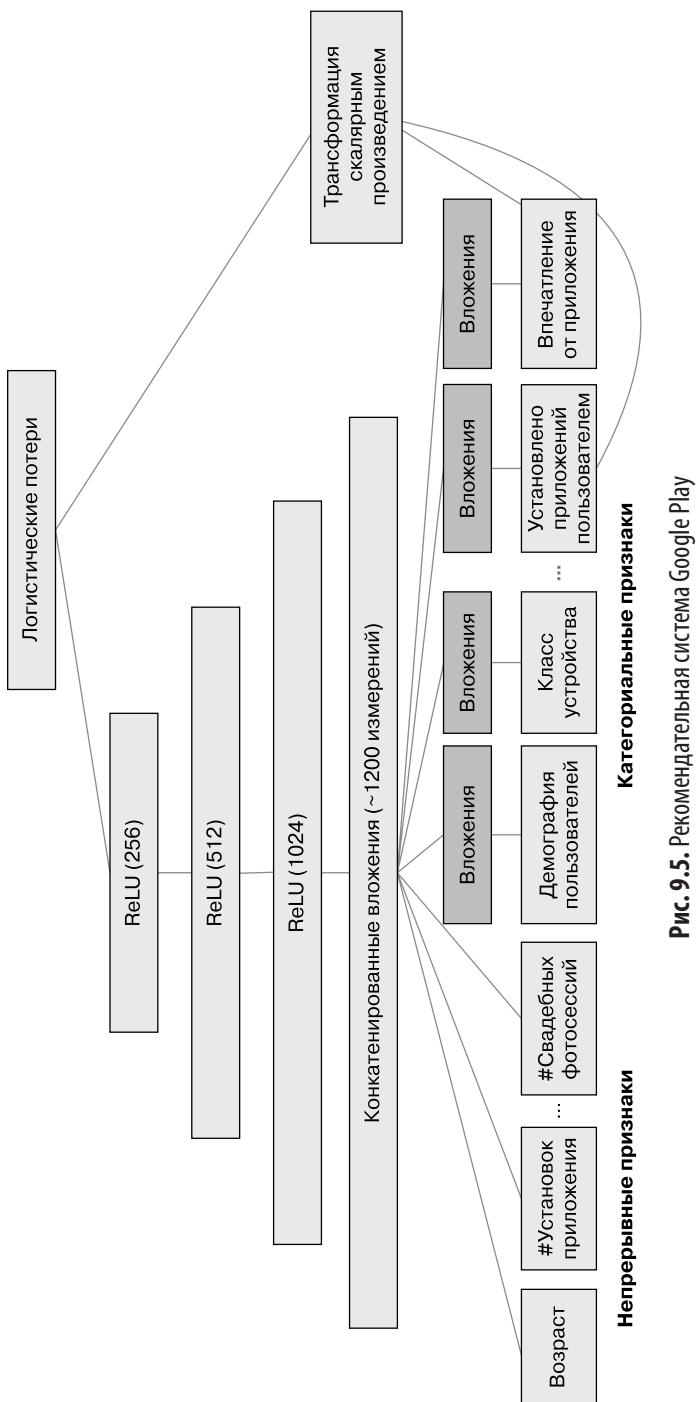


Рис. 9.5. Рекомендательная система Google Play

Интересно, что команда Google совместила оба подхода, которые мы видели в предыдущей главе: скалярное произведение (которое они называют *перекрестным произведением*) и подходы с нейронными сетями.

Давайте сделаем небольшую паузу. До этого момента решением для всех задач моделирования было *обучение модели глубокого обучения*. И в самом деле, это очень хорошее правило для сложных неструктурированных данных вроде изображений, звуков, текста естественного языка и т. д. Глубокое обучение также отлично работает для коллаборативной фильтрации. Но для анализа табличных данных оно не всегда является наилучшей отправной точкой.

За гранью глубокого обучения

Большинство курсов по ML буквально закидывают вас десятками алгоритмов — они дают краткое объяснение стоящих за ними математических принципов и приводят игрушечный пример. Вы теряетесь и не понимаете, как именно их нужно применять.

Но мы рады сообщить, что современное машинное обучение можно сузить до двух ключевых техник, которые применимы практически повсеместно. Недавние исследования показали, что огромное множество датасетов можно наилучшим образом смоделировать с помощью всего двух методов.

- Ансамблей деревьев решений (то есть случайных лесов и градиентного бустинга), используемых в основном для структурированных данных (какие можно встретить в таблицах баз данных большинства компаний).
- Многослойных нейронных сетей, обученных с помощью SGD (то есть не-глубокого и/или глубокого обучения) и используемых в основном для не-структурных данных (аудио, изображения и естественный язык).

Несмотря на то что глубокое обучение практически всегда показывает себя лучше в отношении неструктурных данных, эти два подхода склонны давать схожие результаты для многих видов структурированных данных. Тем не менее ансамбли деревьев решений, как правило, обучаются быстрее, зачастую легче интерпретируются, не нуждаются в специальном GPU-оборудовании и обычно требуют меньше настроек гиперпараметров. К тому же популярность они завоевали намного раньше техник глубокого обучения, поэтому имеют более зрелую экосистему инструментов и документации.

Более существен тот факт, что важнейший этап процесса, а именно интерпретация модели табличных данных, в ансамблях деревьев решений намного проще. Для этого существуют инструменты и методы, позволяющие получать ответы на актуальные вопросы: какие столбцы датасета оказались наиболее важными для прогнозов? Как они связаны с зависимыми переменными? Как они взаимо-

действуют друг с другом? Какие конкретные признаки были наиболее важны для конкретных наблюдений?

Следовательно, ансамбли деревьев решений — это наш первый подход для анализа нового табличного датасета.

Исключением из этого правила будут случаи, когда датасет отвечает одному из следующих условий.

- Наличие важных категориальных переменных с высокой кардинальностью (кардинальность подразумевает число дискретных уровней, представляющих категории, то есть категориальная переменная с высокой кардинальностью — это нечто вроде почтового индекса, который может иметь тысячи возможных уровней).
- Присутствуют столбцы, содержащие данные, которые будет легче понять с помощью нейронной сети, например простой текст.

На практике, когда мы имеем дело с датасетами, отвечающими этим условиям, то всегда пробуем и метод ансамблей деревьев решений, и глубокое обучение, проверяя, какой сработает лучше. Глубокое обучение наверняка окажется более полезным в нашем примере совместной фильтрации, так как в нем есть не менее двух категориальных переменных с высокой кардинальностью: пользователи и фильмы. Но на практике все выходит не столь однозначно, и зачастую будет присутствовать смесь из категориальных переменных как с высокой, так и с низкой кардинальностью вместе с непрерывными переменными.

Так или иначе очевидно, что нам нужно добавить в наш арсенал моделирования ансамбли деревьев решений.

До этого момента мы использовали PyTorch и fastai почти для всех сложных задач. Но эти библиотеки главным образом спроектированы для алгоритмов, которые выполняют много операций матричного умножения и получения производных (как раз то, что относится к глубокому обучению). Деревья решений на эти операции совсем не опираются, поэтому PyTorch здесь не особо пригодится.

Вместо этого мы будем в основном опираться на библиотеку под названием *scikit-learn* (также известную как *sklearn*). Scikit-learn — это популярная библиотека для создания моделей ML, использующая подходы, не относящиеся к глубокому обучению. Дополнительно нам понадобится выполнять обработку табличных данных и запросы, для чего мы используем библиотеку Pandas. И наконец, нам также потребуется NumPy, так как это основная числовая библиотека программирования, на которую опирается и *sklearn*, и Pandas.

У нас не хватает времени на подробное знакомство со всеми этими библиотеками, поэтому мы просто будем затрагивать самые основные элементы каждой из них. Для более подробного ознакомления мы настоятельно рекомендуем книгу

Уэса Маккинни (Wes McKinney) *Python for Data Analysis* (<http://shop.oreilly.com/product/0636920050896.do>) («Python для аналитики данных») (O'Reilly). Маккинни является создателем Pandas, так что вы можете быть уверены в точности изложенной в его книге информации.

Перейдем к сбору нужных данных.

Датасет

Датасет, который мы будем использовать в этой главе, взят из соревнования Kaggle по созданию справочника цен (Голубой книги) на подержанные бульдозеры. Целью соревнования было спрогнозировать стоимость продажи конкретной единицы спецтехники на аукционе, исходя из ее износа, типа и комплектации. Представленные для рассмотрения данные взяты из опубликованных результатов аукционов и включают в себя информацию о состоянии спецтехники, а также ее комплектации.

Это очень распространенный тип датасета и задачи прогнозирования, аналогичные которым вы можете встретить и в своем проекте или организации. Рассматриваемый набор данных доступен для загрузки на Kaggle, сайте, проводящем соревнования в области data science.

Соревнования Kaggle

Kaggle — это замечательный ресурс для целеустремленных специалистов по данным и всех тех, кто ищет возможность отточить свои навыки в машинном обучении. Ничто так не поможет повысить мастерство, как решение практических задач с одновременным получением обратной связи.

Kaggle предоставляет:

- интересные датасеты;
- обратную связь по вашим результатам;
- таблицу лидеров, по которой можно оценить хорошие результаты: достижимые и топовый эталонный уровень;
- публикации в блоге от победителей соревнований, дающих полезные советы и рассказывающих об эффективных техниках.

До сих пор все наши датасеты были доступны для скачивания через интегрированную систему датасетов fastai. Текущий набор данных доступен только на Kaggle. Поэтому вам понадобится зарегистрироваться на их сайте, перейти на страницу соревнования (<https://oreil.ly/B9wfd>), щелкнуть на Rules, а затем на

I Understand and Accept. (Несмотря на то что данное соревнование завершилось и участвовать вы в нем не будете, вам все равно нужно принять правила для получения доступа к загрузке данных.)

Самый простой способ скачивания датасетов Kaggle – это использовать Kaggle API. Вы можете установить его с помощью команды `pip`, выполнив в ячейке блокнота следующую инструкцию:

```
!pip install kaggle
```

Для использования API Kaggle вам понадобится ключ API, который можно получить, щелкнув на картинке профиля на сайте, выбрав **My Account**, а затем щелкнув на **Create New API Token**, – в результате на ваш ПК будет сохранен файл `kaggle.json`. Далее нужно скопировать этот ключ на GPU-сервер. Для этого откройте скачанный файл, скопируйте содержимое и вставьте его в одиночные кавычки в следующей ячейке блокнота этой главы (например, `creds = '{"username": "xxx", "key": "xxx"}'`):

```
creds = ''
```

Затем выполните эту ячейку (выполнение потребуется только один раз):

```
cred_path = Path('~/kaggle/kaggle.json').expanduser()
if not cred_path.exists():
    cred_path.parent.mkdir(exist_ok=True)
    cred_path.write(creds)
    cred_path.chmod(0o600)
```

Теперь вы можете скачивать датасеты, для этого осталось только выбрать нужный путь их сохранения:

```
path = URLs.path('bluebook')
path
Path('/home/sgugger/.fastai/archive/bluebook')
```

После чего можете использовать Kaggle API для скачивания и распаковки датасета:

```
if not path.exists():
    path.mkdir()
    api.competition_download_cli('bluebook-for-bulldozers', path=path)
    file_extract(path/'bluebook-for-bulldozers.zip')

path.ls(file_type='text')

(#7) [Path('Valid.csv'), Path('Machine_Appendix.csv'), Path('ValidSolution.csv'),
Path('TrainAndValid.csv'), Path('random_forest_benchmark_test.csv'),
Path('Test.csv'), Path('median_benchmark.csv')]
```

Теперь, когда мы загрузили датасет, изучим его!

Знакомство с данными

На странице Data (<https://oreil.ly/oSrBi>) сайта Kaggle дается общая информация по датасету и указывается, что основные его поля, расположенные в `train.csv`, это:

- `SalesID` — уникальный идентификатор продажи;
- `MachineID` — уникальный идентификатор машины. Машина может быть продана несколько раз;
- `saleprice` — цена, по которой машина ушла с аукциона (указано только в `train.csv`);
- `saledate` — дата продажи.

При любом виде работы, связанной с исследованием данных, важно *непосредственно изучить эти данные*, убедившись, что вы понимаете их формат, как они сохранены, какие типы значений содержат и т. д. Имейте в виду, что даже если вы прочли описание, фактические данные могут оказаться не такими, как вы ожидаете. Мы начнем с чтения обучающей выборки в Pandas DataFrame. Обычно рекомендуется сразу установить параметр `low_memory=False` и не менять его до тех пор, пока Pandas реально не исчерпает ресурс памяти и не выдаст ошибку. Этот параметр, который по умолчанию установлен как `True`, дает Pandas команду при определении типа данных каждого столбца просматривать только несколько его строк за раз. Это может привести к тому, что Pandas в разных строках будет использовать разные типы данных, что, в свою очередь, приведет либо к ошибкам обработки данных, либо к последующим проблемам при обучении модели.

Загрузим наши данные и взглянем на столбцы:

```
df = pd.read_csv(path/'TrainAndValid.csv', low_memory=False)

df.columns

Index(['SalesID', 'SalePrice', 'MachineID', 'ModelID', 'datasource',
       'auctioneerID', 'YearMade', 'MachineHoursCurrentMeter', 'UsageBand',
       'saledate', 'fiModelDesc', 'fiBaseModel', 'fiSecondaryDesc',
       'fiModelSeries', 'fiModelDescriptor', 'ProductSize',
       'fiProductClassDesc', 'state', 'ProductGroup', 'ProductGroupDesc',
       'Drive_System', 'Enclosure', 'Forks', 'Pad_Type', 'Ride_Control',
       'Stick', 'Transmission', 'Turbocharged', 'Blade_Extension',
       'Blade_Width', 'Enclosure_Type', 'Engine_Horsepower', 'Hydraulics',
       'Pushblock', 'Ripper', 'Scarifier', 'Tip_Control', 'Tire_Size',
       'Coupler', 'Coupler_System', 'Grouser_Tracks', 'Hydraulics_Flow',
       'Track_Type', 'Undercarriage_Pad_Width', 'Stick_Length', 'Thumb',
       'Pattern_Changer', 'Grouser_Type', 'Backhoe_Mounting', 'Blade_Type',
       'Travel_Controls', 'Differential_Type', 'Steering_Controls'],
      dtype='object')
```

Многовато столбцов! Лучше просмотреть весь датасет, чтобы понять, какая информация находится в каждом из них. Вскоре мы увидим, как «сосредоточиться» на самых интересных.

Правильным будет начать с обработки *порядковых столбцов*. Порядковые столбцы — это те, которые содержат строки или аналогичные данные, имеющие естественный порядок. Например, вот уровни `ProductSize`:

```
df['ProductSize'].unique()  
  
array([nan, 'Medium', 'Small', 'Large / Medium', 'Mini', 'Large', 'Compact'],  
      dtype=object)
```

О подходящем порядке этих уровней можно сообщить Pandas так:

```
sizes = 'Large','Large / Medium','Medium','Small','Mini','Compact'  
  
df['ProductSize'] = df['ProductSize'].astype('category')  
df['ProductSize'].cat.set_categories(sizes, ordered=True, inplace=True)
```

Самым важным столбцом данных является зависимая переменная: та, которую мы будем прогнозировать. Вспомните, что метрика модели — это функция, отражающая качество прогнозов. Само собой разумеется, что ее выбор является значительной частью настройки проекта. Во многих случаях выбор удачной метрики станет больше напоминать процесс проектирования — будет недостаточно подобрать уже существующую переменную. Вам нужно внимательно подумать о том, какая метрика или набор метрик реально измеряют важный для вас показатель эффективности модели. Если эту метрику не представляет никакая переменная, то следует оценить возможность ее создания из доступных переменных.

Как бы то ни было, но в данном случае выбор метрики уже предписан Kaggle: мы берем среднеквадратичную логарифмическую ошибку (RMLSE) между фактическими и спрогнозированными ценами аукциона. Вычисляется же она достаточно простым способом: мы получаем логарифм цен, для значения которого `m_rmse` и дает нам то, что в итоге нужно:

```
dep_var = 'SalePrice'  
  
df[dep_var] = np.log(df[dep_var])
```

Теперь мы подготовились к изучению нашего первого алгоритма машинного обучения для табличных данных: деревьев решений.

Деревья решений

В основе ансамблей деревьев решений, как и предполагает само название, лежат деревья решений, поэтому с них мы и начнем. Дерево решений задает

в отношении данных серию двоичных (Yes или No) вопросов. После каждого вопроса данные в текущей части дерева разделяются между ветками Yes и No, как показано на рис. 9.6. По итогу ответов на один или несколько вопросов либо делается конечный прогноз, либо задается дополнительный вопрос.

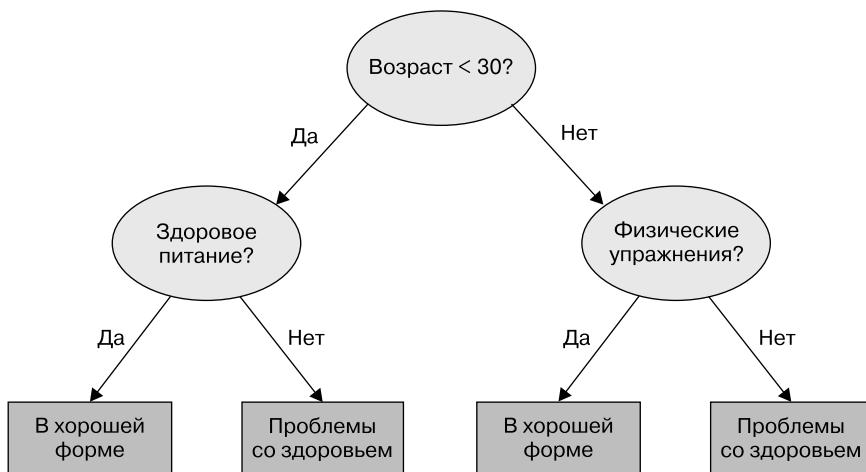


Рис. 9.6. Пример дерева решений

Далее эта последовательность вопросов будет служить процедурой получения любого элемента данных, хоть из обучающей выборки, хоть нового, и присваивания этого элемента группе. То есть после постановки вопроса и получения ответа мы можем сказать, что элемент принадлежит к той же группе, что и все остальные элементы обучающих данных, давшие на вопросы тот же набор ответов. Но какая в этом польза? Цель нашей модели — прогнозировать значения для элементов, а не присваивать их разным группам обучающего датасета. Ценность же в том, что теперь мы можем присваивать значение прогноза для каждой из этих групп — для регрессии мы берем среднее целевые элементов группы.

Рассмотрим, как мы находим нужные вопросы. Конечно же, создавать их все самим нам не захочется, и для этого у нас есть компьютеры! При этом основные шаги обучения дерева решений представить очень просто.

1. Поочередно перебрать каждый столбец датасета.
2. В каждом столбце поочередно перебрать каждый возможный уровень.
3. Попробовать разделить данные на две группы в зависимости от того, больше они или меньше выбранного среднего значения (если же это категориальная переменная, то в зависимости от того, равны они этому уровню или нет).

4. Найти среднюю цену продажи для каждой из этих групп и посмотреть, насколько она близка к фактической цене каждой единицы спецтехники в этой группе. Рассмотреть это как очень примитивную «модель», в которой прогнозы просто представляют среднюю цену продажи в группе этого элемента.
5. После перебора всех столбцов и всех их возможных уровней выбрать точку разделения, в которой с помощью этой простой модели были получены наилучшие прогнозы.
6. Теперь на основе выбранной точки разделения у нас есть для данных две группы. Далее нужно рассмотреть каждую группу как отдельный датасет и также найти в каждом из них наилучшую точку разделения путем возвращения к шагу 1 для каждой этой группы.
7. Продолжать этот процесс рекурсивно, пока не будет достигнут критерий остановки для каждой группы: например, прекратить дальнейшее разделение группы, когда в ней останется всего 20 элементов.

Несмотря на то что это достаточно простой алгоритм для самостоятельной реализации (что было бы также неплохим упражнением), мы можем сэкономить время, использовав реализацию, созданную sklearn.

Но для начала нам нужно подготовить данные.



СЛОВО АЛЕКСИСУ

Вот вам полезный вопрос, над которым стоит подумать. Если вы представите, что процедура по определению дерева решений, по сути, выбирает одну последовательность разделяющих вопросов о переменных, то можете поинтересоваться: откуда нам знать, что она выбирает верную последовательность? Правило подразумевает выбор вопроса, производящего наилучшее разделение (то есть наиболее точно разделяющего элементы на разные категории), и последующее повторяющееся применение этого же правила к группам, получаемым в процессе этого разделения. В компьютерной науке данная техника называется «жадным» подходом. Сможете ли вы представить себе сценарий, в котором постановка «менее эффективного» разделяющего вопроса приводила бы к лучшему последующему разделению и к улучшению результатов в целом?

Обработка дат

Первым делом для подготовки данных нам нужно обогатить представление дат. Основой только что описанного дерева решений является бисекция — разделение группы на две. Мы рассматриваем порядковые переменные и разделяем датасет на основе того, выше или ниже порога находятся их значения. Мы также рассматриваем категориальные переменные, разделяя датасет на основе соответствия их уровней конкретному уровню. Так что этот алгоритм

предоставляет способ разделения датасета на основе как порядковых, так и категориальных данных.

Но как это применить к такому стандартному типу данных, как дата? Вы можете предположить, что данные нужно рассматривать как порядковые значения, потому что можно уверенно сказать, что одна дата больше другой. Тем не менее даты несколько отличаются от большинства порядковых значений. Если говорить конкретнее, то некоторые из них качественно отличаются от других, что зачастую наблюдается в системах, которые мы моделируем.

Для содействия алгоритму в обработке дат нам нужно, чтобы модель знала больше, чем просто то, следует одна дата до либо после другой. Мы можем захотеть, чтобы модель принимала решения на основе, например, дня недели, того, является ли день выходным, какой это месяц и т. д. Для этого мы заменяем каждый столбец дат на набор столбцов метаданных дат, таких как праздники, день недели и месяц. Эти столбцы предоставляют категориальные данные, которые нам пригодятся.

В fastai есть функция, которая сделает это все за нас, нужно лишь передать название столбца с датами:

```
df = add_datepart(df, 'saledate')
```

Давайте сразу сделаем то же самое для тестовой выборки:

```
df_test = pd.read_csv(path/'Test.csv', low_memory=False)
df_test = add_datepart(df_test, 'saledate')
```

Теперь у нас в DataFrame появилось много новых столбцов:

```
' '.join(o for o in df.columns if o.startswith('sale'))
'saleYear saleMonth saleWeek saleDay saleDayofweek saleDayofyear
> saleIs_month_end saleIs_month_start saleIs_quarter_end saleIs_quarter_start
> saleIs_year_end saleIs_year_start saleElapsed'
```

Это хорошее начало, но нам еще нужно будет заняться очисткой данных, для чего мы используем fastai-объекты `TabularPandas` и `TabularProc`.

Использование `TabularPandas` и `TabularProc`

Вторым этапом подготовки мы убедимся, что можем обрабатывать строки и пропуски данных. По умолчанию sklearn не может ни того ни другого. Вместо этого мы используем fastai-класс `TabularPandas`, который обертывает Pandas DataFrame и предоставляет ряд удобств. Для заполнения `TabularPandas` мы действуем два `TabularProc` — `Categorify` и `FillMissing`. `TabularProc` аналогичен обычному `Transform`, за исключением следующего.

- Он возвращает тот же объект, что был ему передан, после изменения этого объекта на месте.
- Он выполняет преобразование один раз при изначальной передаче ему данных, а не начинает работу только при обращении к этим данным.

Categorify — это **TabularProc**, заменяющий столбец на численный категориальный столбец. **FillMissing** — это **TabularProc**, заменяющий отсутствующие значения на медиану этого столбца и создающий новый логический столбец, устанавливающий **True** для любой строки, где это значение отсутствовало. Эти два преобразования необходимы практически для любых табличных данных, которые вы будете использовать, поэтому начать обработку данных лучше всего именно с этого:

```
procs = [Categorify, FillMissing]
```

TabularPandas также выполнит за нас разделение датасета на обучающую и контрольную выборки. Тем не менее нам следует с осторожностью подойти к контрольному набору. Нужно спроектировать его так, как будто это *тестовая выборка*, которую Kaggle будет использовать для оценки соревнования.

Вспомните отличие между контрольной выборкой и тестовой, о котором мы говорили в главе 1. *Контрольная выборка* — это данные, которые мы «отстраняем» от обучения, чтобы в течение этого процесса не произошло переобучение на обучающих данных. *Тестовая выборка* — это данные, которые мы «отстраняем» еще дальше, то есть даже от самих себя, чтобы в ходе изучения разных архитектур модели и гиперпараметров не допустить переобучения в отношении уже контрольной выборки.

Тестовую выборку мы не видим, но нам нужно определить контрольную таким образом, чтобы она имела определенную связь как с обучающими данными, так и с тестовыми.

В некоторых случаях этого можно добиться простым случайным определением точек данных. Но это не наш случай, потому что речь идет о временных рядах.

Если вы посмотрите на диапазон дат, представленный в тестовой выборке, то поймете, что он охватывает период в шесть месяцев, начиная с мая 2012 года — наиболее ранней точки всей обучающей выборки. Это правильная временная структура данных, потому что спонсор соревнования захочет убедиться в том, что модель способна прогнозировать данные в будущем. Но это также значит, что если мы хотим получить полезную контрольную выборку, то нужно, чтобы она представляла более поздний период времени, чем обучающая. Обучающие данные Kaggle оканчиваются апрелем 2012 года, поэтому мы определим более узкий обучающий датасет, заполненный данными за период до ноября 2011 года, а для контрольной выборки возьмем период после ноября 2011-го.

Для этого мы задействуем `np.where`, полезную функцию, возвращающую (как первый элемент кортежа) индексы всех значений `True`:

```
cond = (df.saleYear<2011) | (df.saleMonth<10)
train_idx = np.where( cond)[0]
valid_idx = np.where(~cond)[0]

splits = (list(train_idx),list(valid_idx))
```

Нужно сообщить `TabularPandas`, какие столбцы непрерывные, а какие — категориальные. Это можно сделать автоматически с помощью вспомогательной функции `cont_cat_split`:

```
cont,cat = cont_cat_split(df, 1, dep_var=dep_var)

to = TabularPandas(df, procs, cat, cont, y_names=dep_var, splits=splits)
```

Поведение `TabularPandas` во многом аналогично `fastai`-объекту `Datasets`, включая предоставление атрибутов `train` и `valid`:

```
len(to.train),len(to.valid)

(404710, 7988)
```

Здесь мы видим, что данные по-прежнему отображаются для категорий в виде строк (приведено только несколько столбцов, потому что вся таблица слишком велика для страницы):

```
to.show(3)2
```

	state	ProductGroup	Drive_System	Endosure	SalePrice
0	Alabama	WL	#na#	EROPS w AC	11.097410
1	North Carolina	WL	#na#	EROPS w AC	10.950807
2	New York	SSL	#na#	OROPS	9.210340

Тем не менее все базовые элементы представлены числами:

```
to.items.head(3)
```

	state	ProductGroup	Drive_System	Endosure
0	1	6	0	3
1	33	6	0	3
2	32	3	0	6

Преобразование категориальных столбцов в числа выполняется простым замещением каждого уникального уровня числом. Эти связываемые с уровнями числа выбираются последовательно согласно порядку появления в столбцах, поэтому после преобразования у чисел категориальных столбцов никакого особых значения нет. Другое дело, если вы сначала преобразуете столбец в упорядоченную категорию Pandas (как мы делали для `ProductSize` ранее), тогда будет использован выбранный вами порядок. Это отображение можно просмотреть с помощью атрибута `classes`:

```
to.classes['ProductSize']
```

```
(#7) ['#na#', 'Large', 'Large / Medium', 'Medium', 'Small', 'Mini', 'Compact']
```

Так как на обработку данных для достижения текущего результата уходит около минуты, нам следует их сохранить. Это позволит в дальнейшем продолжить работу с данной точки, не повторяя все предыдущие шаги. В fastai есть метод `save`, который для сохранения практически любого объекта Python использует Python-систему *pickle* (консервация):

```
(path/'to.pkl').save(to)
```

Для прочтения этих данных в дальнейшем нужно будет ввести следующее:

```
to = (path/'to.pkl').load()
```

Теперь, когда подготовительный этап завершен, можно приступить к созданию дерева решений.

Создание дерева решений

Для начала определим независимую и зависимую переменные:

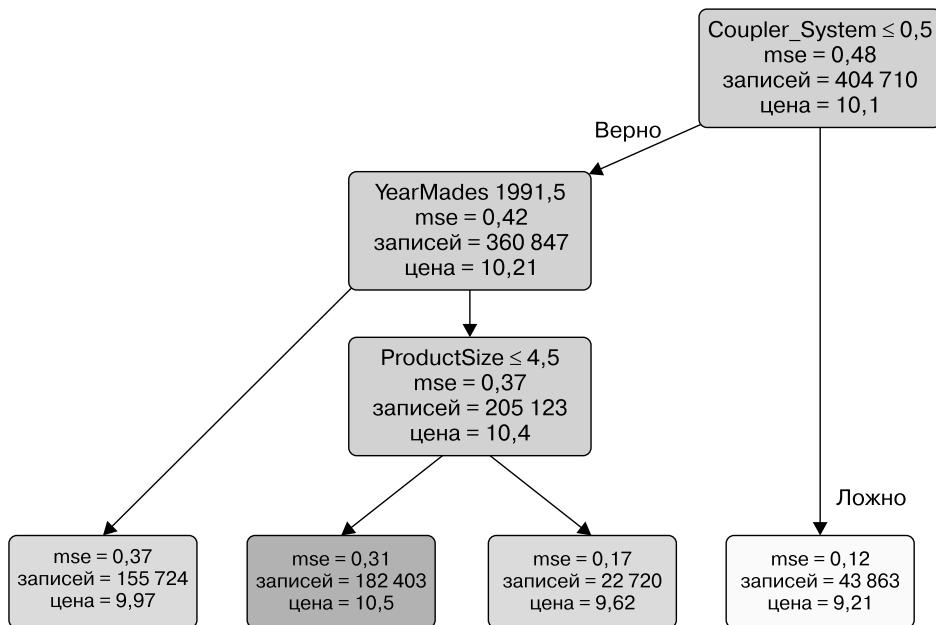
```
xs,y = to.train.xs,to.train.y  
valid_xs,valid_y = to.valid.xs,to.valid.y
```

Теперь, когда все данные представлены в численном формате и отсутствуют пропуски значений, можно создать дерево решений:

```
m = DecisionTreeRegressor(max_leaf_nodes=4)  
m.fit(xs, y);
```

Чтобы не усложнять, мы дали sklearn команду создать только четыре *конечных узла*. Для наглядного представления результатов обучения можно отобразить дерево:

```
draw_tree(m, xs, size=7, leaves_parallel=True, precision=2)
```



Понимание данной схемы — один из лучших способов понять деревья решений в принципе, поэтому мы будем объяснять все шаг за шагом, начиная с верхнего узла.

Верхний узел представляет *начальную модель* до выполнения разделений, когда все данные находятся в одной группе. Это наимпростейшая возможная модель, в которой еще не давались ответы ни на какие вопросы, поэтому она всегда будет прогнозировать значение как среднее для всего датасета. В этом случае мы можем видеть, что она прогнозирует значение логарифма цены продаж как 10,1. Среднеквадратичная ошибка определяется как 0,48, квадратный корень из которой равен 0,69. (Помните, что если вы не видите надпись `m_rmse` или *root mean squared error*, это означает, что перед вами значение до извлечения квадратного корня, то есть оно просто является средним квадрата разностей.) Мы также можем видеть, что в данной группе содержится 404 710 записей с аукционов — это общий размер обучающей выборки. Последний элемент информации — наилучшее найденное значение разделяющего критерия, которое подразумевает разделение на основе столбца `coupler_system`.

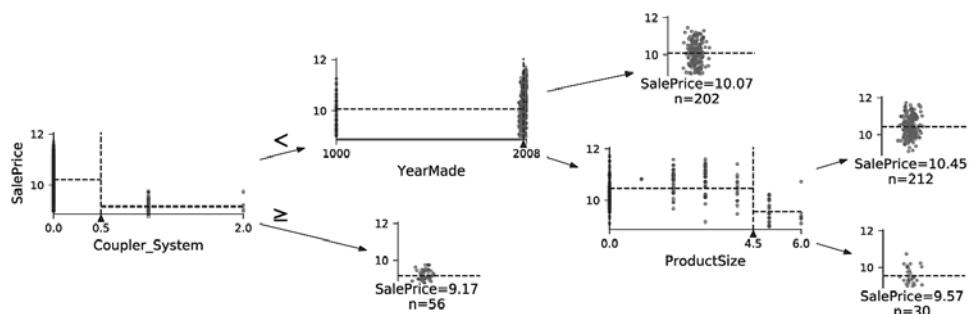
Следующий узел ниже слева показывает, что есть 360 847 записей аукционов для оборудования, в которых `coupler_system` оказался меньше 0,5. Среднее значение зависимой переменной в этой группе — 10,21. Теперь переходим ниже вправо под начальную модель и оказываемся в группе записей, где `coupler_system` оказался выше 0,5.

Нижний ряд содержит *конечные узлы*: узлы, не дающие ответов, потому что вопросов для них не осталось. Крайний правый узел в этом ряду содержит записи, где `coupler_system` оказался больше 0,5. Среднее значение здесь зафиксировано как 9,21, и это показывает, что алгоритм дерева решений нашел одно двоичное решение, которое отделило высокие результаты аукционов от низких. Спрашивая только о `coupler_system`, модель спрогнозировала значение 9,21 против 10,1.

Возвращаясь к верхнему узлу, следующему за первой точкой принятия решения, мы видим, что второе двоичное решение о разделении было принято на основе вопроса о том, является `YearMade` меньшим или равным 1991,5. Для группы, где ответ был `true` (помните, что теперь она следует двум двоичным решениям, принятым на основе `coupler_system` и `YearMade`), среднее значение — 9,97, а записей в этой группе содержится 155 724. Для группы аукционов, где это решение оказалось `false`, среднее значение определено как 10,4, присутствует же в ней 205 123 записи. И снова видно, что алгоритм дерева решений успешно разделил наиболее дорогостоящие записи аукционов еще на две группы, существенно отличающиеся значением.

Теперь можно представить ту же информацию с помощью детища Теренса Парра (Terence Parr) — мощной библиотеки dtreeviz (<https://oreil.ly/e9KrM>):

```
samp_idx = np.random.permutation(len(y))[:500]
dtreeviz(m, xs.iloc[samp_idx], y.iloc[samp_idx], xs.columns, dep_var,
         fontname='DejaVu Sans', scale=1.6, label fontsize=10,
         orientation='LR')
```



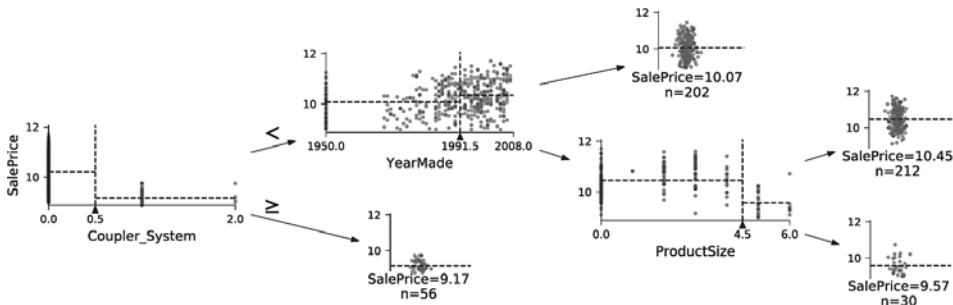
Здесь у нас график распределения данных для каждой точки разделения. Видно, что с данными `YearMade` есть проблема: очевидно, что есть бульдозеры, выпущенные в 1000 году. Можно предположить, что это просто код для отсутствующего значения (значение, не появляющееся в данных и используемое в качестве плейсхолдера на случай отсутствия значения). В целях моделирования 1000 год вполне подойдет, но как вы видите, этот выброс

существенно усложняет визуализацию интересующих нас значений. Поэтому заменим его на 1950:

```
xs.loc[xs['YearMade'] < 1900, 'YearMade'] = 1950
valid_xs.loc[valid_xs['YearMade'] < 1900, 'YearMade'] = 1950
```

Внесенное изменение делает разделение намного более ясным, несмотря на то что результат модели существенно оно не меняет. Это отличный пример отказоустойчивости деревьев решений в отношении проблем с данными.

```
m = DecisionTreeRegressor(max_leaf_nodes=4).fit(xs, y)
dtreeviz(m, xs.iloc[samp_idx], y.iloc[samp_idx], xs.columns, dep_var,
         fontname='DejaVu Sans', scale=1.6, label_fontsize=10,
         orientation='LR')
```



Теперь поручим алгоритму построить дерево побольше. Здесь мы не передаем такие критерии остановки, как `max_leaf_nodes`:

```
m = DecisionTreeRegressor()
m.fit(xs, y);
```

Мы создадим небольшую функцию для проверки среднеквадратичной ошибки модели (`m_rmse`), поскольку именно так определяли победителя соревнований:

```
def r_mse(pred,y): return round(math.sqrt(((pred-y)**2).mean()), 6)
def m_rmse(m, xs, y): return r_mse(m.predict(xs), y)
```

```
m_rmse(m, xs, y)
```

```
0.0
```

Итак, наша модель идеальна, не правда ли? Но не будем спешить с выводами, ведь нужно проверить модель на контрольной выборке, чтобы убедиться в отсутствии переобучения:

```
m_rmse(m, valid_xs, valid_y)
```

```
0.337727
```

Мда... существенное переобучение налицо, и вот почему:

```
m.get_n_leaves(), len(xs)
(340909, 404710)
```

У нас практически столько же конечных узлов, сколько точек данных. Похоже, мы сильно перестарались. В самом деле, предустановленные настройки sklearn позволяют ей продолжать разделение узлов до тех пор, пока в каждом из них не останется по одному элементу. Изменим правило остановки, чтобы sklearn обеспечила наличие в каждом узле не менее 25 записей:

```
m = DecisionTreeRegressor(min_samples_leaf=25)
m.fit(to.train_xs, to.train_y)
m_rmse(m, xs, y), m_rmse(m, valid_xs, valid_y)
(0.248562, 0.32368)
```

Вот так намного лучше. Давайте еще раз проверим количество конечных узлов:

```
m.get_n_leaves()
12397
```

Это намного разумнее!



СЛОВО АЛЕКСИСУ

Вот что я думаю по этому поводу. Представьте игру «Двадцать вопросов», в которой загадывающий тайком загадывает объект (например, «наш телевизор»), а угадывающий должен задать 20 вопросов, на которые можно ответить «да» или «нет», стараясь угадать, что было загадано (например, «Это больше, чем хлебница?»). Отгадывающий старается спрогнозировать не численное значение, а просто определить конкретный объект из набора всего, что может представить. Когда ваше дерево решений содержит больше конечных узлов, чем возможных объектов в предметной области, оно «по сути» становится очень хорошо обученным угадывающим. Оно выучило последовательность вопросов, необходимых для определения конкретного элемента данных в обучающей выборке, и делает «прогноз», просто описывая значение этого элемента. Именно так происходит запоминание обучающей выборки, то есть переобучение.

Построение дерева решений — хороший способ создания модели для наших данных. Оно очень гибкое, так как может ясно обрабатывать нелинейные связи и взаимодействия между переменными. Но здесь присутствует принципиальный компромисс между тем, как хорошо оно обобщается (чего можно достичь, создавая небольшие деревья), и тем, насколько оно точно в отношении обучающей выборки (достигается с помощью больших деревьев).

Так как же найти тот самый идеальный компромисс? Это мы вам покажем сразу после того, как разберемся с одной важной упущененной деталью: обработкой категориальных переменных.

Категориальные переменные

В предыдущей главе при работе с сетями глубокого обучения мы обрабатывали категориальные переменные, переводя их в унитарную кодировку и передавая в слой вложений. Этот слой вложений помогал модели выявлять значение разных уровней этих переменных (уровни категориальной переменной не имеют внутреннего значения, если только мы не зададим их порядок вручную с помощью Pandas). Но в дереве решений у нас нет слоев вложений, так как же могут необработанные переменные приносить какую-либо пользу в нем? Например, как можно использовать код продукта?

Говоря коротко, он просто работает! Представьте ситуацию, в которой один код продукта имеет на аукционе гораздо более высокую цену, чем другой. В этом случае любое двоичное разделение будет приводить к определению одного кода продукта в некую группу, которая в итоге будет представлять более высокую стоимость, чем другая группа. Следовательно, наш простой алгоритм построения дерева решений выберет это разделение. Позднее в процессе обучения он сможет дальше разделять эту подгруппу, содержащую код дорогого продукта, и со временем дерево сосредоточится именно на этом продукте.

Также можно использовать унитарную кодировку для замены одной категориальной переменной на несколько унитарно закодированных столбцов, где каждый столбец будет представлять возможный уровень переменной. В Pandas для этого используется специальный метод `get_dummies`.

Но нет никаких свидетельств того, что этот подход улучшит конечный результат. Поэтому по возможности мы его избегаем, так как в итоге он усложняет работу с датасетом. В 2019 году эта проблема разбиралась в работе Марвина Райта (Marvin Wright) и Инке Кенига (Inke König) *Splitting on Categorical Predictors in Random Forests* (<https://oreil.ly/ojzKJ>) («Разделение по категориальным предикторам в случайных лесах»):

Стандартный подход для номинальных предикторов заключается в рассмотрении всех $(2^{k-1} - 1)$ 2-разбиений категорий k предикторов. Тем не менее это экспоненциальное отношение производит для вычисления потенциально большое число разделений, повышая вычислительную сложность и ограничивая возможное количество категорий в большинстве реализаций. Для двоичной классификации и регрессии было показано, что упорядочивание категорий предикторов в каждом разделении ведет к точно таким же разделениям, что и стандартный подход. Это снижает вычислительную сложность, потому что для номинального предиктора с k категориями требуется рассмотрение только $k - 1$ разделений.

Теперь, когда вы поняли принцип работы деревьев решений, пора перейти к рассмотрению того самого идеального компромисса — случайных лесов.

Случайные леса

В 1994 году профессор Лео Брейман (Leo Breiman) из Беркли, штат Калифорния, спустя год после выхода на пенсию опубликовал небольшой технический доклад под названием *Bagging Predictors* (<https://oreil.ly/6gMuG>) («Бэггинг предикторов»), который выразил одну из наиболее важных идей в современном машинном обучении. Начинался доклад так:

Бэггинг (бутстрэп-агрегирование) предикторов — это метод создания нескольких версий предиктора с последующим их использованием для получения агрегированного предиктора. Агрегация усредняет эти версии... Несколько версий создаются путем формирования реплик обучающей выборки с помощью бутстрэпа, после чего эти реплики используются в качестве новых обучающих выборок. Тесты... показывают, что бэггинг может давать существенный прирост к точности. Важнейший элемент — это нестабильность метода прогнозирования. Если вмешательство в обучающую выборку может привести к существенным изменениям в создаваемом предикторе, тогда бэггинг может привести к повышению точности.

Вот какую процедуру предложил Брейман.

1. Случайным образом создавать подвыборку из двух строк данных (то есть «создавать реплики обучающей выборки с помощью бутстрэпа»).
2. Обучать модель на этой подвыборке.
3. Сохранять модель, а затем возвращаться к шагу 1 несколько раз.
4. В итоге получится несколько обученных моделей, которые нужно использовать для прогнозирования, а затем взять среднее значение прогнозов каждой из этих моделей.

Этот процесс известен как *бэггинг*. Он основан на глубоком и важном открытии: несмотря на то что каждая из моделей, обученная на подвыборке данных, будет допускать больше ошибок, чем модель, обученная на всем датасете, эти ошибки не будут связаны между собой. Разные модели будут допускать разные ошибки. Следовательно, среднее для этих ошибок будет равно нулю. Поэтому если мы возьмем среднее прогнозов всех моделей, то должны получить прогноз, тем больше приближающийся к верному ответу, чем больше моделей мы используем. Это невероятный результат — он означает, что мы можем повышать точность практически любого алгоритма машинного обучения, обучая его несколько раз, каждый раз на разной случайной подвыборке данных, и усредняя полученные прогнозы.

В 2001 году Брейман продемонстрировал, что этот подход к созданию моделей, будучи примененным к алгоритмам построения деревьев решений, оказывается

особенно эффективным. При этом профессор не ограничился случайным выбором строк для обучения каждой модели — при выборе каждого ветвления в каждом дереве решений он также стал делать случайный выбор из подмножества столбцов. Этот метод он назвал *случайным лесом*. На сегодня он является, вероятно, наиболее широко используемым и важным в практическом отношении методом машинного обучения.

По сути, случайный лес — это модель, усредняющая прогнозы большого числа деревьев решений, которые генерируются путем случайного изменения различных параметров, определяющих, какие данные используются для обучения дерева, и других параметров дерева. Бэггинг — это особый подход *ансамблирования*, или совмещения результатов нескольких моделей. Чтобы увидеть его в деле, перейдем к созданию собственного случайного леса.

Создание случайного леса

Случайный лес можно создать так же, как мы создавали дерево решений, только теперь нужно будет дополнительno определить параметры, указывающие количество деревьев в лесу, способ создания подвыборки элементов (строк) и подвыборки полей (столбцов).

В нижеприведенном определении функции параметр `n_estimators` задает нужное нам количество деревьев, `max_samples` указывает, сколько строк отбирать для обучения каждого дерева, а `max_features` определяет количество отбираемых на каждой точке разделения столбцов (где 0.5 означает «брать половину от общего числа столбцов»). С помощью параметра `min_samples_leaf`, который мы использовали в предыдущем разделе, можно также указать, когда нужно прекратить разделение узлов дерева, ограничив тем самым его максимальную высоту. В конце мы передаем `n_jobs=-1`, давая `fastai` команду использовать все CPU для построения деревьев в параллельном режиме. Создав для всего этого столь небольшую функцию, мы можем быстрее опробовать разные вариации на протяжении оставшейся части этой главы:

```
def rf(xs, y, n_estimators=40, max_samples=200_000,
       max_features=0.5, min_samples_leaf=5, **kwargs):
    return RandomForestRegressor(n_jobs=-1, n_estimators=n_estimators,
                                max_samples=max_samples, max_features=max_features,
                                min_samples_leaf=min_samples_leaf, oob_score=True).fit(xs, y)

m = rf(xs, y);
```

Итоговая RMSE для контрольной выборки существенно улучшилась по сравнению с последним результатом, достигнутым `DecisionTreeRegressor`, который из всех доступных данных создал лишь одно дерево:

```
m_rmse(m, xs, y), m_rmse(m, valid_xs, valid_y)
```

```
(0.170896, 0.233502)
```

Одно из важнейших свойств случайных лесов в том, что они не особо чувствительны к выбору гиперпараметров, таких как `max_features`. Вы можете установить для `n_estimators` настолько высокое значение, насколько готовы ждать завершения процесса обучения, — чем больше деревьев у вас будет, тем более точной получится модель. `max_samples` зачастую можно оставить со значением по умолчанию, если только у вас не более 200 000 точек данных, в случае чего определение для него значения в 200 000 ускорит обучение при минимальном влиянии на точность. Параметры `max_features=0.5` и `min_samples_leaf=4` работают достаточно неплохо, хотя и предустановленные в sklearn значения тоже вполне подходят.

В документации sklearn приведен пример (<https://oreil.ly/E0Och>) влияния разных установок `max_features` при увеличивающемся числе деревьев. На рис. 9.7 ниже синим показан результат при использовании наименьшего количества признаков, а зеленым, наоборот, наибольшего, то есть всех признаков. В результате на этом графике видно, что наиболее точные модели получаются, когда используется лишь подмножество признаков, но при большем числе деревьев.

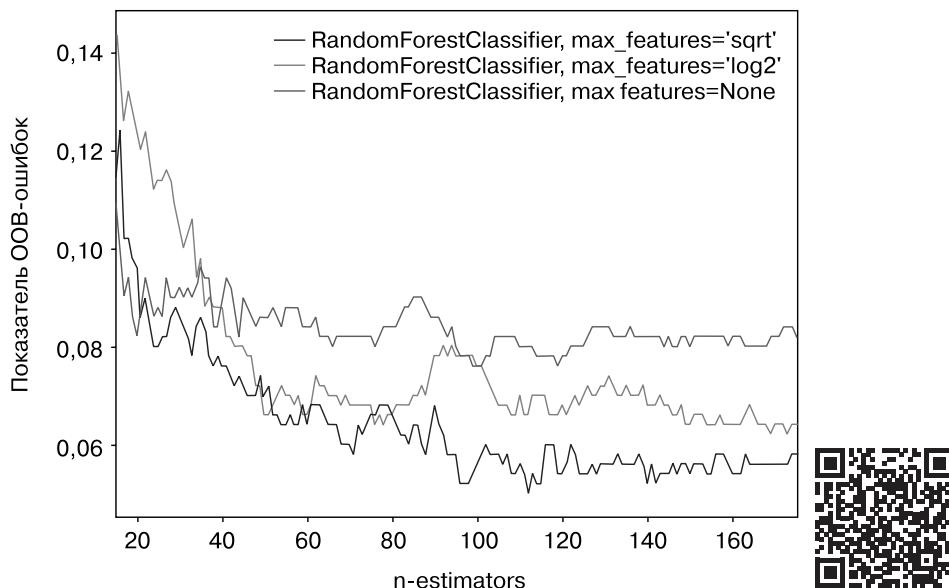


Рис. 9.7. Показатель ошибок на основе количества деревьев и максимального числа признаков (источник: <https://oreil.ly/E0Och>)

Чтобы увидеть влияние `n_estimators`, давайте получим прогнозы от каждого отдельного дерева нашего леса (они находятся в атрибуте `estimators_`):

```
preds = np.stack([t.predict(valid_xs) for t in m.estimators_])
```

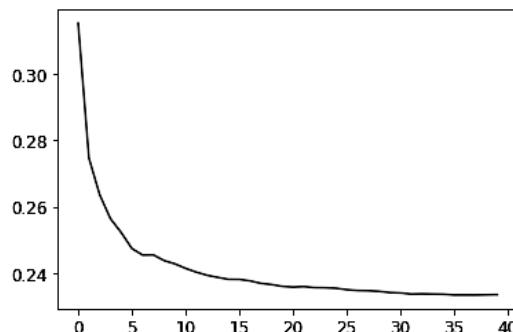
Здесь мы видим, что `preds.mean(0)` дает тот же результат, что и наш случайный лес:

```
r_mse(preds.mean(0), valid_y)
```

```
0.233502
```

Взглянем, как добавление дополнительных деревьев сказывается на RMSE. Как вы видите, спустя примерно 30 деревьев улучшение становится едва заметным.

```
plt.plot([r_mse(preds[:i+1].mean(0), valid_y) for i in range(40)]);
```



Эффективность на контрольной выборке хуже, чем на обучающей. Но в чем причина? Может, в переобучении или в том, что контрольная выборка охватывает другой временной промежуток? А может, и в том и в другом? На основе имеющейся информации мы этого понять не сможем. Тем не менее в технике случайных лесов есть очень умный прием, называемый ошибкой *out-of-bag* (вне пакета) (OOB), которая как раз поможет в данной ситуации, да и не только в ней.

Ошибка Out-of-Bag

Вспомните, что в случайному лесу каждое дерево обучается на разной подвыборке обучающих данных. OOB-ошибка — это способ измерения ошибки прогнозов в обучающем датасете путем ее вычисления для не попавших в обучение строк набора. Это позволяет увидеть, переобучается ли модель, без использования отдельной контрольной выборки.



СЛОВО АЛЕКСИСУ

Я понимаю это так. Поскольку каждое дерево обучается на разном случайнym образом выбранном наборе строк, то вычисление ошибки `out-of-bag` как бы подразумевает, что у дерева также есть собственная контрольная выборка. Эта контрольная выборка представлена строками, которые не были выбраны для обучения этого дерева.

Этот способ особенно пригождается в случаях, когда в нашем распоряжении есть лишь небольшое количество обучающих данных, так как позволяет увидеть, обобщается ли модель, не извлекая данные для создания контрольной выборки. Прогнозы ОOB доступны в атрибуте `oob_prediction_`. Обратите внимание, что мы сравниваем их с обучающими метками, поскольку вычисление делается на деревьях, использующих обучающую выборку:

```
r_mse(m.oob_prediction_, y)
```

```
0.210686
```

Здесь мы видим, что ОOB-ошибка намного меньше, чем ошибка для контрольной выборки. Это означает, что *помимо* обычной ошибки обобщаемости есть и другая причина, и несколько позже в данной главе мы рассмотрим возможные варианты таких причин.

Это был один из способов интерпретирования прогнозов модели — теперь познакомимся с другими.

Интерпретация модели

В случае с табличными данными интерпретация модели особенно важна. Что касается конкретно рассматриваемой нами модели, то интересует нас в первую очередь следующее.

- Насколько мы уверены в прогнозах при использовании конкретной строки данных?
- Какие факторы оказались наиболее важными при прогнозировании с использованием конкретной строки данных и как они повлияли на прогноз?
- Какие столбцы являются сильнейшими предикторами, а какие можно проигнорировать?
- Какие столбцы, по сути, дублируют друг друга с точки зрения прогнозирования?
- Как при изменении этих столбцов изменяются прогнозы?

Как мы позже увидим, случайные леса особенно хорошо подходят для ответа на эти вопросы. Давайте начнем с первого.

Дисперсия деревьев для уверенного прогнозирования

Мы видели, как модель усредняет прогнозы отдельного дерева для получения общего прогноза, то есть приблизительного значения. Но как нам узнать степень уверенности в приблизительном значении? Один из способов подразумевает использование не просто среднего значения, а стандартного отклонения прогнозов среди деревьев, которое сообщает нам *относительную* уверенность в прогнозах. Обычно мы стараемся быть очень осторожными при использовании результатов для строк, где деревья дают очень разные результаты (более высокие стандартные отклонения) по сравнению со случаями, где они более согласованы (низкие стандартные отклонения).

В разделе «Создание случайного леса» мы видели, как получать прогнозы для контрольной выборки, осуществляя этот процесс для каждого дерева леса с помощью генератора списков Python:

```
preds = np.stack([t.predict(valid_xs) for t in m.estimators_])  
preds.shape  
(40, 7988)
```

Теперь у нас есть прогноз для каждого дерева и каждого аукциона контрольной выборки (40 деревьев и 7988 аукционов).

С помощью этого мы можем получить стандартное отклонение прогнозов по всем деревьям для каждого аукциона:

```
preds_std = preds.std(0)
```

Вот стандартные отклонения для прогнозов по пяти первым аукционам, то есть пяти первым строкам контрольной выборки:

```
preds_std[:5]  
array([0.21529149, 0.10351274, 0.08901878, 0.28374773, 0.11977206])
```

Очевидно, что уверенность в прогнозах существенно колеблется. Для некоторых аукционов значение стандартного отклонения низкое, так как деревья согласуются. Для других же стандартное отклонение выше, поскольку в этом случае между деревьями есть разногласие. Эта информация окажется полезной в продакшне. Например, если вы использовали модель для решения, на какие лоты делать ставки на аукционе, то прогнозы с низкой уверенностью

предполагают, что вам следует более внимательно оценить лот, прежде чем ставить на него.

Важность признаков

Обычно недостаточно просто знать, что модель может делать точные прогнозы, нам также нужно знать, *как* она их делает. *Важности признаков* как раз дают нам это понимание. Получить их мы можем непосредственно из случайного леса sklearn, заглянув в атрибут `feature_importances_`. Вот простая функция, с помощью которой мы можем добавить их в DataFrame и отсортировать:

```
def rf_feat_importance(m, df):
    return pd.DataFrame({'cols':df.columns, 'imp':m.feature_importances_}
                        ).sort_values('imp', ascending=False)
```

Здесь мы видим, что несколько первых наиболее важных столбцов имеют большие показатели важности, чем остальные. При этом неудивительно, что лидируют в списке `YearMade` и `ProductSize`:

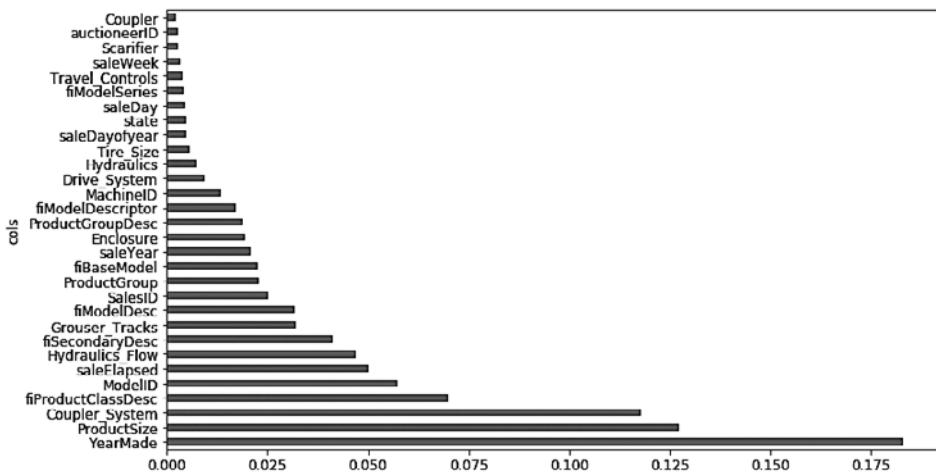
```
fi = rf_feat_importance(m, xs)
fi[:10]
```

	cols	imp
69	YearMade	0.182890
6	ProductSize	0.127268
30	Coupler_System	0.117698
7	fiProductClassDesc	0.069939
66	ModelID	0.057263
77	saleElapsed	0.050113
32	Hydraulics_Flow	0.047091
3	fiSecondaryDesc	0.041225
31	Grouser_Tracks	0.031988
1	fiModelDesc	0.031838

График важностей признаков показывает эту информацию более наглядно:

```
def plot_fi(fi):
    return fi.plot('cols', 'imp', 'barh', figsize=(12,7), legend=False)

plot_fi(fi[:30]);
```



Эти важности вычисляются достаточно простым, но в то же время изящным способом. Алгоритм важности признаков перебирает каждое дерево, а затем рекурсивно исследует каждую ветвь. В каждой ветви он выясняет, на основе какого признака делалось разделение и насколько модель в результате этого разделения улучшилась. Улучшение (взвешенное по числу строк в этой группе) прибавляется к показателю важности рассматриваемого признака. Все эти данные суммируются по всем веткам во всех деревьях, и в конце полученные показатели нормализуются так, чтобы вместе складываться в единицу.

Удаление переменных с низкой важностью

По всей вероятности, мы можем удалить переменные низкой важности и использовать только подмножество столбцов, получая при этом по-прежнему хорошие результаты. Давайте попробуем оставить только те, показатель важности признаков которых превышает 0,005:

```
to_keep = fi[fi.imp>0.005].cols
len(to_keep)
```

21

Мы можем повторно обучить модель, используя только это подмножество столбцов:

```
xs_imp = xs[to_keep]
valid_xs_imp = valid_xs[to_keep]

m = rf(xs_imp, y)
```

И вот результат:

```
m_rmse(m, xs_imp, y), m_rmse(m, valid_xs_imp, valid_y)
(0.181208, 0.232323)
```

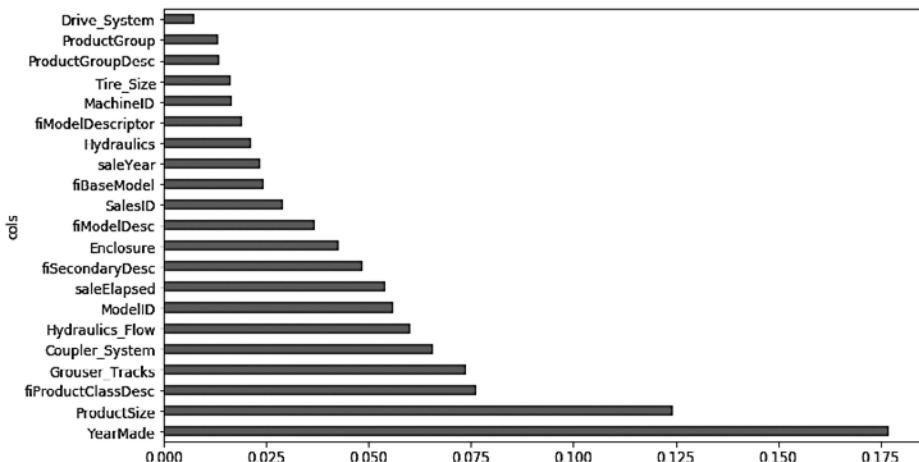
Точность практически такая же, но изучать уже требуется меньшее число столбцов:

```
len(xs.columns), len(xs_imp.columns)
(78, 21)
```

По опыту мы поняли, что обычно первым шагом улучшения модели становится ее упрощение: ведь всех 78 столбцов было бы многовато для изучения. Более того, на практике зачастую более простые и легко интерпретируемые модели легче поддаются внедрению и дальнейшему обслуживанию.

При этом график важностей признаков также стало проще трактовать. Давайте еще раз на него взглянем:

```
plot_fi(rf_feat_importance(m, xs_imp));
```

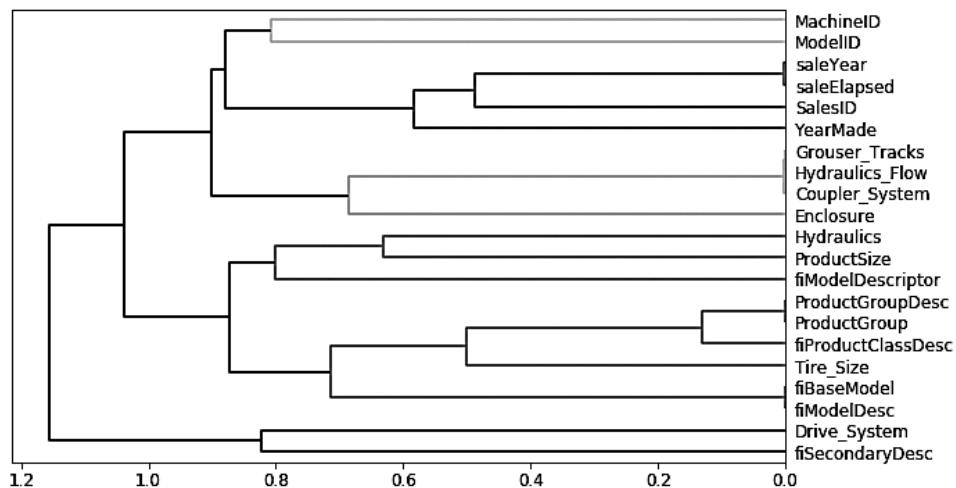


Единственное, что усложняет интерпретацию, это то, что присутствует ряд переменных с очень похожим значением: например, `ProductGroup` и `ProductGroupDesc`. Попробуем удалить все лишние признаки.

Удаление лишних признаков

Начнем с этого:

```
cluster_columns(xs_imp)
```



На этой схеме чем раньше, то есть правее от «корня» дерева слева, соединяются пары столбцов, тем более они схожи. Неудивительно, что поля `ProductGroup` и `ProductGroupDesc` соединились очень рано, что можно сказать и о `saleYear`/`saleElapsed` и `fiModelDesc`/`fiBaseModel`. Они настолько близки по смыслу, что практически выступают синонимами друг друга.



ОПРЕДЕЛЕНИЕ СХОДСТВА

Наиболее схожие пары находятся путем вычисления *ранговой корреляции*, которая подразумевает замену всех значений на их *ранг* (первый, второй, третий и т. д.) с последующим вычислением их корреляции. (Можете смело пропустить эту мелкую деталь, так как в книге она больше не встретится.)

Попробуем удалить некоторые из этих тесно связанных признаков и посмотрим, удастся ли таким образом упростить модель, не снизив ее точность. Сначала мы создаем функцию, которая быстро обучает случайный лес и возвращает показатель ОOB, используя уменьшенное количество `max_samples` и увеличенное количество `min_samples_leaf`. Показатель ОOB — это возвращаемое `sklearn` число в диапазоне между 1 для идеальной модели и 0 — для случайной. (В статистике это называется R^2 , хотя для текущего пояснения подробности не столь важны.) Нам не нужно, чтобы она была очень точной: мы хотим просто использовать ее для сравнения разных моделей в условиях удаления некоторых, вероятно излишних, столбцов:

```
def get_oob(df):
    m = RandomForestRegressor(n_estimators=40, min_samples_leaf=15,
```

```
    max_samples=50000, max_features=0.5, n_jobs=-1, oob_score=True)
m.fit(df, y)
return m.oob_score_
```

Вот наш базовый результат:

```
get_oob(xs_imp)

0.8771039618198545
```

Теперь мы пробуем по одной удалить каждую из потенциально лишних переменных:

```
{c:get_oob(xs_imp.drop(c, axis=1)) for c in (
    'saleYear', 'saleElapsed', 'ProductGroupDesc', 'ProductGroup',
    'fiModelDesc', 'fiBaseModel',
    'Hydraulics_Flow', 'Grouser_Tracks', 'Coupler_System')}

{'saleYear': 0.8759666979317242,
 'saleElapsed': 0.8728423449081594,
 'ProductGroupDesc': 0.877877012281002,
 'ProductGroup': 0.8772503407182847,
 'fiModelDesc': 0.8756415073829513,
 'fiBaseModel': 0.8765165299438019,
 'Hydraulics_Flow': 0.8778545895742573,
 'Grouser_Tracks': 0.8773718142788077,
 'Coupler_System': 0.8778016988955392}
```

Попробуем выбросить несколько переменных. Мы выбросим по одной из каждой тесно связанной пары, которые отметили ранее. Вот что получится:

```
to_drop = ['saleYear', 'ProductGroupDesc', 'fiBaseModel', 'Grouser_Tracks']
get_oob(xs_imp.drop(to_drop, axis=1))

0.8739605718147015
```

Выглядит неплохо! Однозначно не хуже, чем модель со всеми этими полями. Давайте создадим и сохраним DataFrame без этих колонок:

```
xs_final = xs_imp.drop(to_drop, axis=1)
valid_xs_final = valid_xs_imp.drop(to_drop, axis=1)

(path/'xs_final.pkl').save(xs_final)
(path/'valid_xs_final.pkl').save(valid_xs_final)
```

Позже мы сможем загрузить их обратно:

```
xs_final = (path/'xs_final.pkl').load()
valid_xs_final = (path/'valid_xs_final.pkl').load()
```

Теперь можно еще раз проверить RMSE, чтобы подтвердить отсутствие существенного изменения точности:

```
m = rf(xs_final, y)
m_rmse(m, xs_final, y), m_rmse(m, valid_xs_final, valid_y)

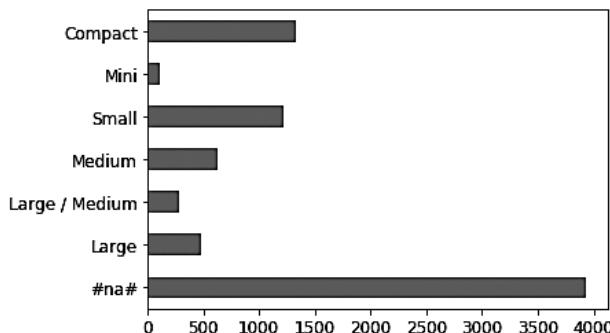
(0.183263, 0.233846)
```

Сосредоточившись на наиболее важных переменных и удалив лишние, мы существенно упростили модель. Теперь с помощью графиков частных зависимостей посмотрим, как эти переменные влияют на наши прогнозы.

Частичная зависимость

Как мы видели, двумя важнейшими предикторами являются `ProductSize` и `YearMade`. Нам нужно понять связь между этими предикторами и ценой продажи. Неплохо будет сначала проверить количество значений в каждой категории (получаемое с помощью Pandas-метода `value_counts`), чтобы увидеть, как часто каждая из этих категорий встречается:

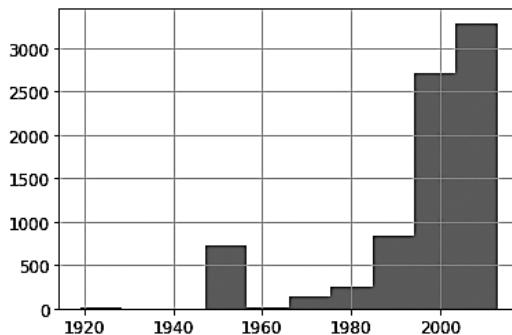
```
p = valid_xs_final['ProductSize'].value_counts(sort=False).plot.barh()
c = to.classes['ProductSize']
plt.yticks(range(len(c)), c);
```



Здесь самая крупная группа — `#na#`, которая является меткой, применяемой fastai к отсутствующим значениям.

Сделаем то же самое для `YearMade`. Поскольку это численный признак, нам потребуется нарисовать гистограмму, группирующую значения года в несколько дискретных интервалов:

```
ax = valid_xs_final['YearMade'].hist()
```



За исключением особого значения 1950, которым мы закодировали отсутствующие значения года, большая часть данных относится к периоду после 1990 года.

Теперь мы можем перейти к рассмотрению *графиков частичных зависимостей*. Задача этих графиков ответить на вопрос: если строка изменяется только рассматриваемым признаком, то как она влияет на зависимую переменную?

Например, как YearMade влияет на цену продажи при прочих равных? Для ответа на этот вопрос мы не можем просто взять среднюю цену продажи для каждого YearMade. Проблема такого подхода в том, что в зависимости от года изменяются и многие другие вещи, например, то, какие товары были проданы, сколько товаров были оснащены кондиционером, размер инфляции и т. д. Поэтому простое усреднение по всем аукционам с одинаковым YearMade также получит эффект от изменения и других связанных с YearMade полей вместе с итоговым влиянием этих изменений на стоимость.



СЛОВО АЛЕКСИСУ

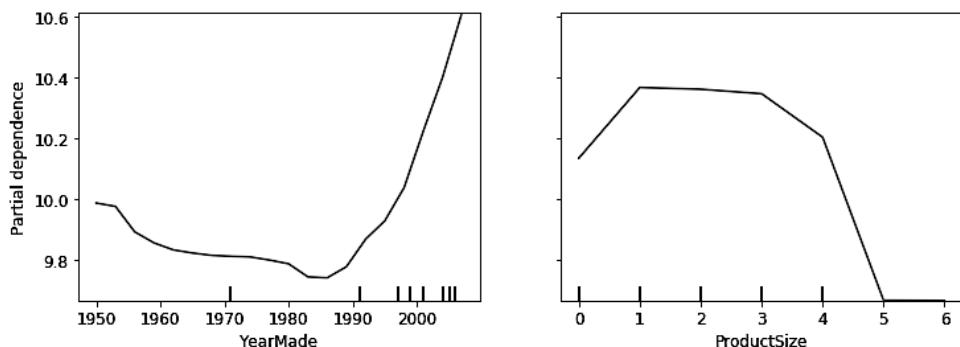
Если у вас философский склад ума, то может несколько озадачивать отслеживание разных видов гипотетических предположений, которыми мы жонглируем, чтобы произвести эти расчеты. Во-первых, нужно учесть тот факт, что каждый прогноз является гипотетическим, потому что мы учитываем не эмпирические данные. Во-вторых, суть в том, что нас не просто интересует вопрос о том, как изменится цена продажи с изменением YearMade и всего остального наряду с ним. Вместо этого мы конкретно спрашиваем, как цена продажи изменится в гипотетическом мире, где изменился только YearMade. Возможность ставить такие вопросы весьма впечатляет. Если вас интересует более глубокий анализ этих утонченных формальностей, то рекомендую к прочтению недавнюю книгу Джуды Перла (Judea Pearl) и Дэни Маккензи (Dana Mackenzie) о причинно-следственных связях: *The Book of Why*, выпущенную Basic Books.

Вместо этого мы заменим каждое значение в столбце YearMade на 1950, а затем вычислим прогнозируемую цену продажи для каждого аукциона и полу-

чим среднее значение по всем аукционам. Затем мы сделаем то же самое для 1951, 1952 и т. д. до последнего рассматриваемого 2011 года. Это избавит нас от эффекта использования одного только `YearMade` (даже если для этого используется усреднение по воображаемым записям, в которых мы присваиваем такое значение `YearMade`, которое может никогда и не существовать вместе с другими значениями).

Получив эти средние значения, мы можем отразить каждый год на оси x, а каждый прогноз — на оси y. В итоге у нас получится график частной зависимости. Взглянем на него:

```
from sklearn.inspection import plot_partial_dependence
fig,ax = plt.subplots(figsize=(12, 4))
plot_partial_dependence(m, valid_xs_final, ['YearMade', 'ProductSize'],
                        grid_resolution=20, ax=ax);
```



Глядя в первую очередь на график `YearMade`, особенно на часть, охватывающую годы после 1990 (поскольку, как мы заметили, именно здесь сосредоточена большая часть данных), мы видим практически линейную связь между годом и ценой. Помните, что зависимая переменная представлена значением после вычисления логарифма, то есть в реальности наблюдается экспоненциальное увеличение цены. Именно этого и можно было ожидать: обычно амортизация техники с течением времени признается как мультипликативный фактор, поэтому для заданной даты продажи изменение года выпуска должно показывать экспоненциальную связь с ценой продажи.

График для `ProductSize` несколько настораживает. Он показывает, что последняя группа, которая, как мы видели, предназначена для отсутствующих значений, имеет наименьшую цену. Чтобы использовать это открытие на практике, нам нужно выяснить, *почему* значения отсутствуют столь часто и что это *означает*. Отсутствующие значения иногда могут служить полезными предикторами,

в зависимости от того, что вызывает их отсутствие. При этом иногда они могут указывать на *утечку данных*.

Утечка данных

В своей работе *Leakage in Data Mining: Formulation, Detection and Avoidance* (<https://oreil.ly/XwvYf>) («Утечки при анализе данных: формулировка, обнаружение и предотвращение») Шахар Кауфман (Shachar Kaufman) и др. описывают утечку следующим образом:

Введение такой информации о цели задачи по добыче данных, которая в процессе этой добычи не должна быть доступна. Простым примером утечки будет модель, использующая в качестве входных данных саму цель, делая, таким образом, заключение, что, например, «в дождливые дни идет дождь». На практике введение недопустимой информации происходит ненамеренно и обуславливается процессом сбора, агрегирования и подготовки данных.

Авторы приводят пример:

Реальный проект бизнес-аналитики в IBM, в котором на основе ключевых слов и прочих находящихся на их сайтах факторов определялись потенциальные покупатели определенных продуктов. В этом проекте произошла утечка, поскольку использованное для обучения содержимое сайта отбиралось в тот момент времени, когда потенциальный покупатель уже стал покупателем и когда сайт содержал данные о купленных продуктах IBM, такие как слово «WebSphere» (например, в пресс-релизе о покупке самого продукта или о конкретной его особенности, которую использует клиент).

Утечка данных — это очень тонкая проблема, которая может принимать множество форм. В частности, ее нередко представляют именно отсутствующие значения.

Например, Джереми участвовал в соревновании Kaggle, посвященном прогнозированию того, какие исследователи получат гранты на исследования. Эта информация предоставлялась университетом и включала тысячи примеров исследовательских проектов наряду с информацией о привлеченных исследователях и данными о том, какие гранты в результате были одобрены, а какие нет. Университет рассчитывал получить возможность использовать модели, разработанные в этом соревновании, для ранжирования заявок на грант по вероятности их успеха, чтобы определить приоритетность их обработки.

Джереми задействовал для моделирования этих данных случайный лес, а затем использовал важность признаков, чтобы выявить наиболее предиктивные из них. При этом он отметил три удивительных фактора.

- Модель могла в 95 % случаев верно прогнозировать, кто получит гранты.
- Очевидно бессмысленные столбцы идентификаторов оказались наиболее важными предикторами.
- Столбцы дней недели и года тоже показали высокую предиктивность. Например, было одобрено подавляющее большинство заявок, поданных в воскресенье, помимо этого, многие одобренные заявки датировались первым января.

В отношении столбцов идентификаторов график частичной зависимости показал, что в случае отсутствия этой информации заявка практически всегда отклонялась. Как в итоге выяснилось, на практике университет заполнял большую часть этой информации только *после* одобрения заявки. Зачастую в отказных заявках эти поля просто оставались пустыми. Следовательно, эта информация не являлась доступной на момент регистрации заявки и, естественно, не могла быть доступна для предиктивной модели — это была утечка данных.

Аналогичным образом итоговая регистрация успешных заявок зачастую выполнялась автоматически одним пакетом в конце недели или года. В итоге в данных отражалась эта итоговая дата регистрации, которая, опять же, будучи предиктивной, фактически на момент получения заявки была недоступна.

Этот пример показывает наиболее практические и простые подходы к определению утечки данных, которые подразумевают построение модели и выполнение следующего.

- Проверить, не является ли точность модели *слишком высокой для реальной*.
- Выявить важные предикторы, которые на практике не имеют смысла.
- Выявить результаты графика частной зависимости, которые на практике бессмысленны.

Возвращаясь к нашему детектору медведей, это отражает совет, который мы давали в главе 2, — всегда сначала создавать модель и лишь затем производить очистку данных, а не наоборот. Эта модель поможет вам определить потенциальные проблемы с данными.

Кроме того, с ее помощью вы сможете определить факторы, влияющие на конкретные прогнозы, используя интерпретаторы деревьев.

Интерпретатор деревьев

В начале этого раздела мы говорили, что хотим получить ответы на следующие пять вопросов.

- Насколько мы уверены в прогнозах при использовании конкретной строки данных?
- Какие факторы оказались наиболее важными при прогнозировании с использованием конкретной строки данных и как они повлияли на прогноз?
- Какие столбцы являются сильнейшими предикторами, а какие можно проигнорировать?
- Какие столбцы, по сути, дублируют друг друга с точки зрения прогнозирования?
- Как при изменении этих столбцов изменяются прогнозы?

Четыре из них мы уже разобрали, и остался всего один. Чтобы ответить на него, потребуется использовать библиотеку *treeinterpreter*, а также библиотеку *waterfallcharts*, с помощью которой мы будем отображать результаты в виде графиков. Установить обе эти библиотеки можно так:

```
!pip install treeinterpreter  
!pip install waterfallcharts
```

Мы уже видели, как вычислять важности признаков среди всего случайного леса. Основная идея состоит в переборе каждой ветви каждого дерева и оценке вклада каждой найденной переменной в улучшение модели с последующим суммированием этих вкладов для каждой переменной.

Теперь мы можем сделать все то же самое, но для одной строки данных. Предположим, что ищем конкретный лот аукциона. Наша модель может спрогнозировать, что этот лот будет очень дорогим, и мы хотим узнать почему. Поэтому мы берем эту одну строку данных и передаем ее через первое дерево решений, оценивая, какое разделение в каких точках дерева производится. Для каждого разделения мы находим увеличение или уменьшение в сложении по сравнению с родительским узлом. Этот процесс мы проделываем для каждого дерева и суммируем общее изменение важности по разделяющей переменной.

В качестве примера возьмем несколько строк из контрольной выборки:

```
row = valid_xs_final.iloc[:5]
```

Затем передадим их в *treeinterpreter*:

```
prediction,bias,contributions = treeinterpreter.predict(m, row.values)
```

Здесь *prediction* — это просто прогноз, который делает случайный лес. *bias* — это прогноз, основанный на получении среднего зависимой переменной (то есть *модели*, являющейся корнем каждого дерева). *contributions* — это самая интересная часть, которая сообщает общее изменение в прогнозе, вызванное

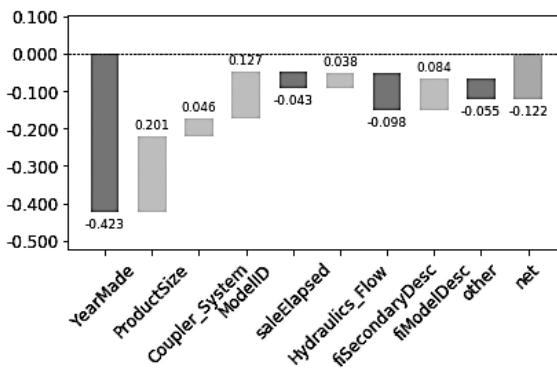
каждой независимой переменной. Следовательно, сумма `contributions` и `bias` для каждой строки должна равняться `prediction`. Взглянем на первую строку:

```
prediction[0], bias[0], contributions[0].sum()

(array([9.98234598]), 10.104309759725059, -0.12196378442186026)
```

Для наиболее ясного отображения вкладов переменных лучше использовать *диаграмму waterfall*. Она показывает, как положительные и отрицательные вклады от всех независимых переменных суммируются, создавая итоговый прогноз, который в данном случае отражен крайним правым столбцом `net`:

```
waterfall(valid_xs_final.columns, contributions[0], threshold=0.08,
          rotation_value=45, formatting='{:,.3f}');
```



Такая информация наиболее ценна в продакшене, но не в процессе создания модели. Ее можно использовать в качестве полезных сведений для пользователей вашего продукта, помогая им лучше понять внутреннюю суть формирования прогнозов.

Теперь, когда мы рассмотрели применение к данной задаче ряда классических техник машинного обучения, настал черед опробовать для нее глубокое обучение!

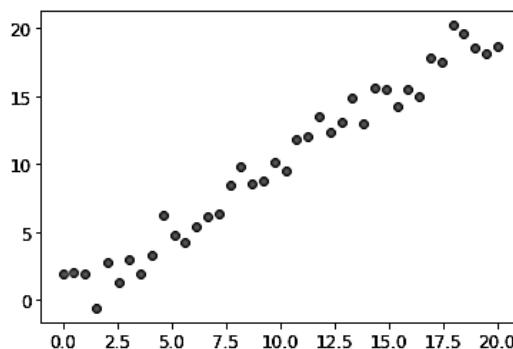
Экстраполяция и нейронные сети

Проблема случайных лесов, как и любых других алгоритмов машинного и глубокого обучения, в том, что они не всегда хорошо обобщаются на новые данные. Далее мы посмотрим, в каких ситуациях нейронные сети обобщаются лучше, но сперва разберем присущую случайнм лесам проблему экстраполяции и то, как они могут помочь в определении *не соответствующих областей данных*.

Проблема экстраполяции

Рассмотрим простую задачу получения прогнозов на основе 40 точек данных, показывающих немного шумную линейную связь:

```
x_lin = torch.linspace(0,20, steps=40)
y_lin = x_lin + torch.randn_like(x_lin)
plt.scatter(x_lin, y_lin);
```



Несмотря на то что у нас всего одна независимая переменная, sklearn ожидает матрицу независимых переменных, а не один вектор. Значит, нам нужно преобразовать имеющийся вектор в матрицу с одним столбцом. Другими словами, необходимо изменить *форму* с [40] на [40, 1]. Это можно сделать с помощью метода `unsqueeze`, который добавляет в тензор новую единичную ось в требуемом измерении:

```
xs_lin = x_lin.unsqueeze(1)
x_lin.shape, xs_lin.shape

(torch.Size([40]), torch.Size([40, 1]))
```

Более гибким подходом будет разделить массив или тензор с помощью специального значения `None`, которое вводит в это местоположение дополнительную единичную ось:

```
x_lin[:,None].shape

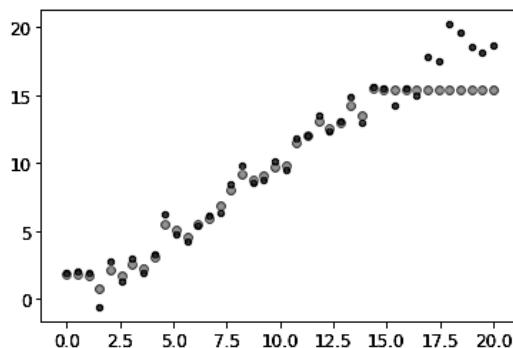
torch.Size([40, 1])
```

Теперь можно создать для этих данных модель случайного леса, для обучения которой мы задействуем только первые 30 строк:

```
m_lin = RandomForestRegressor().fit(xs_lin[:30],y_lin[:30])
```

Затем мы протестируем модель на полном датасете. Синие точки — это обучающие данные, а красные — это прогнозы.

```
plt.scatter(x_lin, y_lin, 20)
plt.scatter(x_lin, m_lin.predict(xs_lin), color='red', alpha=0.5);
```



У нас серьезная проблема! Все прогнозы вне предметной области, охватываемой нашими обучающими данными, очень низкие. Как вы думаете, почему?

Помните, что случайный лес просто усредняет прогнозы заданного числа деревьев. А простое дерево просто прогнозирует среднее значение строк в конечном узле. Из этого следует, что дерево и случайный лес никогда не смогут спрогнозировать значения вне диапазона обучающих данных. Это особенно проблематично для данных, демонстрирующих тенденцию во времени, таких как инфляция, когда вам нужно делать прогноз на будущее. В результате такие прогнозы будут систематически низкими.

Но эта проблема выходит за рамки временных переменных. Случайные леса не могут в более общем смысле экстраполироваться за рамки типов данных, которые они встречали. Именно поэтому необходимо убедиться, что наша контрольная выборка не содержит *не соответствующих областей данных*.

Поиск не соответствующих областей данных

Иногда сложно понять, отличается ли распределение данных в контрольной выборке от их распределения в обучающей, и если да, то какие столбцы отражают это отличие. Выяснить же это можно достаточно просто с помощью случайного леса.

Но в этом случае мы используем случайный лес не для прогнозирования действительных зависимых переменных. Вместо этого мы стараемся спрогнозиро-

вать, находится строка в контрольной или обучающей выборке. Чтобы увидеть этот процесс в действии, совместим контрольную и обучающую выборки, создадим зависимую переменную, отражающую, из какого набора взята каждая строка, построим на основе этих данных случайный лес и определим важность его признаков:

```
df_dom = pd.concat([xs_final, valid_xs_final])
is_valid = np.array([0]*len(xs_final) + [1]*len(valid_xs_final))

m = rf(df_dom, is_valid)
rf_feat_importance(m, df_dom)[:6]
```

	cols	imp
6	saleElapsed	0.859446
9	SalesID	0.119325
13	MachineID	0.014259
0	YearMade	0.001793
8	fiModelDesc	0.001740
11	Enclosure	0.000657

Здесь мы видим, что три столбца существенно отличаются между обучающей и контрольной выборками: `saleElapsed`, `SalesID` и `MachineID`. При этом очевидно, почему это касается `saleElapsed`: это количество дней между началом датасета и каждой строкой, поэтому столбец непосредственно кодирует дату. Разница в `SalesID` предполагает, что идентификаторы продаж аукциона могут со временем увеличиваться. `MachineID` предполагает, что подобное может происходить и с отдельными продаваемыми с этих аукционов лотами.

Получим базовый показатель RMSE модели случайного леса, а затем определим эффект от поочередного удаления этих столбцов:

```
m = rf(xs_final, y)
print('orig', m_rmse(m, valid_xs_final, valid_y))

for c in ('SalesID','saleElapsed','MachineID'):
    m = rf(xs_final.drop(c,axis=1), y)
    print(c, m_rmse(m, valid_xs_final.drop(c,axis=1), valid_y))

orig 0.232795
SalesID 0.23109
saleElapsed 0.236221
MachineID 0.233492
```

Похоже, что мы можем, не снизив точность, удалить SalesID и MachineID. Давайте это проверим:

```
time_vars = ['SalesID', 'MachineID']
xs_final_time = xs_final.drop(time_vars, axis=1)
valid_xs_time = valid_xs_final.drop(time_vars, axis=1)

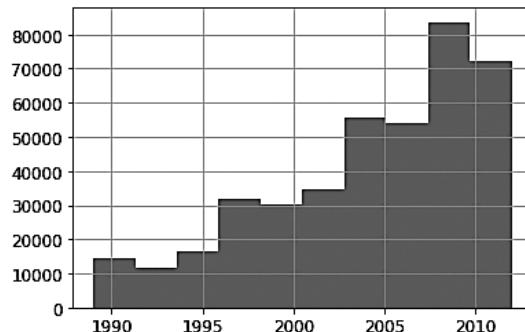
m = rf(xs_final_time, y)
m_rmse(m, valid_xs_time, valid_y)

0.231307
```

Удаление этих переменных несколько улучшило точность модели, но более важно здесь то, что с течением времени это должно сделать ее более отказоустойчивой, а также удобной для понимания и обслуживания. Мы рекомендуем для всех датасетов создавать модель, в которой зависимой переменной выступает `is_valid`, как мы и сделали здесь. Это нередко помогает раскрыть тонкие проблемы *сдвига области*, которые иначе могут остаться незамеченными.

В нашем случае может оказаться полезным просто избегать использования старых данных. Устаревшие данные нередко отражают связи, утратившие свою актуальность. Попробуем просто использовать данные за последние несколько лет:

```
xs['saleYear'].hist();
```



Вот результат обучения на этой подвыборке:

```
filt = xs['saleYear']>2004
xs_filt = xs_final_time[filt]
y_filt = y[filt]

m = rf(xs_filt, y_filt)
m_rmse(m, xs_filt, y_filt), m_rmse(m, valid_xs_time, valid_y)

(0.17768, 0.230631)
```

Немного лучше, а значит, можно сделать вывод, что не всегда стоит использовать весь датасет и иногда лучше взять только его подвыборку.

Посмотрим, насколько в нашей задаче поможет нейронная сеть.

Использование нейронной сети

Мы можем использовать тот же подход для построения модели нейронной сети. Сначала повторим шаги настройки объекта `TabularProc`:

```
df_nn = pd.read_csv(path/'TrainAndValid.csv', low_memory=False)
df_nn['ProductSize'] = df_nn['ProductSize'].astype('category')
df_nn['ProductSize'].cat.set_categories(sizes, ordered=True, inplace=True)
df_nn[dep_var] = np.log(df_nn[dep_var])
df_nn = add_datepart(df_nn, 'saledate')
```

Можно задействовать уже проделанную нами работу по удалению ненужных столбцов в случайном лесу, используя тот же набор столбцов для нейронной сети:

```
df_nn_final = df_nn[list(xs_final_time.columns) + [dep_var]]
```

Категориальные столбцы в нейронных сетях обрабатываются совсем не так, как в деревьях решений. Как мы видели в главе 8, в нейронной сети их лучше всего обрабатывать с помощью вложений. Для создания вложений fastai нужно определить, какие столбцы следует рассматривать как категориальные переменные. Для этого она сравнивает число различных уровней в переменной со значением параметра `max_card`. Если оно оказывается меньше, fastai причислит такую переменную к категориальным. Размеры вложений, превышающие 10 000, как правило, должны использоваться только тогда, когда вы проверили, что они являются лучшим способом группировки переменной, поэтому в качестве значения `max_card` мы будем использовать 9000:

```
cont_nn, cat_nn = cont_cat_split(df_nn_final, max_card=9000, dep_var=dep_var)
```

Тем не менее в этом случае есть одна переменная, которая совершенно не желает рассматриваться в качестве категориальной: `saleElapsed`. По определению категориальная переменная не может экстраполироваться вне диапазона значений, которые она видела, но нам нужно иметь возможность прогнозировать цены аукциона в будущем. Следовательно, необходимо преобразовать эту переменную в непрерывную:

```
cont_nn.append('saleElapsed')
cat_nn.remove('saleElapsed')
```

Взглянем на кардинальность каждой отобранный нами категориальной переменной:

```
df_nn_final[cat_nn].nunique()
```

```
YearMade          73
ProductSize       6
Coupler_System    2
fiProductClassDesc 74
ModelID          5281
Hydraulics_Flow   3
fiSecondaryDesc   177
fiModelDesc        5059
ProductGroup      6
Enclosure         6
fiModelDescriptor 140
Drive_System       4
Hydraulics         12
Tire_Size          17
dtype: int64
```

Тот факт, что есть две переменные, относящиеся к «модели» оборудования, имеющие очень схожую высокую кардинальность, предполагает, что они могут содержать схожую излишнюю информацию. Обратите внимание, что мы не обязательно это уловим при анализе излишних признаков, поскольку этот вывод опирается на схожие переменные, отсортированные в одинаковом порядке (то есть у них должны быть уровни с одинаковыми названиями). Наличие столбцов с 5000 уровнями означает, что в матрице вложений должно присутствовать 5000 столбцов, чего лучше избегать. Посмотрим, какое влияние окажет удаление одного из этих столбцов модели на случайный лес:

```
xs_filt2 = xs_filt.drop('fiModelDescriptor', axis=1)
valid_xs_time2 = valid_xs_time.drop('fiModelDescriptor', axis=1)
m2 = rf(xs_filt2, y_filt)
m_rmse(m, xs_filt2, y_filt), m_rmse(m2, valid_xs_time2, valid_y)

(0.176706, 0.230642)
```

Влияние оказано минимальное, так что удалим его из числа предикторов для нашей нейронной сети:

```
cat_nn.remove('fiModelDescriptor')
```

Мы можем создать объект `TabularPandas` тем же способом, что и при создании случайного леса, внеся одно важное дополнение: нормализацию. Случайный лес не нуждается в нормализации — для процедуры построения дерева важен лишь порядок значения в переменной, а не то, как они масштабируются. Но как мы видели, для нейронной сети это имеет значение. Поэтому при построении объекта `TabularPandas` мы добавляем обработчик `Normalize`:

```
procs_nn = [Categorify, FillMissing, Normalize]
to_nn = TabularPandas(df_nn_final, procs_nn, cat_nn, cont_nn,
                      splits=splits, y_names=dep_var)
```

Табличные модели и данные обычно не требуют большого объема памяти GPU, поэтому можно использовать пакеты большего размера:

```
dls = to_nn.dataloaders(1024)
```

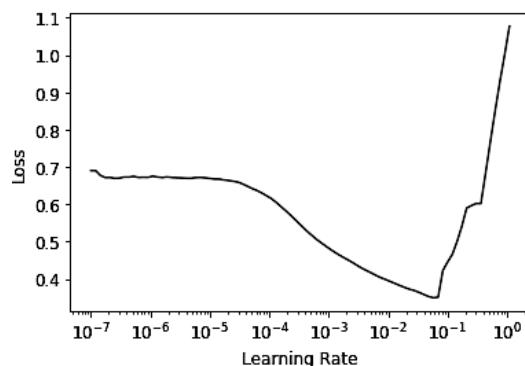
Как мы уже говорили, для моделей регрессии рекомендуется настроить `y_range`, так что найдем минимум и максимум нашей зависимой переменной:

```
y = to_nn.train.y
y.min(), y.max()
(8.465899897028686, 11.863582336583399)
```

Теперь можно создать `Learner` для построения этой табличной модели. Как обычно, мы используем обучающую функцию, соответствующую задаче, чтобы задействовать ее специфичные возможности. В качестве функции потерь мы устанавливаем `MSE`, так как именно она использовалась в соревновании.

По умолчанию для табличных данных `fastai` создает нейронную сеть с двумя скрытыми слоями, имеющими 200 и 100 активаций соответственно. Эта установка хорошо работает для небольших датасетов, но в данном случае мы работаем с достаточно большим набором данных, поэтому увеличим размеры слоев до 500 и 250:

```
from fastai.tabular.all import *
learn = tabular_learner(dls, y_range=(8,12), layers=[500,250],
                        n_out=1, loss_func=F.mse_loss)
learn.lr_find()
(0.005754399299621582, 0.0002754228771664202)
```



В использовании `fine_tune` необходимости нет, поэтому мы будем проводить обучение с `fit_one_cycle` в течение нескольких эпох и посмотрим, что в итоге получится:

```
learn.fit_one_cycle(5, 1e-2)
```

epoch	train_loss	valid_loss	time
0	0.069705	0.062389	00:11
1	0.056253	0.058489	00:11
2	0.048385	0.052256	00:11
3	0.043400	0.050743	00:11
4	0.040358	0.050986	00:11

Текущий результат можно сравнить с результатом случайного леса, полученным ранее, используя функцию `r_mse`:

```
preds,targs = learn.get_preds()
r_mse(preds,targs)
```

0.2258

Улучшение весьма ощутимо (хотя обучение заняло больше времени и потребовалась дополнительная настройка гиперпараметров).

ТАБЛИЧНЫЕ КЛАССЫ FASTAI

В fastai табличная модель — это просто модель, получающая столбцы непрерывных или категориальных данных и прогнозирующая категорию (модель классификации) или непрерывное значение (модель регрессии). Как мы видели в нейронной сети, которую использовали для коллаборативной фильтрации, категориальные независимые переменные передаются через вложение и конкатенируются, после чего конкатенируются и непрерывные переменные.

Модель, создаваемая в `tabular_learner`, является объектом класса `TabularModel`. Взгляните на исходный код `tabular_learner` теперь (помните, что это `tabular_learner??` в Jupyter). Вы увидите, что, аналогично `collab_learner`, он сначала вызывает `get_emb_sz` для вычисления подходящих размеров вложений (вы можете переопределить их с помощью параметра `emb_szs` — это словарь, содержащий любые названия столбцов, для которых вам нужно установить размеры вручную) и устанавливает несколько других значений по умолчанию. Помимо этого, он создает `TabularModel` и передает его в `TabularLearner` (обратите внимание, что `TabularLearner` идентичен `Learner`, за исключением настроенного метода `predict`).

Это означает, что, по существу, вся работа происходит в `TabularModel`, поэтому сейчас вам стоит рассмотреть ее исходный код. За исключением слоев `BatchNorm1d` и `Dropout` (о которых мы вскоре узнаем подробнее), теперь у вас есть знания, необходимые для понимания всего этого класса. Загляните в конец предыдущей главы, где мы рассказывали про `EmbeddingNN`. Вспомните, что он передавал `n_cont=0` в `TabularModel`. Теперь можно понять, что это означало: там было указано ноль непрерывных переменных (в `fastai` префикс `n_` означает `number of` (количество), а `cont` — это сокращение для `continuous` (непрерывных)).

Прежде чем двигаться далее, сохраним полученную модель на случай, если захотим вернуться к ней позже:

```
learn.save('nn')
```

Для лучшей обобщаемости также можно использовать несколько моделей с последующим усреднением их прогнозов — технику, которая упоминалась ранее и называется *ансамблированием*.

Ансамблирование

Вспомните, что изначально стоит за столь эффективными прогнозами случайных лесов: у каждого дерева есть свои ошибки, но эти ошибки не коррелируют друг с другом, поэтому как только в лесу появляется достаточно деревьев, среднее значение их ошибок начинает стремиться к нулю. Аналогичные рассуждения можно применить и к усреднению прогнозов моделей, обученных с использованием разных алгоритмов.

В нашем случае мы имеем две сильно отличающиеся модели, обученные с помощью совершенно разных алгоритмов: случайного леса и нейронной сети. Логично предположить, что и типы совершаемых ими ошибок тоже будут различными. Следовательно, можно ожидать, что среднее прогнозов этих моделей будет лучше, чем отдельно взятый прогноз любой из них.

Как мы видели, случайный лес сам по себе уже является ансамблем. Но мы можем включить такой случайный лес и в *другой* ансамбль: ансамбль случайного леса и нейронной сети. Несмотря на то что такое ансамблирование не повлияет на успешность самого процесса моделирования, оно определенно может внести небольшое приятное улучшение в любые созданные вами модели.

При этом нужно лишь обязательно учитывать, что модель PyTorch и модель sklearn создают данные разных типов: PyTorch дает нам тензор ранга 2 (матрицу-столбец), а NumPy дает массив ранга 1 (вектор). Метод `squeeze` удаляет из тензора все единичные оси, а `to_np` преобразует его в массив NumPy:

```
rf_preds = m.predict(valid_xs_time)
ens_preds = (to_np(preds.squeeze()) + rf_preds) / 2
```

В итоге мы получаем лучший результат, чем любая из моделей могла бы достичь самостоятельно.

```
r_mse(ens_preds, valid_y)
```

```
0.22291
```

В действительности этот результат даже лучше, чем любой показатель, представленный в таблице лидеров Kaggle. Тем не менее его нельзя сравнивать с ними

непосредственно, потому что таблица лидеров использует отдельный датасет, к которому у нас доступа нет. Kaggle не позволяет нам зарегистрироваться в этом старом соревновании, чтобы выяснить, каких показателей мы смогли бы в нем добиться. Тем не менее эти полученные нами результаты выглядят весьма вдохновляюще!

Бустинг

До сих пор в качестве подхода к ансамблированию мы использовали *бэггинг*, который подразумевает совмещение множества моделей (каждая из которых обучается на отдельной подвыборке) и их усреднение. Как мы видели, при применении этого подхода к деревьям решений получается *случайный лес*.

Есть и еще один важный подход ансамблирования: *бустинг (усиление)*. В нем мы уже не усредняем модели, а складываем их. Вот как он работает.

1. Создаем небольшую модель, которая будет недообучена.
2. Вычисляем прогнозы этой модели для обучающей выборки.
3. Вычитаем полученные прогнозы из целей, получая таким образом *остатки*, представляющие ошибку для каждой точки обучающей выборки.
4. Возвращаемся к шагу 1, но вместо исходных целей используем в их качестве остатки.
5. Продолжаем этот процесс до момента достижения критерия остановки, например максимума деревьев или начала ухудшения ошибки на контрольной выборке.

При использовании этого подхода каждое новое дерево будет стараться скорректировать ошибку всех предыдущих деревьев, взятых вместе. Так как мы непрерывно создаем новые остатки путем вычитания прогнозов каждого нового дерева из остатков предыдущего дерева, получаемые остатки становятся все меньше.

Чтобы получать прогнозы с помощью ансамбля усиленных деревьев, мы вычисляем прогнозы от каждого дерева, после чего их складываем. Этот простой подход используется во множестве моделей, которые, будучи одинаковыми, могут называться по-разному. Вам наверняка будут попадаться такие выражения, как *градиентный бустинг* (gradient boosting machines – GBM) и *деревья решений с градиентным бустингом* (gradient boosted decision trees – GBDT). Помимо этого, вы можете встречать названия специальных реализующих их библиотек. На момент написания книги наиболее популярной является XGBoost.

Обратите внимание, что, в отличие от случайных лесов, при использовании данного подхода ничто не препятствует переобучению. Использование большего

числа деревьев в случайному лесу не ведет к переобучению, потому что каждое дерево независимо от всех других. Но в усиленном ансамбле, чем больше деревьев вы создаете, тем ниже показатель ошибки обучения, и в конечном счете вы получаете переобучение на контрольной выборке.

Мы не будем углубляться в подробности того, как обучать ансамбль деревьев с градиентным бустингом, потому что данная область очень быстро меняется и любое руководство по этому поводу, скорее всего, утратит свою актуальность к моменту, когда вы это прочтете. Пока мы писали книгу, в sklearn появился класс `HistGradientBoostingRegressor`, обеспечивающий превосходную эффективность. Для этого класса есть множество настраиваемых параметров, как и для всех известных нам методов деревьев с градиентным бустингом. В отличие от случайных лесов, эти деревья чрезвычайно чувствительны к выборам гиперпараметров. Поэтому в большинстве случаев для поиска наиболее эффективных гиперпараметров используется цикл, который поочередно перебирает их возможные варианты.

Еще одна очень эффективная техника — использование в модели ML-вложений, обученных с помощью нейронной сети.

Совмещение вложений с другими методами

Аннотация к работе о вложении сущностей, о которой мы говорили в начале главы, гласит: «Вложения, полученные из обученной нейронной сети, существенно повышают качество всех протестированных методов машинного обучения, если используются как входные признаки». К этому прилагается очень интересная таблица, показанная на рис. 9.8.

Метод	MAPE	MAPE (средняя абсолютная ошибка в процентах)
KNN	0,290	0,116
Случайный лес	0,158	0,108
Деревья с усилением градиента	0,152	0,115
Нейронная сеть	0,101	0,093

Рис. 9.8. Эффекты от использования вложений нейронной сети в качестве ввода в других методах ML (таблица любезно предоставлена Ченом Го и Феликсом Бекханом)

Здесь показана средняя абсолютная ошибка в процентах (MAPE) в сравнении между четырьмя техниками моделирования. Три из них мы с вами уже видели,

четвертая является очень простым базовым методом k-ближайших соседей (KNN). Первый численный столбец содержит результаты использования этих методов на данных, предоставленных на соревновании. Второй столбец показывает, что происходит, когда вы сначала обучаете нейронную сеть с категориальными вложениями, а затем используете эти вложения в модели вместо необработанных категориальных столбцов. Здесь мы видим, что в каждом случае модели существенно улучшаются при использовании вложений вместо необработанных категорий.

Это действительно важный результат, так как он показывает, что можно добиваться от нейронной сети большей эффективности, не используя ее при выводе. Можно совместно с небольшим ансамблем деревьев решений использовать вложение, которое буквально является простым поиском по массиву.

Эти вложения даже не обязательно обучать раздельно для каждой модели или задачи организации. Вместо этого после однократного обучения набора вложений для столбца в конкретной задаче их можно сохранить в общем пространстве и повторно использовать в нескольких моделях. К тому же из личного общения с аналитиками данных крупных компаний мы поняли, что такой способ уже активно применяется.

Резюме

Мы рассмотрели два подхода к табличному моделированию: ансамбли деревьев решений и нейронные сети. Мы также затронули два вида ансамблей деревьев решений: случайные леса и машины повышения градиента. Каждый из них по-своему эффективен, но при этом также подразумевает и компромиссы.

- *Случайные леса* легче всего обучать, потому что они чрезвычайно устойчивы к выбору гиперпараметров и не требуют особой предварительной обработки. Они быстро обучаются и при наличии достаточного числа деревьев не должны допускать переобучения. Однако они могут быть несколько менее точны, особенно при необходимости экстраполяции, как, например, в прогнозировании будущих временных периодов.
- *Градиентный бустинг* в теории так же быстро обучается, как случайные леса, но на практике при этом вам придется перебрать множество гиперпараметров. В таком случае возможно переобучение, но при этом градиентный бустинг зачастую несколько более точен, чем случайные леса.
- *Нейронным сетям* нужно больше времени на обучение, и они требуют дополнительной предварительной обработки, а именно нормализации. Нормализацию также необходимо использовать и при выводе. Нейронные сети могут обеспечить отличные результаты и при этом хорошо экстраполири-

ются, но только если вы осторожно выберете гиперпараметры и проследите, чтобы не произошло переобучение.

Мы предлагаем начинать аналитику с помощью случайного леса. Так вы сформируете устойчивую основу и сможете быть уверены, что она послужит уверенной отправной точкой. После этого можно применить эту модель для выбора признаков и анализа частной зависимости, чтобы лучше разобраться в имеющихся данных.

Основываясь на полученных результатах, вы сможете попробовать уже нейронные сети и GBM. Если же они в течение разумного количества времени дадут намного лучшие показатели на контрольной выборке, то можете переключиться на них. Если в вашем случае хорошо сработают ансамбли деревьев решений, то попробуйте добавить в данные вложения для категориальных переменных и посмотрите, поможет ли это улучшить обучаемость ваших деревьев решений.

Вопросник

1. Что такое непрерывная переменная?
2. Что такое категориальная переменная?
3. Назовите два слова, используемые для возможных значений категориальной переменной.
4. Что такое dense layer?
5. Каким образом вложения сущностей снижают использование памяти и ускоряют нейронные сети?
6. Для каких видов датасетов вложения сущностей особенно полезны?
7. Каковы два основных семейства алгоритмов машинного обучения?
8. Почему некоторые категориальные столбцы требуют для своих классов особого порядка? Как это реализуется в Pandas?
9. Опишите в общем виде, что делает алгоритм дерева решений.
10. Чем отличается дата от обычной категориальной или непрерывной переменной и как можно ее предварительно обработать для использования в модели?
11. Стоит ли вам задействовать случайную контрольную выборку в соревновании по оценке стоимости бульдозеров? Если нет, то какую контрольную выборку следует использовать?
12. Что такое pickle (консервация) и для чего она используется?
13. Как вычисляются `mse`, `samples` и `values` в деревьях решений, приведенных в этой главе?

14. Что мы делаем с отклонениями перед построением дерева решений?
15. Как в дереве решений обрабатываются категориальные переменные?
16. Что такое бэггинг?
17. В чем разница между `max_samples` и `max_features` при создании случайного леса?
18. Может ли увеличение `n_estimators` до очень высокого значения привести к переобучению? Почему?
19. Почему в разделе «Создание случайного леса» после рис. 9.7 `preds.mean(0)` выдал тот же результат, что и наш случайный лес?
20. Что такое ошибка out-of-bag?
21. Перечислите причины, почему ошибка модели на контрольной выборке может быть хуже, чем ОOB-ошибка. Как можно проверить эту гипотезу?
22. Объясните, почему случайные леса отлично подходят для ответа на каждый из следующих вопросов.
 - Насколько мы уверены в прогнозах при использовании конкретной строки данных?
 - Какие факторы оказались наиболее важными при прогнозировании с использованием конкретной строки данных и как они повлияли на прогноз?
 - Какие столбцы являются сильнейшими предикторами, а какие можно проигнорировать?
 - Какие столбцы, по сути, дублируют друг друга с точки зрения прогнозирования?
23. Как при изменении этих столбцов изменяются прогнозы?
24. В чем смысл удаления несущественных переменных?
25. Какой тип графика хорошо подойдет для отображения результатов интерпретации дерева?
26. В чем заключается проблема экстраполяции?
27. Как можно понять, что в тестовой или контрольной выборке данные распределяются не так, как в обучающей?
28. Почему мы делаем `saleElapsed` непрерывной переменной, несмотря на то что в ней менее 9000 различных значений?
29. Что такое бустинг?
30. Как можно использовать вложения в связке со случайнм лесом? В чем это может помочь?
31. Почему мы не всегда можем использовать нейронную сеть для табличного моделирования?

Дополнительные задания

1. Выберите на Kaggle соревнование с использованием табличных данных (текущее или прошедшее) и попробуйте подстроить техники из этой главы для получения наилучших результатов. Сравните эти результаты с таблицей лидеров.
2. Реализуйте алгоритм дерева решений этой главы с нуля самостоятельно и опробуйте его на датасете, использованном в первом упражнении.
3. Используйте вложения из нейронной сети этой главы в случайному лесу и посмотрите, можете ли вы улучшить результаты случайногого леса, полученные нами.
4. Объясните, что делает каждая строка исходного кода `TabularModel` (за исключением слоев `BatchNormId` и `Dropout`).

ГЛАВА 10

Погружение в NLP: рекуррентные нейронные сети

В главе 1 мы видели, что глубокое обучение можно использовать для получения отличных результатов на датасетах естественного языка. Наш пример опирался на предварительно обученную языковую модель, настроенную для классификации отзывов. Тот пример подчеркивал различия между переносом обучения в NLP и в компьютерном зрении: обычно предварительно обученная NLP-модель изначально обучается на другой задаче.

Так называемая *языковая модель* – это модель, обученная угадывать следующее слово текста (прочитав предыдущие). Такой вид задачи называется *самообучением*: нам не требуется предоставлять модели метки, нужно лишь передать ей очень много текстов. В ней наложен процесс автоматического получения меток из данных, что является весьма непростой задачей: чтобы угадать следующее слово предложения, модели потребуется выработать понимание английского или иного языка. Самообучение также может использоваться и в других предметных областях. Загляните, например, в *Self-Supervised Learning and Computer Vision* (<https://oreil.ly/ECjfJ>) («Самообучение и компьютерное зрение») для ознакомления с применением этой техники в области компьютерного зрения. Самообучение, как правило, используется не для непосредственного, а для предварительного обучения модели, впоследствии применяемой в переносе обучения.



ТЕРМИН: САМООБУЧЕНИЕ

Обучение модели с использованием не внешних меток, а меток, вложенных в независимую переменную. Например, обучение модели прогнозированию следующего слова текста.

Языковая модель, которую в главе 1 мы использовали для классификации отзывов на IMDb, была предварительно обучена на Википедии. Мы добились

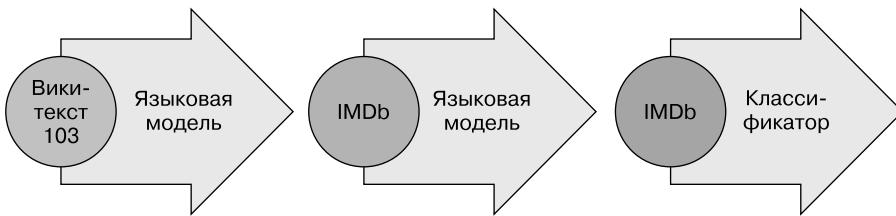
отличных результатов, тонко настроив эту модель непосредственно для классификации отзывов, но если добавить в этот процесс еще один этап, то можно добиться даже большего. Английский в Википедии несколько отличается от английского на IMDb, поэтому вместо того, чтобы переходить сразу к классификатору, мы можем тонко настроить предварительно обученную языковую модель на корпусе IMDb, а затем уже использовать полученный вариант в качестве основы для классификатора.

Даже если наша языковая модель знает только основы используемого нами в задаче языка (например, рассматриваемая модель разработана на английском), это помогает ей подстроиться к стилю корпуса текста, который мы берем в качестве цели. Это может быть более неформальный язык или более технический, где присутствуют новые слова для изучения или другие способы построения предложений. В случае с датасетом IMDb будет наблюдаться множество имен режиссеров и актеров, а также зачастую менее официальный стиль языка, чем в Википедии.

Мы уже видели, что с помощью fastai можно скачать предварительно обученную английскую языковую модель и использовать ее для получения эталонных результатов в NLP-классификации. (Мы ожидаем, что вскоре появится множество предварительно обученных на разных языках моделей; по правде говоря, они вполне могут оказаться доступными даже к моменту прочтения вами этих строк.) Итак, зачем же мы разбираем обучение языковой модели в таких подробностях?

Одна из причин, конечно же, в том, что это поможет понять основы используемых вами моделей. Но есть и еще одна очень практическая причина, заключающаяся в том, что вы получите лучшие результаты, если тонко настроите (на основе последовательности) языковую модель, прежде чем тонко настраивать модель классификации. Например, для задачи анализа настроений датасет включает 50 000 дополнительных отзывов, которые не размечены как положительные или отрицательные. А поскольку в обучающей и контрольной выборке присутствует по 25 000 размеченных отзывов, то в общей сложности получается уже 100 000. Мы можем использовать все эти отзывы для тонкой настройки языковой модели, которая предварительно обучалась только на статьях Википедии. В таком случае у нас получится языковая модель, которая будет особенно хороша в прогнозировании следующего слова в отзыве к фильму.

Такой подход называется тонкой настройкой универсальной языковой модели (Universal Language Model Fine-tuning – ULMFiT). Работа, в которой он предлагается (<https://oreil.ly/rET-C>), показала, что эта дополнительная стадия тонкой настройки языковой модели, предваряющая перенос обучения для задачи классификации, приводит к наиболее точным прогнозам. При использовании этого подхода у нас получается три стадии переноса обучения в NLP, которые отображены на рис. 10.1.

**Рис. 10.1.** Процесс ULMFiT

Далее мы рассмотрим, как применить нейронную сеть к этой задаче моделирования языка, используя принципы, с которыми мы познакомились в двух предыдущих главах. Но прежде чем продолжать чтение, приостановитесь и подумайте, какой подход к этой задаче выбрали бы *вы*.

Предварительная обработка текста

Пока не совсем понятно, как мы собираемся использовать все изученное для построения языковой модели. Предложения могут отличаться длиной, а документы могут быть достаточно объемными. Так как же нам прогнозировать следующее слово предложения с помощью нейронной сети? Давайте это выясним!

Мы уже видели, как для нейронной сети в качестве независимых переменных можно задействовать категориальные. Ниже приведен подход, примененный нами для одной категориальной переменной.

1. Создать список всех возможных уровней этой категориальной переменной (назовем его словарем).
2. Заменить каждый уровень на его индекс в словаре.
3. Создать для этого матрицу вложений, содержащую строку для каждого уровня (то есть для каждого элемента словаря).
4. Использовать эту матрицу вложений в качестве первого слоя нейронной сети. (Выделенная матрица вложений может получать в качестве ввода необработанные индексы словаря, созданные в шаге 2; она равнозначна матрице, которая получает в качестве ввода векторы индексов, закодированные унитарно, но при этом более быстра и эффективна.)

Практически то же самое можно сделать и с текстом. При этом изменится лишь последовательность. Сначала мы конкатенируем все документы датасета в одну длинную строку, которую разделяем на слова (или *токены*), получая огромный список слов. Нашей независимой переменной будет последовательность этих слов, начиная с первого и заканчивая предпоследним. Зависимой же переменной будет последовательность слов, начиная со второго и заканчивая последним.

Словарь будет состоять из смеси общих слов, уже находящихся в словаре предварительно обученной модели, и новых слов, являющихся специфичными для нашего корпуса (например, кинематографических выражений и имен актеров). Матрица вложений будет строиться соответственно: для слов, находящихся в словаре предварительно обученной модели, мы будем брать соответствующую строку ее матрицы, но для новых слов у нас ничего не будет, поэтому мы просто будем инициализировать соответствующую строку со случайным вектором.

Каждый из этих шагов создания языковой модели имеет свою характерную терминологию из мира обработки естественного языка, а также использует вспомогательные классы fastai и PyTorch. Вот эти шаги.

Токенизация

Преобразование текста в список слов (или знаков, или подстрок в зависимости от детализации модели).

Нумеризация

Составление списка всех уникальных слов словаря и преобразование каждого из них в число, соответствующее его индексу в словаре.

Создание загрузчика данных языковой модели

fastai предоставляет класс `LMDataloader`, который автоматически обрабатывает создание зависимой переменной, смещенной от независимой на один токен. Он также обрабатывает ряд важных деталей, например перемешивание обучающих данных с сохранением требуемой структуры независимой и зависимой переменных.

Создание языковой модели

Нам нужен особый вид модели, которая делает то, чего мы с вами еще не видели, а именно обрабатывает входные списки, которые могут быть произвольно большими или малыми. Это можно осуществить несколькими способами, и в данной главе мы будем использовать *рекуррентные нейронные сети* (RNN). Подробно с ними мы познакомимся в главе 12, а пока вы можете рассматривать их просто как еще один вид глубокой нейронной сети.

Теперь подробно рассмотрим каждый из шагов.

Токенизация

Когда мы говорили о «преобразовании текста в список слов», то упустили несколько деталей. Например, как обработать пунктуацию? Как поступать с такими словами, как «don't» (do not)? Рассматривать их как одно или два слова? А что

насчет слов из медицинской или химической области? Стоит ли их разделять на отдельные смысловые части? А слова, написанные через дефис? Как насчет таких языков, как немецкий и польский, в которых могут образовываться очень длинные слова, состоящие из множества частей? Под вопросом оказываются также и японский с китайским, где основы вообще не используются и нет отчетливого представления идеи слова в принципе.

Единого верного ответа на все эти вопросы нет, соответственно, возможно несколько подходов к токенизации. Рассмотрим три основных.

На основе слов

Разделение предложения на пробелы одновременно с применением характерных для языка правил с целью отделить смысловые части даже при отсутствии пробелов (например, преобразовать *don't* в *do n't*). Обычно знаки препинания тоже разделяются в отдельные токены.

На основе подслов

Разделение слова на меньшие части на основе наиболее часто встречающихся подстрок. Например, *occasion* можно токенизировать как *o c ca sion*.

На основе знаков

Разделение предложения на отдельные знаки.

Здесь мы рассмотрим токенизацию слов и подслов, а подход на основе знаков оставим вам для самостоятельной реализации в разделе вопросника.



ТЕРМИН: ТОКЕН

Элемент списка, созданного в процессе токенизации. Он может быть словом, частью слова (подсловом) или одним знаком.

Токенизация слов с помощью fastai

Fastai предлагает согласованный интерфейс для различных вариантов токенизаторов, которые заимствуются из внешних библиотек. Область токенизации активно исследуется, поэтому постоянно появляются новые усовершенствованные токенизаторы, так что используемые в fastai предустановки тоже со временем меняются. Тем не менее API и настройки не должны подвергаться существенным изменениям, поскольку fastai стремится поддерживать согласованный API даже в условиях изменения лежащей в основе технологии.

Проверим это на датасете IMDb, который использовали в главе 1:

```
from fastai.text.all import *
path = untar_data(URLs.IMDB)
```

Для использования токенизатора нам потребуется изъять текстовые файлы. Так же как `get_image_files` (который мы уже неоднократно использовали) получает все файлы изображений из указанного пути, так и `get_text_files` получает все текстовые файлы. При этом можно дополнительно передать параметр `folders`, чтобы ограничить поиск до конкретного списка подкаталогов:

```
files = get_text_files(path, folders = ['train', 'test', 'unsup'])
```

Вот отзыв, который мы будем токенизировать (в целях экономии места приводится только его начало):

```
txt = files[0].open().read(); txt[:75]
'This movie, which I just discovered at the video store, has apparently sit '
```

По мере написания этой книги предустановленный токенизатор английских слов в fastai использует библиотеку `spaCy`. В ней заложен сложный механизм, включающий особые правила для URL, отдельных английских слов и многое другого. Тем не менее вместо непосредственного использования `SpacyTokenizer` мы задействуем `WordTokenizer`, поскольку он будет всегда указывать на текущий предустановленный токенизатор слов fastai (которым может не обязательно быть `spaCy`; это будет зависеть от момента прочтения вами данного материала).

Давайте его опробуем. Для вывода результатов мы используем fastai-функцию `coll_repr(collection, n)`. Она отображает первые `n` элементов `collection` вместе с ее полным размером — это настройки `L` по умолчанию. Обратите внимание, что токенизаторы fastai принимают для токенизации коллекцию документов, поэтому нам нужно обернуть `txt` в список:

```
spacy = WordTokenizer()
toks = first(spacy([txt]))
print(coll_repr(toks, 30))
(#201) ['This', 'movie', ',', 'which', 'I', 'just', 'discovered', 'at', 'the', 'video', 's
tore', ',', 'has', 'apparently', 'sit', 'around', 'for', 'a', 'couple', 'of', 'years', 'wi
thout', 'a', 'distributor', '.', 'It', "'s", 'easy', 'to', 'see'...]
```

Как вы видите, `spaCy`, по большому счету, просто разделила слова и знаки препинания. Но она сделала и еще кое-что, а именно разделила `it's` на `it` и `'s`. В этом есть очевидный смысл, ведь это и в самом деле раздельные слова. Токенизация оказывается на удивление тонкой задачей, если учесть все мелкие детали, которые в ходе этого процесса необходимо обработать. К счастью, `spaCy` делает это за нас достаточно хорошо: например, здесь мы видим, что `«.»` отделяется, только когда выступает окончанием предложения, но не элементом акронима или числа:

```
first(spacy(['The U.S. dollar $1 is $1.00.']))
(#9) ['The', 'U.S.', 'dollar', '$', '1', 'is', '$', '1.00', '.']
```

После этого `fastai` с помощью класса `Tokenizer` привносит в процесс некоторую дополнительную функциональность:

```
tkn = Tokenizer(spacy)
print(coll_repr(tkn(txt), 31))

(#228) ['xxbos','xxmaj','this','movie','','which','i','just','discovered','at',
'the','video','store','','has','apparently','sit','around','for','a','couple','
of','years','without','a','distributor','.','xxmaj','it','''s','easy'...]
```

Обратите внимание, что теперь здесь присутствуют токены, начинающиеся со знаков «`xx`», которые не относятся к общепринятым приставкам английского языка. Это *особые токены*.

Например, первый элемент списка, `xxbos`, является особым токеном, указывающим на начало нового текста (BOS — это стандартный акроним в NLP, обозначающий *beginning of stream* — «начало потока»). Встретив этот стартовый токен, модель сможет понять, что ей нужно «забыть» предыдущий текст и сфокусироваться на последующих словах.

Эти особые токены не привносятся непосредственно из `spaCy`. Они появляются здесь, потому что `fastai` добавляет их по умолчанию в ходе применения ряда правил при обработке текста. Задача этих правил — упростить для модели распознавание важных частей предложения. В некотором смысле мы переводим исходную английскую последовательность в упрощенный токенизованный язык — язык, спроектированный для удобства обучения модели.

К примеру, согласно этим правилам, последовательность из четырех восклицательных знаков будет заменена на один восклицательный знак, сопровождаемый токеном *повторяющегося знака* и цифрой четыре. Таким образом, матрица вложений модели может кодировать информацию об основных принципах, например повторяющуюся пунктуацию, вместо того чтобы использовать отдельный токен для каждого числа повторений каждого знака препинания. Аналогичным образом слово, написанное с заглавной буквы, будет замещено особым токеном заглавной буквы, сопровождаемым этим же словом, но уже в нижнем регистре. В итоге матрице вложений требуются только версии слов в нижнем регистре, что экономит вычислительные ресурсы и память, но при этом позволяет ей также выучить принцип использования заглавных букв.

Вот некоторые из основных специальных токенов, которые вам встретятся:

- `xxbos` — указывает на начало текста (в данном случае обзора);
- `xxmaj` — указывает на то, что следующее слово начинается с заглавной буквы (поскольку мы перевели все в нижний регистр);
- `xxunk` — указывает, что это слово неизвестно.

Чтобы узнать, какие правила были использованы, можете проверить их список по умолчанию:

```
defaults.text_proc_rules  
[<function fastai.text.core.fix_html(x)>,  
<function fastai.text.core.replace_rep(t)>,  
<function fastai.text.core.replace_wrep(t)>,  
<function fastai.text.core.spec_add_spaces(t)>,  
<function fastai.text.core.rm_useless_spaces(t)>,  
<function fastai.text.core.replace_all_caps(t)>,  
<function fastai.text.core.replace_maj(t)>,  
<function fastai.text.core.lowercase(t, add_bos=True, add_eos=False)>]
```

Как обычно, вы можете заглянуть в исходный код каждого из них в блокноте, введя следующее:

```
??replace_rep
```

Вот краткая сводка по назначению каждого:

- `fix_html` — заменяет специальные HTML-знаки на их читаемую версию (в отзывах на IMDb таких достаточно много).
- `replace_rep` — заменяет любые знаки, повторяющиеся три и более раза, на специальный токен повторения (`xxrep`), сопровождаемый количеством повторений и самим знаком.
- `replace_wrep` — заменяет любое слово, повторяющееся три и более раза, на специальный токен повторения слов (`xxwrep`), сопровождаемый количеством повторений и самим словом.
- `spec_add_spaces` — добавляет пробелы вокруг / и #.
- `rm_useless_spaces` — удаляет все повторения знака пробела.
- `replace_all_caps` — переводит слово, написанное заглавными буквами, в нижний регистр и добавляет перед ним специальный токен заглавного написания (`xxcap`).
- `replace_maj` — переводит слово, написанное с заглавной буквы, в нижний регистр и добавляет перед ним специальный токен заглавной буквы (`xxmaj`).
- `lowercase` — переводит весь текст в нижний регистр и добавляет специальный токен в начале (`xxbos`) и/или в конце (`xxeos`).

Теперь посмотрим на некоторые из них в действии:

```
coll_repr(tkn('&copy;      Fast.ai www.fast.ai/INDEX'), 31)  
"(#11) ['xxbos','@','xxmaj','fast.ai','xxrep','3','w','.fast.  
ai','/','xxup','index']"
```

А далее разберем принцип работы токенизации подслов.

Токенизация подслов

В дополнение к только что рассмотренному подходу *токенизации слов* есть еще один популярный метод, называемый *токенизацией подслов*. Первый исходит из предположения, что эффективное разделение смысловых компонентов предложения производится пробелами. Тем не менее такое предположение не всегда оказывается верным. Взгляните, к примеру, на следующее предложение: 我的名字是郝杰瑞 (здесь написано «Меня зовут Джереми Ховард» на китайском). Использование токенизатора слов в данном случае нам не поможет, потому что пробелы в этом предложении отсутствуют. Языки, подобные китайскому и японскому, не используют пробелы. Более того, в них даже нет четкого определения понятию «слово» как таковому. Другие языки, например турецкий и венгерский, могут объединять множество подслов вместе без пробелов, формируя очень длинные слова, которые включают много раздельных элементов информации.

Для обработки этих случаев, как правило, лучше всего использовать токенизацию подслов, которая выполняется в два этапа.

1. Анализ корпуса документа в поиске наиболее часто встречающихся групп букв, из которых формируется словарь.
2. Токенизация этого корпуса на основе словаря *элементов подслов*.

Рассмотрим пример. Для нашего корпуса мы задействуем первые 2000 отзывов:

```
txts = L(o.open().read() for o in files[:2000])
```

Мы инстанцируем токенизатор, передавая ему размер словаря, который хотим создать, после чего нужно будет его «обучить». Это означает, что он должен прочесть наши документы и найти общие последовательности знаков, создав из них словарь. Выполняется это с помощью `setup`. Как мы вскоре увидим, `setup` – это особый метод `fastai`, вызываемый автоматически в стандартных конвейерах обработки данных. Тем не менее поскольку пока что мы все делаем вручную, то и вызывать его нам нужно самим. Вот функция, которая выполняет эти шаги для заданного размера словаря и показывает пример вывода:

```
def subword(sz):  
    sp = SubwordTokenizer(vocab_sz=sz)  
    sp.setup(txts)  
    return ' '.join(first(sp([txt]))[:40])
```

Давайте его проверим:

```
subword(1000)  
  
'_This _movie , _which _I _just _dis c over ed _at _the _video _st or e , _has_a  
p par ent ly _s it _around _for _a _couple _of _years _without _a _dis tri but  
or . _It'
```

При использовании токенизатора подслов fastai спецсимвол `_` представляет знак пробела в исходном тексте.

Если мы задействуем словарь меньшего размера, каждый токен будет выражать меньшее число знаков, и для представления предложения потребуется больше токенов:

```
subword(200)
```

```
'_T h i s _m o v i e , _w h i c h _I _j u s t _d i s c o v e r e d _a t _t h e _v i d e o _s t o r e , _h a s '
```

С другой стороны, если мы используем более крупный словарь, то наиболее распространенные английские слова сами окажутся в словаре, и нам не потребуется такого количества токенов для представления предложения:

```
subword(10000)
```

```
"_T h i s _m o v i e , _w h i c h _I _j u s t _d i s c o v e r e d _a t _t h e _v i d e o _s t o r e , _h a s _a p p a r e n t l y _s i t _a r o u n d _f o r _a _c o u p l e _o f _y e a r s _w i t h o u t _a _d i s t r i b u t o r . _I t ' s _e a s y _t o _s e e _w h y . _T h e _s t o r y _o f _t w o _f r i e n d s _l i v i n g "
```

Выбор размера словаря подслов представляет компромисс: более крупный размер означает меньшее число токенов в предложении, что ведет к ускоренному обучению, меньшим затратам памяти и меньшему количеству запоминаемых моделью состояний. Но с другой стороны, это означает увеличение матриц вложений, что потребует изучения большего объема данных.

В целом токенизация подслов предоставляет возможность легкого масштабирования между токенизацией знаков (то есть использования небольшого словаря подслов) и токенизацией слов (то есть применения большего словаря подслов) и способна обрабатывать любой человеческий язык, не требуя разработки под него специализированных алгоритмов. Она может обрабатывать даже другие «языки», например геномную последовательность или нотную запись MIDI.

По этой причине популярность данного подхода токенизации за последний год существенно выросла, и он постепенно становится самым популярным среди своих аналогов.

Как только наши тексты разделены на токены, нужно преобразовать их в числа, чем мы и займемся далее.

Нумеризация с помощью fastai

Нумеризация — это процесс отображения токенов в целые числа. Этапы в данном случае, по сути, идентичны созданию переменной `Category`, такой как зависимая переменная цифр в MNIST.

1. Создать список возможных уровней этой категориальной переменной (словаря).
2. Заменить каждый уровень на его индекс в словаре.

Рассмотрим этот процесс в действии на примере токенизированного по словам текста, который видели ранее:

```
toks = tkn(txt)
print(coll_repr(tkn(txt), 31))

(#228) ['xxbos','xxmaj','this','movie','','which','i','just','discovered','at',
'the','video','store','','has','apparently','sit','around','for','a','couple',
'of','years','without','a','distributor','.','xxmaj','it',"s","easy"]
```

Как и в случае с `SubwordTokenizer`, чтобы создать словарь, нам нужно вызвать для `Numericalize` метод `setup`. Это означает, что сначала нам необходим токенизованный корпус. Поскольку токенизация занимает некоторое время, выполняется она в `fastai` параллельно. Но для текущего выполняемого нами вручную примера мы задействуем небольшое подмножество:

```
toks200 = txts[:200].map(tkn)
toks200[0]

(#228)
['xxbos','xxmaj','this','movie','','which','i','just','discovered','at']
```

Теперь для создания словаря можно передать полученный результат в `setup`:

```
num = Numericalize()
num.setup(toks200)
coll_repr(num.vocab,20)

(#2000) ['xxunk','xxpad','xxbos','xxeos','xxfld','xxrep','xxwrep','xxup','xxmaj',
'the','..','a','and','of','to','is','in','i','it']
```

Сначала мы видим токены правил, после которых каждое слово появляется по одному разу в порядке своей частотности. Для `Numericalize` по умолчанию установлены параметры `min_freq=3` и `max_vocab=60000`. Здесь `max_vocab=60000` означает, что `fastai` заменяет все слова, выходящие за пределы наиболее распространенных 60 000, на специальный токен *неизвестного слова*, `xxunk`. Это помогает избежать создания чрезмерно большой матрицы вложений, поскольку таковая может замедлить обучение и затребовать слишком много памяти. Кроме того, это может привести к тому, что данных для обучения полезных представлений редких слов окажется недостаточно. При всем при этом последнюю из перечисленных проблем лучше всего обработать с помощью параметра `min_freq`. Его предустановленное значение `min_freq=3` указывает, что любое слово, встречающееся реже трех раз, заменяется на `xxunk`.

fastai также может нумеризовать датасет, используя словарь, предоставленный вами, путем передачи списка слов в качестве параметра `vocab`.

Создав объект `Numericalize`, мы можем использовать его как функцию:

```
nums = num(toks)[:20]; nums  
tensor([ 2, 8, 21, 28, 11, 90, 18, 59, 0, 45, 9, 351, 499, 11, 72, 533, 584,  
146, 29, 12])
```

На этот раз наши токены преобразовались в тензор целых чисел, которые может получить наша модель. Мы можем убедиться, что они отображаются обратно в исходный текст:

```
' '.join(num.vocab[o] for o in nums)  
'xxbos xxmaj this movie , which i just xxunk at the video store , has apparently  
sit around for a'
```

Теперь, когда у нас есть числа, можно поместить их в пакеты для нашей модели.

Разделение текстов на пакеты

При работе с изображениями нам приходилось преобразовывать их все в одинаковый размер с одинаковой высотой и шириной и только затем группировать в мини-пакеты, чтобы они могли эффективно разместиться в одном тензоре. В данном случае процесс будет несколько отличаться, потому что нельзя так просто взять и изменить размер текста до желаемой длины. Кроме того, для эффективного прогнозирования следующего слова необходимо, чтобы наша модель читала текст по порядку. Это означает, что каждый новый пакет должен начинаться именно с того, на чем закончился предыдущий.

Предположим, у нас есть следующий текст:

In this chapter, we will go back over the example of classifying movie reviews we studied in chapter 1 and dig deeper under the surface. First we will look at the processing steps necessary to convert text into numbers and how to customize it. By doing this, we'll have another example of the PreProcessor used in the data block API.

Then we will study how we build a language model and train it for a while.

(В этой главе мы вернемся к примеру классификации отзывов на фильмы, который изучали в главе 1, чтобы более углубленно разобрать его детали. Сначала мы рассмотрим этапы обработки, необходимые для преобразования текста в числа, и их настройку. Таким образом, мы получим еще один пример препроцессора, используемого в этом API блока данных.)

Затем мы разберем построение языковой модели и немного ее обучим.)

В процессе токенизации будут добавлены специальные токены и обработана пунктуация, в результате чего текст приобретет следующий вид:

xxbos xxmaj in this chapter , we will go back over the example of classifying movie reviews we studied in chapter 1 and dig deeper under the surface . xxmaj first we will look at the processing steps necessary to convert text into numbers and how to customize it . xxmaj by doing this , we 'll have another example of the preprocessor used in the data block xxup api . \n xxmaj then we will study how we build a language model and train it for a while.

Теперь у нас есть 90 разделенных пробелами токенов. Предположим, нас интересует размер пакета 6. Тогда этот текст нужно разбить на 6 смежных частей с длиной 15:

xxbos	xxmaj	in	this	chapter	,	we	will	go	back	over	the	example	of	classifying
movie	reviews	we	studied	in	chapter	1	and	dig	deeper	under	the	surface	.	xxmaj
first	we	will	look	at	the	processing	steps	necessary	to	convert	text	into	numbers	and
how	to	customize	it	.	xxmaj	by	doing	this	,	we	'll	have	another	example
of	the	preprocessor	used	in	the	data	block	xxup	api	.	\n	xxmaj	then	we
will	study	how	we	build	a	language	model	and	train	it	for	a	while	.

В идеальном мире мы бы могли далее передать этот пакет в модель. Но такой подход не удастся масштабировать, потому что за пределами этого игрушечного примера один пакет, содержащий все токены, вряд ли уместится в память GPU (здесь у нас всего 90 токенов, но все отзывы IMDb образуют несколько миллионов).

Поэтому нам нужно более тщательно разделить этот массив на подмассивы с фиксированной длиной последовательности. Очень важно соблюсти порядок внутри этих подмассивов и между ними, так как мы будем использовать модель, поддерживающую состояние, в котором она запоминает прочитанный текст для прогнозирования следующего за ним.

Возвращаясь к нашему предыдущему примеру с шестью мини-пакетами длиной 15 токенов, если мы выберем для последовательности длину 5, это будет означать, что сначала нам нужно передать следующий массив:

xxbos	xxmaj	in	this	chapter
movie	reviews	we	studied	in
first	we	will	look	at
how	to	customize	it	.
of	the	preprocessor	used	in
will	study	how	we	build

А затем этот:

,	we	will	go	back
chapter	1	and	dig	deeper
the	processing	steps	necessary	to
xxmaj	by	doing	this	,
the	data	block	xxup	api
a	language	model	and	train

И наконец, этот:

over	the	example	of	classifying
under	the	surface	.	xxmaj
convert	text	into	numbers	and
we	'll	have	another	example
.	\n	xxmaj	then	we
it	for	a	while	.

Вернемся к нашему датасету отзывов. Первый шаг — преобразование отдельных текстов в единый поток, конкатенировав их вместе. Как и с изображениями, лучше всего рандомизировать порядок входных данных, поэтому в начале каждой эпохи мы будем перемешивать записи, создавая новый поток (перемешивать мы будем порядок документов, а не слов в них, иначе тексты утратят свой смысл).

Затем мы разбиваем этот поток на конкретное количество мини-потоков (это будет *размер пакета*). Например, если в потоке содержится 50 000 токенов и мы устанавливаем размер пакета 10, то в итоге получаем 10 мини-потоков, содержащих по 5000 токенов. При этом важно сохранить последовательность токенов (то есть с 1 по 5000 в первом мини-потоке, затем с 5001 по 10 000 во втором и т. д.), потому что нам нужно, чтобы модель читала непрерывные строки текста (как в предыдущем примере). В процессе предварительной обработки в начало каждого текста добавляется токен `xxbos`, благодаря которому модель во время чтения потока узнает о начале новой записи.

Подведем итог. В начале каждой эпохи мы перемешиваем коллекцию документов и конкатенируем их в поток токенов. Затем мы разделяем этот поток на пакет, состоящий из последовательных мини-потоков фиксированного размера. После этого модель по очереди читает эти мини-потоки и, благодаря внутреннему

состоянию, производит одинаковую активацию, какую бы длину последовательности мы ни выбрали.

Все это выполняется библиотекой fastai за кадром, когда мы создаем `LMDataLoader`. Для этого сначала мы применяем к токенизованным текстам объект `Numericalize`:

```
nums200 = toks200.map(num)
```

А затем передаем его в `LMDataLoader`:

```
d1 = LMDataLoader(nums200)
```

Убедимся в получении ожидаемых результатов, выбрав первый пакет:

```
x,y = first(d1)
x.shape,y.shape
(torch.Size([64, 72]), torch.Size([64, 72]))
```

И посмотрим на первую строку независимой переменной, которая должна быть началом первого текста:

```
'.join(num.vocab[o] for o in x[0][:20])
'xxbos xxmaj this movie , which i just xxunk at the video store , has apparently
sit around for a'
```

Зависимая переменная будет такой же, но смещенной на один токен:

```
'.join(num.vocab[o] for o in y[0][:20])
'xxmaj this movie , which i just xxunk at the video store , has apparently sit
around for a couple'
```

На этом завершаются все шаги препроцессинга данных. Теперь мы готовы переходить к обучению нашего классификатора текста.

Обучение классификатора текста

Как мы видели в начале главы, с помощью переноса обучения эталонный классификатор текста обучается в два этапа: сначала нужно тонко настроить предварительно обученную на Википедии языковую модель на корпусе отзывов IMDb, после чего ее можно будет использовать для обучения классификатора.

Как обычно, начнем мы со сбора данных.

Создание языковой модели с помощью DataBlock

При передаче `TextBlock` в `DataBlock` `fastai` обрабатывает токенизацию и нумеризацию автоматически. Все аргументы, которые могут быть переданы в `Tokenizer` и `Numericalize`, также можно передать и в `TextBlock`. В следующей главе мы обсудим самые простые способы раздельного выполнения каждого из этих шагов для облегчения отладки, которую в ином случае вы можете осуществлять, выполняя их вручную для подмножества ваших данных, как показывалось в предыдущих разделах. При этом не стоит забывать об удобном методе `summary`, используемом в `DataBlock`, который очень полезен для устранения проблем с данными.

Ниже показано применение `TextBlock` для создания языковой модели с использованием установок `fastai` по умолчанию:

```
get_imdb = partial(get_text_files, folders=['train', 'test', 'unsup'])

dls_lm = DataBlock(
    blocks=TextBlock.from_folder(path, is_lm=True),
    get_items=get_imdb, splitter=RandomSplitter(0.1)
).dataloaders(path, path=path, bs=128, seq_len=80)
```

Единственное отличие от ранее использованных в `DataBlock` типов в том, что здесь мы не задействуем класс напрямую (то есть `TextBlock(...)`), а вместо этого вызываем *метод класса*. Метод класса — это метод Python, который, как и предполагает его название, принадлежит *классу*, а не *объекту*. (Если вы еще не знакомы с методами класса, то обязательно поищите по ним дополнительную информацию, потому что они часто используются во многих библиотеках и приложениях Python. Мы уже применяли их несколько раз на протяжении книги, но просто не обращали на это ваше внимание.) Особенность `TextBlock` заключается в том, что настройка словаря нумеризатора может занимать длительное время (так как потребуется прочесть и токенизировать каждый документ).

При этом для максимального повышения эффективности `fastai` производит некоторые оптимизации.

- Сохраняет токенизированные документы во временном каталоге, чтобы их не пришлось токенизировать повторно.
- Выполняет несколько процессов токенизации параллельно, используя возможности всех ядер вашего ПК.

При этом нам необходимо сообщить `TextBlock`, откуда брать тексты, чтобы он мог произвести начальный препроцессинг, — за что отвечает `from_folder`.

После этого `show_batch` работает в стандартном режиме:

```
dls_lm.show_batch(max_n=2)
```

	text	text_
0	xxbos xxmaj it 's awesome ! xxmaj in xxmaj story xxmaj mode , your going from punk to pro . xxmaj you have to complete goals that involve skating , driving , and walking . xxmaj you create your own skater and give it a name , and you can make it look stupid or realistic . xxmaj you are with your friend xxmaj eric throughout the game until he betrays you and gets you kicked off of the skateboard	xxmaj it 's awesome ! xxmaj in xxmaj story xxmaj mode , your going from punk to pro . xxmaj you have to complete goals that involve skating , driving , and walking . xxmaj you create your own skater and give it a name , and you can make it look stupid or realistic . xxmaj you are with your friend xxmaj eric throughout the game until he betrays you and gets you kicked off of the skateboard xxunk
1	what xxmaj i 've read , xxmaj death xxmaj bed is based on an actual dream , xxmaj george xxmaj barry , the director , successfully transferred dream to film , only a genius could accomplish such a task . \n\n xxmaj old mansions make for good quality horror , as do portraits , not sure what to make of the killer bed with its killer yellow liquid , quite a bizarre dream , indeed . xxmaj also , this	xxmaj i 've read , xxmaj death xxmaj bed is based on an actual dream , xxmaj george xxmaj barry , the director , successfully transferred dream to film , only a genius could accomplish such a task . \n\n xxmaj old mansions make for good quality horror , as do portraits , not sure what to make of the killer bed with its killer yellow liquid , quite a bizarre dream , indeed . xxmaj also , this is

Теперь, когда наши данные подготовлены, можно переходить к тонкой настройке предварительно обученной языковой модели.

Тонкая настройка языковой модели

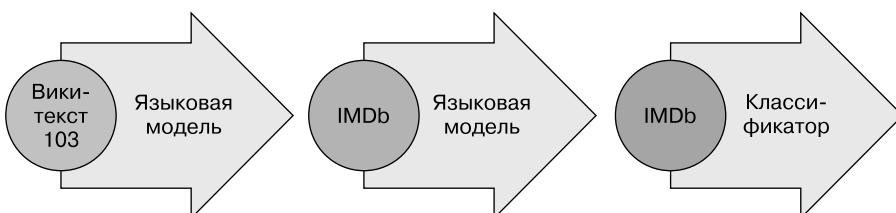
Для преобразования целочисленных индексов слов в активации, которые можно будет использовать для нейронной сети, мы задействуем вложения, так же как делали при колаборативной фильтрации и табличном моделировании. После этого мы передадим эти вложения в *рекуррентную нейронную сеть* (RNN), используя архитектуру под названием AWD-LSTM (в главе 12 мы покажем, как написать такую модель с нуля). Как уже говорилось ранее, вложения в предварительно обученной модели совмещаются со случайными вложениями, добавленными для слов, отсутствующих в словаре предварительного обучения. Это происходит автоматически внутри `language_model_learner`:

```
learn = language_model_learner(  
    dls_lm, AWD_LSTM, drop_mult=0.3,  
    metrics=[accuracy, Perplexity()]).to_fp16()
```

По умолчанию в качестве функции потерь используется перекрестная энтропия, так как, по существу, мы работаем над задачей классификации (разные категории представлены словами в словаре). В качестве метрики в данном случае используется *перплексивность*, которая часто применяется в NLP для языковых моделей и представляет экспоненту потерь (то есть `torch.exp(cross_entropy)`).

Помимо этого, мы включаем метрику точности, чтобы оценить, в каком количестве случаев наша модель угадает следующее слово, поскольку перекрестная энтропия (как мы видели) сложна для понимания и больше говорит об уверенности модели, а не о ее точности.

Вернемся к схеме процесса, приведенной в начале главы. Операция в первой стрелке была произведена за нас и предоставила нам предварительно обученную в fastai модель. Далее мы просто создали `DataLoaders` и `Learner` для второй стадии. Теперь мы готовы к тонкой настройке нашей языковой модели!



Каждая эпоха обучения занимает немало времени, поэтому в ходе этого процесса мы будем сохранять промежуточные результаты модели. Поскольку `fine_tune` за нас этого не сделает, мы используем `fit_one_cycle`. Как и `cnn_learner`, `language_model_learner` при использовании предварительно обученной модели автоматически вызывает `freeze` (это установка по умолчанию), поэтому обучаться будут только вложения (единственная часть модели, содержащая случайным образом инициализированные веса, то есть вложения для слов из словаря IMDb, отсутствующие в словаре предварительно обученной модели):

```
learn.fit_one_cycle(1, 2e-2)
```

epoch	train_loss	valid_loss	accuracy	perplexity	time
0	4.120048	3.912788	0.299565	50.038246	11:39

На обучение этой модели уходит немало времени, так что у нас есть удачная возможность обсудить ее промежуточные результаты.

Сохранение и загрузка моделей

Вы можете легко сохранять состояние модели следующим образом:

```
learn.save('1epoch')
```

В результате в `learn.path/models/` будет создан файл `1epoch.pth`. Если вы хотите загрузить свою модель на другой машине, создав такой же `Learner`, или

просто вернуться к обучению позднее, то загрузить содержимое этого файла можно так:

```
learn = learn.load('1epoch')
```

После завершения начального обучения можно вернуться к тонкой настройке, выполнив разморозку:

```
learn.unfreeze()
learn.fit_one_cycle(10, 2e-3)
```

epoch	train_loss	valid_loss	accuracy	perplexity	time
0	3.893486	3.772820	0.317104	43.502548	12:37
1	3.820479	3.717197	0.323790	41.148880	12:30
2	3.735622	3.659760	0.330321	38.851997	12:09
3	3.677086	3.624794	0.333960	37.516987	12:12
4	3.636646	3.601300	0.337017	36.645859	12:05
5	3.553636	3.584241	0.339355	36.026001	12:04
6	3.507634	3.571892	0.341353	35.583862	12:08
7	3.444101	3.565988	0.342194	35.374371	12:08
8	3.398597	3.566283	0.342647	35.384815	12:11
9	3.375563	3.568166	0.342528	35.451500	12:05

Как только и этот процесс будет завершен, мы сохраним всю модель, за исключением последнего слоя, который преобразует активации в вероятности выбора каждого токена в словаре. Модель без последнего слоя называется *энкодером*. Сохранить ее можно с помощью `save_encoder`:

```
learn.save_encoder('finetuned')
```



ТЕРМИН: ЭНКОДЕР

Модель, не включающая относящийся к конкретной задаче слой (слои). Этот термин означает почти то же, что и «тело» в отношении CNN, но энкодер чаще используется в области NLP и порождающих моделях.

На этом заканчивается вторая стадия процесса классификации текста, а именно тонкая настройка языковой модели. Теперь можно использовать ее для тонкой настройки классификатора с помощью меток тональности. Тем не менее, прежде чем переходить к настройке классификатора, предлагаем между делом по-пробовать кое-что другое — применить нашу модель для генерации случайных отзывов.

Генерация текста

Так как наша модель обучена угадывать следующее слово предложения, ее можно использовать для написания новых отзывов:

```
TEXT = "I liked this movie because"
N_WORDS = 40
N_SENTENCES = 2
preds = [learn.predict(TEXT, N_WORDS, temperature=0.75)
         for _ in range(N_SENTENCES)]

print("\n".join(preds))

i liked this movie because of its story and characters . The story line was very
strong , very good for a sci - fi film . The main character , Alucard , was very
well developed and brought the whole story
i liked this movie because i like the idea of the premise of the movie , the (
very ) convenient virus ( which , when you have to kill a few people , the "evil "
machine has to be used to protect
```

Как видите, мы добавили элемент случайности (мы выбираем случайное слово на основе возвращаемых моделью вероятностей), поэтому мы не получаем абсолютно одинаковых отзывов дважды. В нашей модели не заложено никаких знаний о правилах грамматики или структурах предложений, несмотря на то что об их построении в английском языке она выучила достаточно много: можно заметить, что она правильно использует заглавные буквы (I преобразована в i, потому что, согласно нашим правилам, слово должно содержать не менее двух знаков, чтобы допускалось его написание с заглавной буквы, поэтому нижний регистр в данном случае вполне уместен) и правильно употребляет времена. В целом, на первый взгляд отзыв выглядит вполне осмысленным, и только если вы внимательно в него вчитаетесь, то заметите определенные странности. Не плохо для модели, обученной в течение пары часов!

Но нашей задачей было обучить модель не генерировать отзывы, а классифицировать их... Так что для этого ее и используем.

Создание DataLoaders классификатора

Теперь мы переходим от тонкой настройки языковой модели к тонкой настройке классификатора. Давайте подытожим. Языковая модель прогнозирует следующее слово документа, поэтому внешние метки ей не требуются. А вот классификатор уже прогнозирует внешние метки — в случае с IMDb они представляют тональность документа.

Это означает, что структура нашего `DataBlock` для NLP-классификации будет выглядеть очень знакомой. Это почти то же самое, что мы видели на многих датасетах для классификации изображений.

```
dls_clas = DataBlock(
    blocks=(TextBlock.from_folder(path, vocab=dls_lm.vocab), CategoryBlock),
    get_y = parent_label,
    get_items=partial(get_text_files, folders=['train', 'test']),
    splitter=GrandparentSplitter(valid_name='test')
).dataloaders(path, path=path, bs=128, seq_len=72)
```

Так же как и в случае с классификацией изображений, `show_batch` показывает зависимую переменную (в нашем случае тональность) вместе с каждой независимой переменной (текст отзыва):

```
dls_clas.show_batch(max_n=3)
```

	text	category
0	xxbos i rate this movie with 3 skulls , only coz the girls knew how to scream , this could 've been a better movie , if actors were better , the twins were xxup ok , i believed they were evil , but the eldest and youngest brother , they sucked really bad , it seemed like they were reading the scripts instead of acting them . . . spoiler : if they 're vampire 's why do they freeze the blood ? vampires ca n't drink frozen blood , the sister in the movie says let 's drink her while she is alive . . . but then when they 're moving to another house , they take on a cooler they 're frozen blood . end of spoiler \n\n it was a huge waste of time , and that made me mad coz i read all the reviews of how	neg
1	xxbos i have read all of the xxmaj love xxmaj come xxmaj softly books . xxmaj knowing full well that movies can not use all aspects of the book , but generally they at least have the main point of the book . i was highly disappointed in this movie . xxmaj the only thing that they have in this movie that is in the book is that xxmaj missy 's father comes to xxunk in the book both parents come) . xxmaj that is all . xxmaj the story line was so twisted and far fetch and yes , sad , from the book , that i just could n't enjoy it . xxmaj even if i did n't read the book it was too sad . i do know that xxmaj pioneer life was rough , but the whole movie was a downer . xxmaj the rating	neg
2	xxbos xxmaj this , for lack of a better term , movie is lousy . xxmaj where do i start . . . \n\n xxmaj cinematography — xxmaj this was , perhaps , the worst xxmaj i 've seen this year . xxmaj it looked like the camera was being tossed from camera man to camera man . xxmaj maybe they only had one camera . xxmaj it gives you the sensation of being a volleyball . \n\n xxmaj there are a bunch of scenes , haphazardly , thrown in with no continuity at all . xxmaj when they did the ' split screen ' , it was absurd . xxmaj everything was squished flat , it looked ridiculous . \n\n xxmaj the color tones were way off . xxmaj these people need to learn how to balance a camera . xxmaj this ' movie ' is poorly made , and	neg

При рассмотрении определения `DataBlock` каждая его часть выглядит знакомой по предыдущим собранным нами блокам данных, но есть два важных исключения.

- В `TextBlock.from_folder` отсутствует параметр `is_lm=True`.

- Мы передаем `vocab`, который создали для тонкой настройки языковой модели.

Мы передаем языковой модели `vocab`, чтобы убедиться, что используем то же соответствие токена индексу. В противном случае обученные вложения настроенной языковой модели окажутся бессмысленными для данной модели, и толку от этой тонкой настройки не будет.

Передавая `is_lm=False` (или не передавая `is_lm` совсем, так как по умолчанию он будет `False`), мы сообщаем `TextBlock` о том, что в качестве меток используем стандартные размеченные данные, а не следующие токены. Однако при этом нам нужно урегулировать одну сложность, связанную с объединением нескольких документов в мини-пакет. Рассмотрим эту ситуацию на примере, попробовав создать мини-пакет, содержащий первые десять документов. Для начала сделаем нумерацию:

```
nums_samp = toks200[:10].map(len)
```

Теперь посмотрим, сколько токенов есть в каждом из этих десяти отзывов:

```
nums_samp.map(len)  
(#10) [228, 238, 121, 290, 196, 194, 533, 124, 581, 155]
```

Помните, что `DataLoader` в PyTorch должны объединять все элементы пакета в один тензор, имеющий фиксированную форму (то есть определенную длину каждой оси, в которую должны вписываться все элементы). Это должно прозвучать знакомо, так как подобная проблема у нас была и с изображениями. Тогда для их приведения к единому размеру мы использовали обрезку, заполнение и/или сжимание. Для документов обрезка вряд ли окажется хорошей идеей, так как наверняка приведет к удалению ключевой информации (к слову, то же касается и изображений, но обрезку мы для них используем; аугментация данных для NLP еще не была достаточно изучена, поэтому есть вероятность, что эту технику можно применить и в этой области). «Сжать» документ мы не можем технически, поэтому остается только его заполнение.

Мы расширим самые короткие тексты, сделав их все одного размера. Для этого мы применим токен заполнения, который наша модель в итоге будет игнорировать. Кроме того, для повышения производительности и устранения проблем с памятью мы объединим тексты, имеющие примерно одинаковую длину (с небольшим перемешиванием для обучающей выборки). Для этого мы будем (для обучающей выборки) сортировать документы по их длине перед каждой эпохой. В результате объединяемые в один пакет документы будут иметь примерно одинаковую длину. Мы не будем заполнять каждый пакет до одинакового размера, а вместо этого используем размер самого большого документа каждого пакета в качестве целевого.



ДИНАМИЧЕСКОЕ ИЗМЕНЕНИЕ РАЗМЕРОВ ИЗОБРАЖЕНИЙ

Нечто подобное можно проделать и с изображениями, что особенно полезно для прямоугольных экземпляров нестандартного размера. Но на момент написания книги ни одна библиотека не поддерживает эту возможность, равно как и ни одна исследовательская работа не уделяет этому вопросу внимания. Тем не менее мы планируем внедрить такую возможность в fastai, так что следите за сайтом, где мы непременно об этом сообщим.

При использовании `TextBlock` с параметром `is_lm=False` API блока данных выполняет сортировку и заполнение автоматически. (Эта проблема не распространяется на данные языковой модели, поскольку в этом случае сначала мы конкатенируем все документы, а затем разделяем их на отрезки одного размера.)

Вот теперь мы можем создать модель для классификации текстов:

```
learn = text_classifier_learner(dls_clas, AWD_LSTM, drop_mult=0.5,
                                 metrics=accuracy).to_fp16()
```

Последний шаг перед обучением классификатора — это загрузка энкодера из тонко настроенной языковой модели. Здесь вместо `load` мы используем `load_encoder`, потому что для энкодера в нашем распоряжении есть только предварительно обученные веса, а `load` в случае загрузки незавершенной модели по умолчанию выдает исключение:

```
learn = learn.load_encoder('finetuned')
```

Тонкая настройка классификатора

Последним этапом будет обучение с помощью дискриминативных скоростей и *постепенной разморозки*. В компьютерном зрении мы часто размораживаем модель в один подход, но для NLP-классификаторов более эффективным будет размораживать по нескольку слоев за раз:

```
learn.fit_one_cycle(1, 2e-2)
```

epoch	train_loss	valid_loss	accuracy	time
0	0.347427	0.184480	0.929320	00:33

За одну эпоху мы получили тот же результат, что и в главе 1, — совсем неплохо!

Мы можем передать во `freeze_to` значение 2, чтобы заморозить все параметры, кроме двух последних групп:

```
learn.freeze_to(-2)
learn.fit_one_cycle(1, slice(1e-2/(2.6**4),1e-2))
```

epoch	train_loss	valid_loss	accuracy	time
0	0.247763	0.171683	0.934640	00:37

После этого можно разморозить еще и продолжить обучение:

```
learn.freeze_to(-3)
learn.fit_one_cycle(1, slice(5e-3/(2.6**4),5e-3))
```

epoch	train_loss	valid_loss	accuracy	time
0	0.193377	0.156696	0.941200	00:45

И наконец, всю модель!

```
learn.unfreeze()
learn.fit_one_cycle(2, slice(1e-3/(2.6**4),1e-3))
```

epoch	train_loss	valid_loss	accuracy	time
0	0.172888	0.153770	0.943120	01:01
1	0.161492	0.155567	0.942640	00:57

Мы достигли точности 94,3 %, что еще три года назад считалось эталонным результатом. Если обучить другую модель прочтением всех текстов наоборот и усреднить прогнозы этих двух моделей, то можно получить даже 95,1%-ную точность, что соответствует эталонному результату, приводимому в работе ULMFiT. Этот рекорд был побит всего несколько месяцев назад путем тонкой настройки гораздо большей модели и использованием дорогостоящих техник аугментации данных (перевод предложений в другой язык и обратно с помощью другой модели).

Использование предварительно обученной модели позволило нам создать достаточно эффективную тонко настроенную языковую модель, способную либо генерировать поддельные отзывы, либо классифицировать их. Это удивительно, но нужно помнить, что данная технология может использоваться и во вред.

Дезинформация и языковые модели

Даже простые основанные на правилах алгоритмы во времена отсутствия широкодоступных языковых моделей глубокого обучения могли использоваться для создания мошеннических аккаунтов и попыток воздействия на влиятельных лиц. Джейф Кэо (Jeff Kao), ныне IT-журналист в ProPublica, проанализировал

комментарии, отправленные в Федеральную комиссию по связи США (FCC) в отношении предложения 2017 года об отмене сетевого нейтралитета. В своей статье *More than a Million Pro-Repeal Net Neutrality Comments Were Likely Faked* (<https://oreil.ly/ptq8B>) («Более миллиона комментариев в поддержку отмены сетевого нейтралитета, скорее всего, фейк») он сообщает о том, как раскрыл огромный кластер комментариев против сетевого нейтралитета, которые, по всей видимости, были сгенерированы по принципу составления стандартных писем в стиле Mad Libs. На рис. 10.2 Као выделил фейковые комментарии цветом, подчеркнув таким образом их шаблонность.

"In the matter of restoring Internet freedom. I'd like to recommend the commission to undo The Obama/Wheeler power grab to control Internet access. Americans, as opposed to Washington bureaucrats, deserve to enjoy the services they desire. The Obama/Wheeler power grab to control Internet access is a distortion of the open Internet. It ended a hands-off policy that worked exceptionally successfully for many years with bipartisan support.",
"Chairman Pai: With respect to Title 2 and net neutrality. I want to encourage the FCC to rescind Barack Obama's scheme to take over Internet access. Individual citizens, as opposed to Washington bureaucrats, should be able to select whichever services they desire. Barack Obama's scheme to take over Internet access is a corruption of net neutrality. It ended a free-market approach that performed remarkably smoothly for many years with bipartisan consensus.",
"FCC: My comments re: net neutrality regulations. I want to suggest the commission to overturn Obama's plan to take over the Internet. People like me, as opposed to so-called experts, should be free to buy whatever products they choose. Obama's plan to take over the Internet is a corruption of net neutrality. It broke a pro-consumer system that performed fabulously successfully for two decades with Republican and Democrat support.",
"Mr Pai: I'm very worried about restoring Internet freedom. I'd like to ask the FCC to overturn The Obama/Wheeler policy to regulate the Internet. Citizens, rather than the FCC, deserve to use whichever services we prefer. The Obama/Wheeler policy to regulate the Internet is a perversion of the open Internet. It disrupted a market-based approach that functioned very, very smoothly for decades with Republican and Democrat consensus.",
"FCC: In reference to net neutrality. I would like to suggest Chairman Pai to reverse Obama's scheme to control the web. Citizens, as opposed to Washington bureaucrats, should be empowered to buy whatever products they prefer. Obama's scheme to control the web is a betrayal of the open Internet. It undid a hands-off approach that functioned very, very successfully for decades with broad

Рис. 10.2. Комментарии, полученные FCC в ходе дебатов о статусе сетевого нейтралитета

П оценке Као, «менее 800 000 комментариев из 22+ миллионов можно было уверенно счесть уникальными» и «более 99 % истинных комментариев были написаны в поддержку сетевого нейтралитета».

Учитывая успехи, достигнутые в области моделирования языка после 2017 года, подобные мошеннические кампании теперь обнаружить практически невозможно. На сегодня у вас есть все необходимые инструменты для создания превосходной языковой модели: такой, которая сможет генерировать правдоподобный, соответствующий ситуации текст. Он не обязательно будет идеально точным или верным, но выглядеть будет вполне правдоподобно. Помните о том, какое значение эта технология может обрести при совмещении с кампаниями по дезинформации, подобными тем, о которых мы слышали на протяжении последних лет. Взгляните на рис. 10.3, где приведен взятый с Reddit диалог, в котором языковая модель на основе OpenAI алгоритма

GPT-2 сама с собой ведет беседу о том, стоит ли правительству США снизить расходы на оборону.

↑ [-] **unpopularopinionGPT2** [S] 38 points 15 hours ago
 ↓ You're wrong.

The defense budget is a good example of how badly the US spends money on the military. I've never seen anyone in the US that wouldn't spend money on the military. If you're going to spend that kind of money, then it absolutely makes sense to send the money to the military.

[permalink](#) [embed](#) [save](#) [report](#) [give award](#) [reply](#)

↑ [-] **unpopularopinionGPT2** [S] 9 points 15 hours ago
 ↓ Yeah, but that's already happening. There is a huge increase in the military budget, the Pentagon budget is already increasing, and the Navy is getting two frigates a year. If we just keep cutting military spending, then we're already there.

If we stop paying for the military, there would be no need for an increase in defense spending.

This is all about the money.

[permalink](#) [embed](#) [save](#) [parent](#) [report](#) [give award](#) [reply](#)

↑ [-] **unpopularopinionGPT2** [S] 4 points 15 hours ago
 ↓ I didn't mean to sound like "stop paying for the military". I'm not saying that we cannot pay the bills but I think it would make sense to cut defense spending.

The military isn't a good example of what we could do with the money we have. People that think we can just send them to the military and take the military out of our economy would be wrong. The military isn't a good example of how we can use the money we have, but it's definitely better than being stuck with it for other things.

[permalink](#) [embed](#) [save](#) [parent](#) [report](#) [give award](#) [reply](#)

Рис. 10.3. Алгоритм, беседующий сам с собой на Reddit

В этом случае объясняется, что этот алгоритм использовался для генерации диалога. Но вы только представьте, что может случиться, если злоумышленник решит запустить подобный алгоритм в социальные сети. Ведь это можно сделать медленно и осторожно, чтобы алгоритм постепенно нарабатывал последователей и доверие. При этом не потребуется большого количества ресурсов для создания буквально миллионов подобных аккаунтов. В подобной ситуации легко представить, как мы оказываемся в точке, где подавляющий процент онлайн-бесед ведется ботами, и никто об этом даже не догадывается.

Мы уже начинаем встречать примеры использования машинного обучения для генерации личностей. Например, на рис. 10.4 показан профиль Кэти Джонс (Katie Jones) на LinkedIn.

Кэти Джонс была связана на LinkedIn с несколькими членами основных аналитических центров Вашингтона. Но в реальности она не существовала. Фотография профиля была сгенерирована генеративно-состязательной сетью,

и в действительности никакая Кэти Джонс не оканчивала Центр стратегических и международных исследований (CSIS).

Многие люди предполагают или надеются, что алгоритмы встанут на нашу защиту: что мы разработаем алгоритмы классификации, которые смогут распознать автоматически генерированный контент. Но здесь проблема в том, что это выльется в бесконечную гонку вооружений, в которой все более совершенные алгоритмы классификации (или дискриминаторы) могут использоваться для создания все более совершенных алгоритмов генерации.



Рис. 10.4. Профиль Кэти Джонс на LinkedIn

Резюме

В этой главе мы ознакомились с последней областью, обработка которой заложена в fastai по умолчанию, — текстом. Здесь мы разобрали два вида моделей: языковую, которая способна генерировать тексты, и классификатор, определяющий положительность и отрицательность отзывов. Для создания эталонного классификатора мы использовали предварительно обученную языковую модель, которую тонко настроили на корпусе текста нашей задачи, после чего применили ее тело (энкодер) вместе с новой вершиной для выполнения классификации.

Прежде чем закончить эту часть книги, мы рассмотрим, как библиотека fastai способна помочь вам в сборе данных под вашу конкретную задачу.

Вопросник

1. Что такое самообучение?
2. Что такое языковая модель?
3. Почему языковая модель считается самообучающейся?
4. Для чего обычно используются самообучающиеся модели?
5. Зачем мы тонко настраиваем языковую модель?
6. Какими тремя шагами создается эталонный классификатор текста?
7. Как 50 000 неразмеченных отзывов на фильмы помогают создать более эффективный классификатор текста для датасета IMDb?
8. Каковы три этапа подготовки данных для языковой модели?
9. Что такое токенизация и зачем она нужна?
10. Назовите три подхода токенизации.
11. Что такое `xxbos`?
12. Перечислите четыре правила, применяемые `fastai` к тексту в процессе токенизации.
13. Почему повторяющиеся знаки заменяются токеном, указывающим число повторов, и он сопровождается самим повторяемым знаком?
14. Что такое нумеризация?
15. Почему могут появиться слова, замененные на токен «неизвестного слова»?
16. При размере пакета 64 первая строка тензора, представляющая первый мини-пакет, содержит первые 64 токена для датасета. Что содержит вторая строка этого тензора? Что содержит первая строка второго мини-пакета? (Будьте внимательны — студенты нередко допускают здесь ошибку. Обязательно проверьте свои ответы на сайте книги.)
17. Зачем в классификации текста применяется заполнение? Почему оно не пригождается при языковом моделировании?
18. Что содержит матрица вложений для NLP? Какова ее форма?
19. Что такое перплексивность?
20. Зачем нужно передавать словарь языковой модели в блок данных классификатора?
21. Что такое постепенная разморозка?
22. Почему генерация текста всегда будет на шаг впереди автоматического распознавания таких сгенерированных текстов?

Дополнительные задания

1. Поиските информацию о языковых моделях и дезинформации. Какие языковые модели на сегодня являются наилучшими? Рассмотрите производимые ими выводы. Убедительно ли они выглядят? Какой наихудший сценарий использования такой модели злоумышленниками для порождения конфликтов и вселения неуверенности?
2. С учетом ограничения, из-за которого модели вряд ли смогут успевать за распознаванием машинного текста, какие еще подходы могут потребоваться для обнаружения крупномасштабных кампаний дезинформации, использующих глубокое обучение?

ГЛАВА 11

Преобразование данных с помощью Mid-Level API

Мы видели, что делают с коллекцией текстов `Tokenizer` и `Numericalize` и как они используются внутри API блока данных, который обрабатывает эти преобразования за нас непосредственно с помощью `TextBlock`. Но что делать, если нам нужно применить только одно из этих преобразований для просмотра промежуточных результатов или в случае, когда мы уже токенизировали тексты? Если говорить более обобщенно, то что можно сделать, когда API блока данных недостаточно гибок для реализации нашего конкретного случая? В таких ситуациях для обработки данных нужно использовать предлагаемый `fastai mid-level API` (промежуточный API). API блока данных строится поверх этого слоя, поэтому вы сможете задействовать как его возможности, так и многие дополнительные.

Знакомство с многослойным API

Библиотека `fastai` построена на «многослойном» API. В самом верхнем слое находятся *приложения*, позволяющие нам обучать модель с помощью пяти строк кода, как мы видели в главе 1. Например, в случае создания `DataLoaders` для текстового классификатора мы использовали следующую строку:

```
from fastai.text.all import *
dls = TextDataLoaders.from_folder(unzip_data(URLs.IMDB), valid='test')
```

Фабричный метод `TextDataLoaders.from_folder` очень полезен, если данные организованы точно так же, как в датасете IMDb, но на практике это зачастую будет не так. В данном случае API блока данных предлагает большую гибкость. Как мы видели в предыдущей главе, можно получить тот же результат с помощью следующего кода:

```
path = untar_data(URLs.IMDB)
dls = DataBlock(
    blocks=(TextBlock.from_folder(path), CategoryBlock),
    get_y = parent_label,
    get_items=partial(get_text_files, folders=['train', 'test']),
    splitter=GrandparentSplitter(valid_name='test')
).dataloaders(path)
```

Но и этот способ не всегда оказывается достаточно гибким. Например, в целях отладки нам может потребоваться применить только часть преобразований, присутствующих в этом блоке данных. В другом случае нам может понадобиться создать `DataLoaders` для области применения, которая не поддерживается fastai непосредственно. В этом разделе мы подробно изучим компоненты, используемые внутри fastai для реализации API блока данных. Их понимание позволит вам задействовать всю мощь и гибкость этого промежуточного API.



ПРОМЕЖУТОЧНЫЙ API

Промежуточный API содержит не только функциональность для создания `DataLoaders`. В нем также есть система обратных вызовов, позволяющая настраивать цикл обучения любым нужным нам образом, и глобальный оптимизатор. Оба этих компонента мы рассмотрим в главе 16.

Преобразования

В процессе знакомства с токенизацией и нумеризацией в предыдущей главе мы начинали с извлечения части текстов:

```
files = get_text_files(path, folders = ['train', 'test'])
txts = L(o.open().read() for o in files[:2000])
```

Затем мы показали вам, как применить `Tokenizer` для их токенизации:

```
tok = Tokenizer.from_folder(path)
tok.setup(txts)
toks = txts.map(tok)
toks[0]
```

```
(#374) ['xxbos','x xmaj','well','','''','cube','','','(',')','1997','...']
```

А затем нумеровали их, включая автоматическое создание нового словаря для нашего корпуса:

```
num = Numericalize()
num.setup(toks)
nums = toks.map(num)
nums[0][:10]

tensor([ 2,     8,     76,     10,     23,  3112,     23,     34,  3113,     33])
```

В этих классах также есть метод `decode`. Например, инструкция `Numericalize`.`decode` возвращает нам строковые токены:

```
nums_dec = num.decode(nums[0][:10]); nums_dec  
(#10) ['xxbos', 'xxmaj', 'well', ',', "'", 'cube', "'", '(', '1997', ')']
```

`Tokenizer.decode` преобразует их обратно в одну строку (хотя эта строка может и не быть точно такой, как изначальная; все будет зависеть от того, *обратимый* ли токенизатор, но на момент написания книги токенизатор слов таковым не является):

```
tok.decode(nums_dec)  
'xxbos xxmaj well , " cube " ( 1997 )'
```

`decode` используется в `fastai`-методах `show_batch` и `show_results`, а также в ряде других методов вывода для преобразования прогнозов и мини-пакетов в понятный человеку вид.

Для каждого `tok` или `num` в предыдущих примерах мы создавали объект, называемый `setup` (который при необходимости обучает токенизатор для `tok` и создает словарь для `num`), применяли его к нашим необработанным текстам (вызывая этот объект как функцию) и в заключение декодировали результат обратно в понятное представление. Эти шаги необходимы для большинства задач препроцессинга данных, поэтому `fastai` предоставляет класс для их инкапсуляции, который называется `Transform`. И `Tokenize`, и `Numericalize` оба являются `Transform`.

Как правило, `Transform` — это объект, который действует как функция и содержит необязательный метод `setup`, инициализирующий внутреннее состояние (например, словарь внутри `num`), а также необязательный метод `decode`, преобразующий эту функцию (как мы видели в случае с `tok`, такое преобразование может оказаться не идеальным).

Хороший пример `decode` можно увидеть в преобразовании `Normalize`, которое мы встречали в главе 7: для получения возможности отобразить изображения на графике `decode` отменяет нормализацию (то есть выполняет умножение на стандартное отклонение и обратно прибавляет среднее). С другой стороны, преобразования, выполняющие аугментацию данных, не содержат метод `decode`, поскольку нам нужно отразить их эффекты на изображениях, чтобы убедиться, что аугментация работает должным образом.

Особенность поведения `Transform` в том, что они всегда применяются через кортежи. Как правило, наши данные всегда представлены кортежем (`input, target`) (Иногда в них более одного ввода и более одной цели.) При подобном применении преобразования к элементам, например, `Resize`, мы не хотим изменять размер всего кортежа сразу. Вместо этого нам нужно сделать это для ввода (если

допустимо) и цели (если допустимо) по отдельности. То же касается и преобразований пакетов, выполняющих аугментацию их данных: когда вводом выступает изображение, а целью — маска для сегментации, преобразование нужно применять (тем же способом) к вводу и цели.

Это поведение можно проанаблюдать, если передать кортеж текстов в `tok`:

```
tok((txts[0], txts[1]))
```

```
((#374) ['xxbos','xxmaj','well',' ','','cube','','','(',')','1997','...'],
 (#207) ['xxbos','xxmaj','conrad','xxmaj','hall','went','out','with','a',
 'bang'])
```

Написание собственного преобразования

Если вам нужно написать для данных собственное преобразование, то проще всего это будет сделать через функцию. Как вы увидите в примере ниже, `Transform` будет применяться только к соответствующему типу в случае передачи этого типа (в противном случае он будет применяться всегда). В следующем коде `:int` в сигнатуре функции означает, что `f` применяется только к `ints`. Именно поэтому `tfm(2.0)` возвращает `2.0`, а `tfm(2)` возвращает `3`:

```
def f(x:int): return x+1
tfm = Transform(f)
tfm(2),tfm(2.0)

(3, 2.0)
```

В данном случае `f` преобразовывается в `Transform` без методов `setup` и `decode`.

В Python есть специальный синтаксис для передачи функции (такой, как `f`) в другую функцию (или компонент, действующий как функция; в Python они называются *вызываемыми*), который называется *декоратором*. Для использования декоратора перед вызываемым компонентом добавляют `@` и размещают его перед определением функции (в Сети есть много обучающих материалов, посвященных декораторам, так что рекомендуем ознакомиться, если это понятие для вас в диковинку). Ниже приведен код, идентичный предыдущему:

```
@Transform
def f(x:int): return x+1
f(2),f(2.0)

(3, 2.0)
```

Если вам нужен `setup` или `decode`, то потребуется создать подкласс `Transform`, чтобы реализовать кодирование в `encodes`, а затем (не обязательно) настройку в `setups` и декодирование в `decodes`:

```
class NormalizeMean(Transform):
    def setups(self, items): self.mean = sum(items)/len(items)
    def encodes(self, x): return x-self.mean
    def decodes(self, x): return x+self.mean
```

Здесь `NormalizeMean` инициализирует в процессе настройки конкретное состояние (среднее всех переданных элементов), после чего преобразование заключается в вычитании этого среднего. Для декодирования мы реализуем обратное преобразование, прибавляя среднее. Вот пример `NormalizeMean` в действии:

```
tfm = NormalizeMean()
tfm.setup([1,2,3,4,5])
start = 2
y = tfm(start)
z = tfm.decode(y)
tfm.mean,y,z

(3.0, -1.0, 2.0)
```

Обратите внимание, что вызов и реализация методов выполняются по-разному:

Класс	Вызвать	Реализовать
<code>nn.Module</code> (PyTorch)	<code>()</code> (вызывать как функцию)	<code>forward</code>
<code>Transform</code>	<code>()</code>	<code>encodes</code>
<code>Transform</code>	<code>decode()</code>	<code>decodes</code>
<code>Transform</code>	<code>setup()</code>	<code>setups</code>

Поэтому, к примеру, вы никогда не вызываете напрямую `setups`, а вместо этого вызываете `setup`. Причина в том, что `setup` выполняет действия до и после вызова `setups`. Для лучшего понимания `Transform` и возможностей их использования для реализации разного поведения в зависимости от типа ввода обязательно ознакомьтесь с соответствующим разделом документации fastai.

Pipeline

Для объединения нескольких преобразований fastai предоставляет класс `Pipeline` (конвейер). В процессе его определения мы передаем ему список `Transform`, после чего он объединяет эти преобразования внутри себя. При вызове `Pipeline` для объекта он автоматически вызывает содержащиеся внутри преобразования по порядку:

```
tfms = Pipeline([tok, num])
t = tfms(txts[0]); t[:20]
```

```
tensor([ 2, 8, 76, 10, 23, 3112, 23, 34, 3113, 33,
10, 8, 4477, 22, 88, 32, 10, 27, 42, 14])
```

Затем вы можете вызвать `decode` для результатов кодирования, чтобы получить обратно читаемое представление, которое можно проанализировать.

```
tfms.decode(t)[:100]
```

```
'xxbos xxmaj well , " cube " ( 1997 ) , xxmaj vincenzo \'s first movie , was one
of the most interesti'
```

Единственное, что здесь работает не так, как в `Transform`, это настройка. Для корректной настройки `Pipeline`, содержащего набор `Transform`, в отношении определенных данных необходимо задействовать `TfmdLists`.

TfmdLists и датасеты: преобразованные коллекции

Обычно данные представляют собой набор необработанных элементов (например, имена файлов или строки в `DataFrame`), к которым вам нужно применить череду преобразований. Мы только что видели, что последовательность преобразований в `fastai` реализуется с помощью `Pipeline`. Класс, который группирует этот `Pipeline` с вашими необработанными элементами, называется `TfmdLists`.

TfmdLists

Вот короткий способ реализации преобразований из предыдущего раздела:

```
tls = TfmdLists(files, [Tokenizer.from_folder(path), Numericalize])
```

При инициализации `TfmdLists` автоматически вызывает метод `setup` каждой `Transform` по очереди, по порядку передавая ему не исходные элементы, а их преобразованные предыдущими `Transform` версии. Можно получить результат `Pipeline` для любого необработанного элемента, просто обратившись к его индексу в `TfmdLists`:

```
t = tls[0]; t[:20]
```

```
tensor ([ 2, 8, 91, 11, 22, 5793, 22, 37, 4910, 34, 11, 8, 13042, 23, 107, 30,
11, 25, 44, 14])
```

При этом `TfmdLists` умеет декодировать элементы для их просмотра:

```
tls.decode(t)[:100]
```

```
'xxbos xxmaj well , " cube " ( 1997 ) , xxmaj vincenzo \'s first movie , was one
of the most interesti'
```

В действительности у него даже есть метод `show`:

```
tls.show(t)

xxbos xxmaj well , " cube " ( 1997 ) , xxmaj vincenzo 's first movie , was one
of the most interesting and tricky ideas that xxmaj i 've ever seen when
talking about movies . xxmaj they had just one scenery , a bunch of actors
and a plot . xxmaj so , what made it so special were all the effective
direction , great dialogs and a bizarre condition that characters had to deal
like rats in a labyrinth . xxmaj his second movie , " cypher " ( 2002 ) , was
all about its story , but it was n't so good as " cube " but here are the
characters being tested like rats again .

" nothing " is something very interesting and gets xxmaj vincenzo coming back
to his ' cube days ' , locking the characters once again in a very different
space with no time once more playing with the characters like playing with
rats in an experience room . xxmaj but instead of a thriller sci - fi ( even
some of the promotional teasers and trailers erroneous seemed like that ) ,
" nothing " is a loose and light comedy that for sure can be called a modern
satire about our society and also about the intolerant world we 're living .
xxmaj once again xxmaj xxunk amaze us with a great idea into a so small kind
of thing . 2 actors and a blinding white scenario , that 's all you got most
part of time and you do n't need more than that . xxmaj while " cube " is a
claustrophobic experience and " cypher " confusing , " nothing " is
completely the opposite but at the same time also desperate .
```

xxmaj this movie proves once again that a smart idea means much more than just a millionaire budget . xxmaj of course that the movie fails sometimes , but its prime idea means a lot and offsets any flaws . xxmaj there 's nothing more to be said about this movie because everything is a brilliant surprise and a totally different experience that i had in movies since " cube " .

TfmdLists пишется во множественном числе, потому что он способен обрабатывать и обучающую, и контрольную выборки, для чего используется аргумент `splits`. Нужно лишь передать индексы элементов обучающей и контрольной выборок:

```
cut = int(len(files)*0.8)
splits = [list(range(cut)), list(range(cut,len(files)))]
tls = TfmdLists(files, [Tokenizer.from_folder(path), Numericalize],
                 splits=splits)
```

Затем можно обратиться к ним через атрибуты `train` и `valid`:

```
tls.valid[0][:20]
```

```
tensor ([ 2, 8, 20, 30, 87, 510, 1570, 12, 408, 379, 4196, 10, 8, 20, 30, 16,
13, 12216, 202, 509])
```

Если вы сами написали `Transform`, разом выполняющий всю предварительную обработку, преобразуя необработанные элементы в кортеж с вводами и целями, тогда вам нужен именно класс **TfmdLists**. Можно напрямую преобразовать его

в объект `DataLoaders` с помощью метода `dataloaders`. Именно это мы и будем делать в примере сиамских сетей несколько позже.

Как бы то ни было, обычно у вас будет два (или более) параллельных конвейера преобразований: один для обработки исходных входных элементов, а другой – для исходных целевых элементов. Например, здесь определенный нами конвейер только преобразует сырой текст во входные данные. Если же нам нужно произвести классификацию текста, то также потребуется преобразовать метки в цели.

Для этого необходимо выполнить два действия. Сначала мы берем название метки из родительского каталога, для чего используем функцию `parent_label`:

```
lbls = files.map(parent_label)
lbls
(#50000) ['pos', 'pos', 'pos', 'pos', 'pos', 'pos', 'pos', 'pos', 'pos', 'pos'...]
```

Затем нам нужна функция `Transform`, которая извлечет уникальные элементы и в процессе настройки создаст из них словарь, после чего при вызове преобразует строковые метки в целые числа. В `fastai` для этого используется `Categorize`:

```
cat = Categorize()
cat.setup(lbls)
cat.vocab, cat(lbls[0])
((#2) ['neg', 'pos'], TensorCategory(1))
```

Чтобы произвести всю настройку для списка файлов автоматически, можно, как и ранее, создать `TfmdLists`:

```
tls_y = TfmdLists(files, [parent_label, Categorize()])
tls_y[0]
TensorCategory(1)
```

Но тогда у нас получаются два раздельных объекта для ввода и целей, что нежелательно. В этом случае на выручку приходит `Datasets`.

Datasets

`Datasets` применяет два (или более) конвейера к одному необработанному объекту параллельно и создает результат в виде кортежа. Аналогично `TfmdLists` он автоматически выполняет настройку за нас, а когда мы обращаемся по индексу к `Datasets`, возвращает кортеж с результатами каждого конвейера:

```
x_tfms = [Tokenizer.from_folder(path), Numericalize]
y_tfms = [parent_label, Categorize()]
dsets = Datasets(files, [x_tfms, y_tfms])
x,y = dsets[0]
x[:20],y
```

Как и в случае с `TfmdLists`, для разделения данных между обучающей и контрольной выборками можно передать в `Datasets` метод `splits`:

```
x_tfms = [Tokenizer.from_folder(path), Numericalize]
y_tfms = [parent_label, Categorize()]
dsets = Datasets(files, [x_tfms, y_tfms], splits=splits)
x,y = dsets.valid[0]
x[:20],y

(tensor ([ 2, 8, 20, 30, 87, 510, 1570, 12, 408, 379, 4196, 10, 8, 20, 30, 16,
13, 12216, 202, 509]).
TensorCategory(0))
```

Он также может декодировать любой обработанный кортеж или показать его непосредственно:

```
t = dsets.valid[0]
dsets.decode(t)

('xxbos xxmaj this movie had horrible lighting and terrible camera movements .
 > xxmaj this movie is a jumpy horror flick with no meaning at all . xxmaj the
 > slashes are totally fake looking . xxmaj it looks like some 17 year - old
 > idiot wrote this movie and a 10 year old kid shot it . xxmaj with the worst
 > acting you can ever find . xxmaj people are tired of knives . xxmaj at least
 > move on to guns or fire . xxmaj it has almost exact lines from " when a xxmaj
 > stranger xxmaj calls " . xxmaj with gruesome killings , only crazy people
 > would enjoy this movie . xxmaj it is obvious the writer does n\'t have kids
 > or even care for them . i mean at show some mercy . xxmaj just to sum it up ,
> this movie is a " b " movie and it sucked . xxmaj just for your own sake , do
 > n\'t even think about wasting your time watching this crappy movie .',
 'neg')
```

Последний шаг — преобразование объекта `Datasets` в `DataLoaders`, что можно сделать с помощью метода `dataloaders`. При этом нам также нужно передать специальный аргумент, который решит проблему заполняемости данных (как мы видели в предыдущей главе). Это должно произойти до распределения элементов по пакетам, поэтому мы передаем его в `before_batch`:

```
dls = dsets.dataloaders(bs=64, before_batch=pad_input)
```

`dataloaders` вызывает `DataLoader` непосредственно для каждой подвыборки нашего `Datasets`. `fastai`-объект `DataLoader` расширяет одноименный класс PyTorch и отвечает за объединение элементов датасетов в пакеты. У него есть много возможных настроек, но наиболее важные среди них следующие.

- `after_item` — применяется к каждому элементу после его извлечения в датасет. Эквивалент `item_tfms` в `DataBlock`.
- `before_batch` — применяется к списку элементов до их распределения. Это идеальный момент для заполнения элементов до одного размера.

- `after_batch` — применяется ко всему пакету после его создания. Эквивалент `batch_tfms` в `DataBlock`.

В качестве заключения приведем весь код, необходимый для подготовки данных к классификации текста:

```
tfms = [[Tokenizer.from_folder(path), Numericalize], [parent_label, Categorize]]
files = get_text_files(path, folders = ['train', 'test'])
splits = GrandparentSplitter(valid_name='test')(files)
dsets = Datasets(files, tfms, splits=splits)
dls = dsets.dataloaders(dl_type=SortedDL, before_batch=pad_input)
```

Два отличия от предыдущего кода: использование `GrandparentSplitter` для разделения обучающих и контрольных данных, а также наличие аргумента `dl_type`, который сообщает `dataloaders` о необходимости применить принадлежащий `DataLoader` класс `SortedDL` вместо обычного. `SortedDL` создает пакеты, помещая в них образцы примерно одинакового размера.

Этот код делает то же, что и наш предыдущий `DataBlock`:

```
path = untar_data(URLs.IMDB)
dls = DataBlock(
    blocks=(TextBlock.from_folder(path),CategoryBlock),
    get_y = parent_label,
    get_items=partial(get_text_files, folders=['train', 'test']),
    splitter=GrandparentSplitter(valid_name='test')
).dataloaders(path)
```

Но зато теперь вы знаете, как настраивать каждую его часть.

Попрактикуем то, чему только что научились, используя этот промежуточный API для препроцессинга данных в примере компьютерного зрения.

Использование промежуточного API: SiamesePair

Сиамская модель получает два изображения и должна определить, относятся ли они к одному классу. Для этого примера мы будем снова использовать датасет домашних животных и подготовим данные для модели, которая будет прогнозировать, относятся ли животные на двух изображениях к одной породе. Здесь мы только объясним, как подготовить данные для этой модели, а ее обучением займемся в главе 15.

И опять же начнем мы с получения изображений для датасета:

```
from fastai.vision.all import *
path = untar_data(URLs.PETS)
files = get_image_files(path/"images")
```

Если бы нас не интересовал просмотр объектов, то мы могли бы непосредственно создать одно преобразование для полной предварительной обработки этого списка файлов. Но так как просмотр изображений нас интересует, нужно создать собственный тип. Когда вы вызываете метод `show` для объекта `TfmdLists` или `Datasets`, он начинает декодировать элементы, пока не достигнет типа, содержащего `show`, использовав его для показа этого объекта. `show` передается в `ctx`, который может быть либо осью `matplotlib` для изображений, либо строкой `Dataframe` для текстов.

Здесь мы создаем объект `SiameseImage`, который выступает подклассом `fasttuple` и содержит три компонента: два изображения и логическое значение, которое будет `True`, если изображения окажутся одной породы. Помимо этого, мы реализуем специальный метод `show`, чтобы он конкатенировал эти два изображения, разделив их черной линией посередине. Не обращайте особого внимания на часть в условии `if` (она выполняет показ `SiameseImage`, когда в качестве изображений выступают изображения Python, а не тензоры). Важной же частью являются три последние строки:

```
class SiameseImage(fasttuple):
    def show(self, ctx=None, **kwargs):
        img1,img2,same_breed = self
        if not isinstance(img1, Tensor):
            if img2.size != img1.size: img2 = img2.resize(img1.size)
            t1,t2 = tensor(img1),tensor(img2)
            t1,t2 = t1.permute(2,0,1),t2.permute(2,0,1)
        else: t1,t2 = img1,img2
        line = t1.new_zeros(t1.shape[0], t1.shape[1], 10)
        return show_image(torch.cat([t1,line,t2], dim=2),
                          title=same_breed, ctx=ctx)
```

Создадим первое `SiameseImage` и убедимся, что наш метод `show` работает:

```
img = PILImage.create(files[0])
s = SiameseImage(img, img, True)
s.show();
```

True



Мы также можем попробовать использовать второе изображение, не относящееся к тому же классу:

```
img1 = PILImage.create(files[1])
s1 = SiameseImage(img, img1, False)
s1.show();
```

False



Важный момент — преобразования передаются через кортежи или их подклассы. Именно поэтому мы предпочли создать в этом экземпляре подкласс `fasttuple` — так мы можем применить любое подходящее для изображений преобразование к `SiameseImage`, и оно будет автоматически применено к каждому изображению кортежа:

```
s2 = Resize(224)(s1)
s2.show();
```

False



Здесь преобразование `Resize` применено к каждому изображению, но не к логическому флагу. Даже если мы используем собственный тип, у нас все равно есть возможность воспользоваться преимуществом всех доступных в библиотеке преобразований, выполняющих аугментацию данных.

Теперь можно перейти к созданию функции `Transform`, с помощью которой мы подготовим данные для сиамской модели. Сначала нужно создать функцию для определения классов изображений:

```
def label_func(fname):
    return re.match(r'^(.*)_\d+.jpg$', fname.name).groups()[0]
```

Для каждого изображения преобразование будет с вероятностью 0,5 извлекать изображение того же класса и возвращать `SiameseImage` с меткой `True` либо вытягивать изображение из другого класса и возвращать `SiameseImage` с меткой `False`. Все это выполняется в закрытой функции `_draw`. Между обучающей и контрольной выборками есть одно отличие, из-за которого преобразование необходимо инициализировать с разделением: в обучающем наборе мы будем делать случайный выбор при каждом чтении изображения, в то время как для контрольного набора мы будем делать единовременный выбор всех экземпляров сразу при инициализации. Таким образом мы будем получать более разнообразные образцы в процессе обучения и всегда одинаковую контрольную выборку:

```
class SiameseTransform(Transform):
    def __init__(self, files, label_func, splits):
        self.labels = files.map(label_func).unique()
        self.lbl2files = {l: [f for f in files if label_func(f) == l]
                         for l in self.labels}
        self.label_func = label_func
        self.valid = {f: self._draw(f) for f in files[splits[1]]}

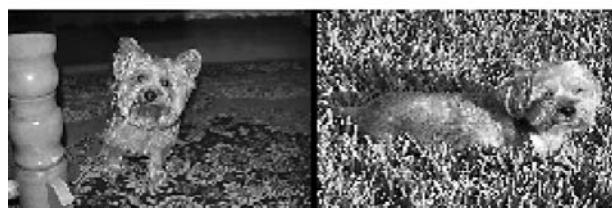
    def encodes(self, f):
        f2,t = self.valid.get(f, self._draw(f))
        img1,img2 = PILImage.create(f),PILImage.create(f2)
        return SiameseImage(img1, img2, t)

    def _draw(self, f):
        same = random.random() < 0.5
        cls = self.label_func(f)
        if not same:
            cls = random.choice([l for l in self.labels if l != cls])
        return random.choice(self.lbl2files[cls]),same
```

После этого можно создать основное преобразование:

```
splits = RandomSplitter()(files)
tfm = SiameseTransform(files, label_func, splits)
tfm(files[0]).show();
```

True



В промежуточном API для сбора данных у нас есть два объекта, с помощью которых можно применить преобразования к набору элементов: `TfmdLists`

и **Datasets**. Если вы запомнили недавнюю информацию, то первый применяет **Pipeline** преобразований, а другой несколько **Pipeline** преобразований параллельно, создавая кортежи. В данном случае наше основное преобразование уже создает кортеж, поэтому мы используем **TfmdLists**:

```
tls = TfmdLists(files, tfm, splits=splits)
show_at(tls.valid, 0);
```

True



И в завершение мы можем передать данные в **DataLoaders**, вызвав метод **dataloaders**. При этом нужно учитывать, что в отличие от **DataBlock** этот метод не принимает **item_tfms** и **batch_tfms**. fastai-объект **DataLoader** имеет несколько хуков, которые именуются согласно событиям. В данном случае хук, применяемый нами к элементам после их извлечения, называется **after_item**, а применяемый к пакету после его создания — **after_batch**:

```
dls = tls.dataloaders(after_item=[Resize(224), ToTensor],
                      after_batch=[IntToFloatTensor, Normalize.from_stats(*imagenet_stats)])
```

Обратите внимание, что нам нужно передать больше преобразований, чем обычно, потому что API блока данных, как правило, добавляет их автоматически:

- **ToTensor** преобразует изображения в тензоры (опять же, применяется он к каждой части кортежа);
- **IntToFloatTensor** преобразует тензор изображений, содержащий целые числа от 0 до 255, в тензор чисел с плавающей запятой и делит их на 255, получая, таким образом, значения между 0 и 1.

Теперь можно перейти к обучению модели с помощью этого **DataLoaders**. Для него потребуется немного больше настроек, чем для обычной модели, предоставляемой **cnn_learner**, поскольку он будет получать не одно, а два изображения. Но мы зайдем в рассмотрением создания такой модели и ее обучением в главе 15.

Резюме

Fastai предоставляет многослойный API. Используя всего одну строку кода, с его помощью можно извлечь данные, когда они находятся в одном из стандартных состояний, что упрощает для начинающих переход к обучению модели и не требует дополнительной подготовки данных. Затем высокоуровневый API блока данных обеспечивает дополнительную гибкость, позволяя смешивать и сопоставлять составляющие части. Под ним промежуточный API дает еще большую гибкость для применения преобразований к элементам. Это, скорее всего, и понадобится в реальных задачах, и мы надеемся, что так этап преобразования данных станет максимально простым.

Вопросник

1. Почему мы говорим, что fastai предоставляет «многослойный» API? Что это значит?
2. Зачем в `Transform` нужен метод `decode`? Что он делает?
3. Зачем в `Transform` нужен метод `setup`? Каково его назначение?
4. Как работает `Transform` при вызове кортежа?
5. Какие методы вам нужно реализовать при написании собственной функции `Transform`?
6. Напишите преобразование `Normalize`, которое полностью нормализует элементы (вычитает среднее и производит деление на стандартное отклонение датасета) и может декодировать эту операцию. Постарайтесь не подглядывать!
7. Напишите `Transform`, выполняющую нумеризацию токенизованных текстов (оно должно автоматически создавать словарь из просмотренного датасета и содержать метод `decode`). Можете заглянуть в исходный код fastai для получения дополнительной информации.
8. Что такое `Pipeline`?
9. Что такое `TfmdLists`?
10. Что такое `Datasets`? Чем он отличается от `TfmdLists`?
11. Почему `TfmdLists` и `Datasets` имеют множественное число?
12. Как можно создать `DataLoaders` из `TfmdLists` или `Datasets`?
13. Как передавать `item_tfms` и `batch_tfms` при построении `DataLoaders` из `TfmdLists` или `Datasets`?
14. Что нужно сделать, если вы хотите, чтобы ваши собственные элементы работали с такими методами, как `show_batch` или `show_results`?

15. Почему мы можем легко применить fastai-преобразования, аугментирующие данные, к созданному нами `SiamesePair`?

Дополнительные задания

1. Используйте промежуточный API для подготовки данных в `DataLoaders` в своих датасетах. Попробуйте это с датасетом Pet и датасетом Adult из главы 1.
2. Ознакомьтесь с инструкцией к Siamese в документации fastai (<https://docs.fast.ai/>), чтобы узнать, как настраивать поведение `show_batch` и `show_results` для новых типов элементов. Реализуйте это в своем проекте.

Сфера применения fastai: обобщение

Поздравляем! Вы прошли все главы книги, охватывающие ключевые практические моменты обучения моделей и использования глубокого обучения. Теперь вы знаете, как применять все поддерживаемые fastai возможности в разных сферах, а также как настраивать их с помощью API блока данных и функций потерь. Кроме того, вы даже научились создавать нейронную сеть с нуля и обучать ее! (Надеемся, что теперь вы также знаете, какие вопросы нужно задавать, чтобы убедиться в положительном вкладе ваших проектов в общество).

Имеющихся сейчас знаний вам будет достаточно для создания полноценных рабочих прототипов многих типов приложений на базе нейронных сетей. Но более важно то, что это поможет вам понять возможности и ограничения моделей глубокого обучения, а также способы проектирования систем, успешно с ними согласующихся.

В оставшейся части книги мы будем разбирать эти приложения по деталям, чтобы понять основы, на которых они построены. Это знание очень важно для практикующего глубокое обучение, так как позволяет инспектировать создаваемые модели и выполнять их отладку, а также создавать новые приложения, соответствующие конкретным нуждам ваших проектов.

ЧАСТЬ III

**Основы глубокого
обучения**

ГЛАВА 12

Языковая модель с нуля

Вот теперь мы готовы к настоящему углублению в... глубокое обучение. Вы уже умеете обучать простую нейронную сеть, но как создавать эталонные модели? Здесь мы раскроем все тайны и начнем с языковых моделей.

В главе 10 вы видели, как тонко настраивать предварительно обученную языковую модель для построения классификатора текста. В этой же главе мы объясним, что именно находится внутри этой модели и что такое RNN. Для начала соберем данные, которые позволяют в ускоренном темпе создавать прототипы различных моделей.

Данные

Работу над каждой новой задачей мы всегда начинаем с создания наиболее простого датасета, который позволит нам быстро и легко опробовать разные методы и интерпретировать результаты. Когда мы несколько лет назад занялись разработкой языковых моделей, то не могли найти ни одного датасета, который бы позволял осуществлять быстрое прототипирование, поэтому создали такой сами. Назвали мы его *Human Numbers* (числа прописью), и содержит он всего-навсего первые 10 000 чисел, написанных на английском.



СЛОВО ДЖЕРЕМИ

Одна из наиболее распространенных ошибок аналитиков заключается в использовании для анализа не соответствующих той или иной ситуации датасетов. В особенности это относится к тому, что большинство специалистов склонны действовать слишком большие и сложные наборы данных.

Скачать, извлечь и рассмотреть датасет можно привычным способом:

```
from fastai.text.all import *
path = untar_data(URLs.HUMAN_NUMBERS)
```

```
path.ls()
(#2) [Path('train.txt'), Path('valid.txt')]
```

Заглянем в эти два файла: для начала мы объединим все тексты и проигнорируем обусловленное датасетом разделение train/valid (к этому мы вернемся позднее):

```
lines = L()
with open(path/'train.txt') as f: lines += L(*f.readlines())
with open(path/'valid.txt') as f: lines += L(*f.readlines())
lines

(#9998) ['one \n', 'two \n', 'three \n', 'four \n', 'five \n', 'six \n', 'seven
\n', 'eight \n', 'nine \n', 'ten \n'...]
```

Мы берем все эти строки и конкатенируем их в один большой поток. Переход от одного числа к другому мы отмечаем разделителем в виде точки:

```
text = ' . '.join([l.strip() for l in lines])
text[:100]

'one . two . three . four . five . six . seven . eight . nine . ten . eleven .
twelve . thirteen . fo'
```

Токенизировать этот датасет можно, разделив его по пробелам:

```
tokens = text.split(' ')
tokens[:10]

['one', '.', 'two', '.', 'three', '.', 'four', '.', 'five', '.']
```

Для нумеризации нужно создать список всех уникальных токенов (*словарь*):

```
vocab = L(*tokens).unique()
vocab

(#30) ['one', '.', 'two', 'three', 'four', 'five', 'six', 'seven', 'eight', 'nine'...]
```

Затем можно преобразовать токены в числа, используя поиск по индексу каждого в словаре:

```
word2idx = {w:i for i,w in enumerate(vocab)}
nums = L(word2idx[i] for i in tokens)
nums

(#63095) [0,1,2,1,3,1,4,1,5,1...]
```

Теперь, когда у нас есть небольшой датасет, на котором языковое моделирование не должно составить трудностей, мы готовы переходить к созданию первой модели.

Первая языковая модель с нуля

Для преобразования имеющегося материала в нейронную сеть потребуется просто прописать, что мы собираемся прогнозировать каждое слово на основе трех предшествующих. Для начала создадим список, состоящий из всех последовательностей трех слов, выступающих в качестве независимых переменных, и следующего за каждой такой последовательностью слова, выступающего в качестве зависимой переменной.

Это можно сделать в чистом Python, для чего мы сначала используем токены, просто чтобы понять, как все будет выглядеть:

```
L((tokens[i:i+3], tokens[i+3]) for i in range(0,len(tokens)-4,3))

(#21031) [(['one', '.', 'two'], '.'), ('.', 'three', '.'), ('four', ('four',
> '.', 'five', '.'), ('.', 'six', '.'), ('seven', ('seven', '.', 'eight'),
> '.'), ('.', 'nine', '.'), ('ten', ('ten', '.', 'eleven'), '.'), ('.',
> 'twelve', '.'), ('thirteen', ('thirteen', '.', 'fourteen'), '.'), ('.',
> 'fifteen', '.'), ('sixteen')...]
```

Теперь мы проделаем это с тензорами нумеризованных значений, которые модель в итоге будет использовать:

```
seqs = L((tensor(nums[i:i+3]), nums[i+3]) for i in range(0,len(nums)-4,3))
seqs

(#21031) [(tensor([0, 1, 2]), 1), (tensor([1, 3, 1]), 4), (tensor([4, 1, 5]),
> 1), (tensor([1, 6, 1]), 7), (tensor([7, 1, 8]), 1), (tensor([1, 9, 1]),
> 10), (tensor([10, 1, 11]), 1), (tensor([1, 12, 1]), 13), (tensor([13, 1,
> 14]), 1), (tensor([1, 15, 1]), 16)...]
```

Можно легко объединить их в пакеты, используя класс `DataLoader`. Пока же мы разделим последовательности произвольно:

```
bs = 64
cut = int(len(seqs) * 0.8)
dls = DataLoaders.from_dsets(seqs[:cut], seqs[cut:], bs=64, shuffle=False)
```

Теперь создадим архитектуру нейронной сети, получающую на входе три слова и возвращающую прогноз вероятности каждого возможного следующего слова словаря. Здесь мы применим три стандартных линейных слоя, но с двумя доработками.

Первая из них подразумевает, что начальный линейный слой будет использовать в качестве активаций только вложение первого слова, второй будет использовать вложение второго слова плюс выходные активации первого слоя, а третий будет использовать вложение третьего слова плюс выходные активации второго слоя. Основной эффект от этого в том, что каждое слово интерпретируется в информационном контексте всех предшествующих ему слов.

Вторая доработка подразумевает, что каждый из этих трех слоев будет использовать одну и ту же матрицу весов. Способ, которым одно слово влияет на активации от предыдущих слов, не должен меняться в зависимости от позиции слова. Другими словами, значения активаций будут меняться по мере продвижения данных через слои, но сами веса слоев от слоя к слою меняться не будут. Таким образом, слой изучает не одну позицию последовательности, а должен обучиться обрабатывать все позиции.

Поскольку веса слоев не меняются, вы можете воспринимать последовательные слои как «один и тот же слой» в повторении. На самом деле PyTorch предлагает здесь еще больше конкретики: с его помощью можно создать всего один слой и использовать его несколько раз.

Языковая модель в PyTorch

Теперь создадим модуль языковой модели, который уже описывали ранее:

```
class LMModel1(Module):
    def __init__(self, vocab_sz, n_hidden):
        self.i_h = nn.Embedding(vocab_sz, n_hidden)
        self.h_h = nn.Linear(n_hidden, n_hidden)
        self.h_o = nn.Linear(n_hidden, vocab_sz)

    def forward(self, x):
        h = F.relu(self.h_h(self.i_h(x[:,0])))
        h = h + self.i_h(x[:,1])
        h = F.relu(self.h_h(h))
        h = h + self.i_h(x[:,2])
        h = F.relu(self.h_h(h))
        return self.h_o(h)
```

Как видите, мы создали три слоя.

- Слой вложений (`i_h` означает *ввод в скрытый*).
- Линейный слой для создания активаций для следующего слова (`h_h` означает передачу *из скрытого в скрытый*).
- Заключительный слой для прогнозирования четвертого слова (`h_o` означает передачу *из скрытого на выход*).

Для наглядности отобразим все это в виде иллюстрации, схематично отражающей простую нейронную сеть. На рис. 12.1 показано, как мы будем представлять такую сеть с одним скрытым слоем.

Каждая фигура представляет активации: прямоугольник для входных данных, круг для активаций скрытого (внутреннего) слоя и треугольник — для выходных

активаций. Эти фигуры (обобщенные на рис. 12.2) мы будем использовать во всех диаграммах текущей главы.

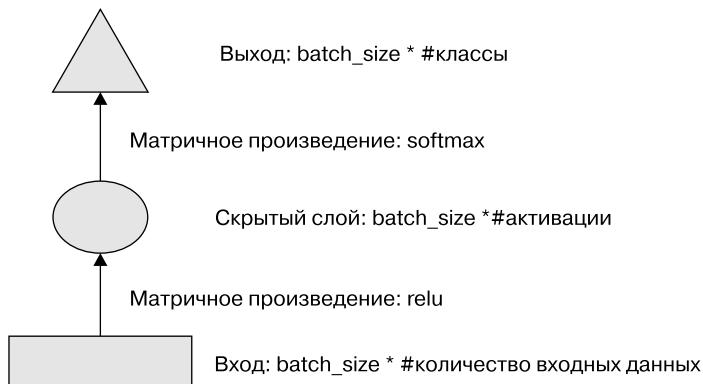


Рис. 12.1. Графическое представление простой нейронной сети

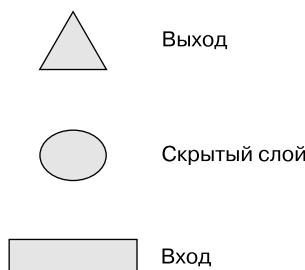


Рис. 12.2. Фигуры, используемые в графических представлениях

Каждая стрелка представляет фактическое вычисление слоя, то есть линейный слой, сопровождаемый функцией активации. Этими обозначениями на рис. 12.3 показана наша простая языковая модель.

В целях упрощения мы удалили из каждой стрелки подробности вычисления слоев. Помимо этого, мы применили к стрелкам цветокодирование, подразумевающее, что стрелки одного цвета имеют одинаковые матрицы весов. Например, все входные слои используют одну матрицу вложений, поэтому для них используется один цвет.

Попробуем обучить эту модель и оценим результат:

```

learn = Learner(dls, LMModel1(len(vocab), 64), loss_func=F.cross_entropy,
                 metrics=accuracy)
learn.fit_one_cycle(4, 1e-3)
  
```

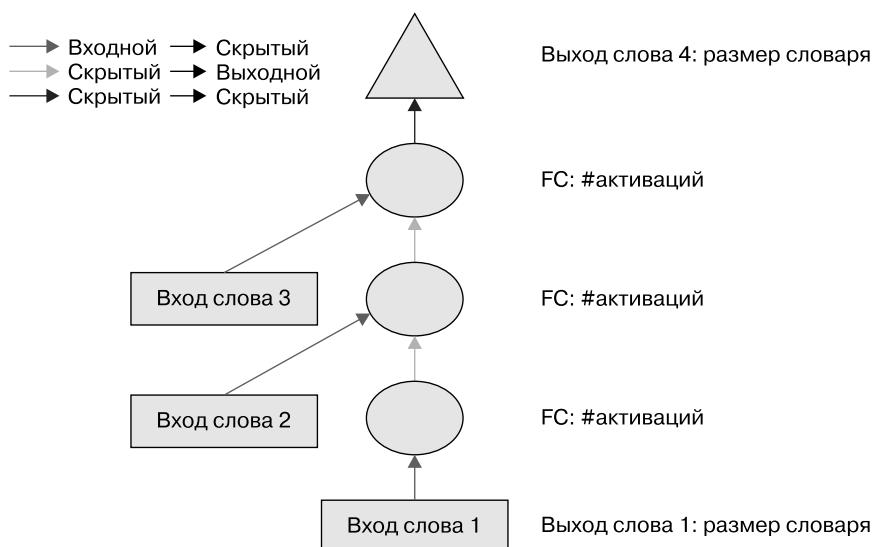


Рис. 12.3. Представление простой языковой модели

epoch	train_loss	valid_loss	accuracy	time
0	1.824297	1.970941	0.467554	00:02
1	1.386973	1.823242	0.467554	00:02
2	1.417556	1.654497	0.494414	00:02
3	1.376440	1.650849	0.494414	00:02

Чтобы понять, насколько хорош текущий результат, посмотрим, какой бы мы получили от самой простой модели. В этом случае всегда можно спрогнозировать наиболее распространенный токен, так что выясним, какой из них чаще других является целью в контрольной выборке:

```
n,counts = 0,torch.zeros(len(vocab))
for x,y in dls.valid:
    n += y.shape[0]
    for i in range_of(vocab): counts[i] += (y==i).long().sum()
idx = torch.argmax(counts)
idx, vocab[idx.item()], counts[idx].item()/n

(tensor(29), 'thousand', 0.15165200855716662)
```

Наиболее распространен индекс 29, и соответствует он токену `thousand`. Постоянное прогнозирование этого токена даст нам точность примерно 15 %, так что наш результат намного лучше.



СЛОВО АЛЕКСИСУ

Изначально я предположил, что самым частым токеном будет разделитель, так как он сопровождает каждое число. Но при рассмотрении tokens я вспомнил, что большие числа пишутся несколькими словами, поэтому в процессе счета до 10 000 приходится много раз писать «thousand»: five thousand, five thousand and one, five thousand and two и т. д. Вот как бывает. Очень полезно просматривать данные, чтобы обнаруживать не только их тонкие детали, но и вполне очевидные.

Мы получили неплохую базовую модель. Теперь выполним ее рефакторинг с помощью цикла.

Первая рекуррентная нейронная сеть

Если рассмотреть код нашего модуля, то мы можем упростить его, заменив повторяющийся код, вызывающий слои, на цикл `for`. Помимо упрощения кода, это также даст нам удобную возможность с той же эффективностью применять модуль к последовательностям токенов разной длины — мы не будем ограничены длиной их списков в три единицы:

```
class LMModel2(Module):
    def __init__(self, vocab_sz, n_hidden):
        self.i_h = nn.Embedding(vocab_sz, n_hidden)
        self.h_h = nn.Linear(n_hidden, n_hidden)
        self.h_o = nn.Linear(n_hidden, vocab_sz)

    def forward(self, x):
        h = 0
        for i in range(3):
            h = h + self.i_h(x[:,i])
            h = F.relu(self.h_h(h))
        return self.h_o(h)
```

Убедимся, что при использовании рефакторинга мы получаем те же результаты:

```
learn = Learner(dls, LMModel2(len(vocab), 64), loss_func=F.cross_entropy,
                 metrics=accuracy)
learn.fit_one_cycle(4, 1e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	1.816274	1.964143	0.460185	00:02
1	1.423805	1.739964	0.473259	00:02
2	1.430327	1.685172	0.485382	00:02
3	1.388390	1.657033	0.470406	00:02

Аналогичным образом можно выполнить рефакторинг графического представления, что и показано на рис. 12.4 (помимо прочего, мы удаляем детали размеров активаций и используем те же цвета стрелок, что и на рис. 12.3).

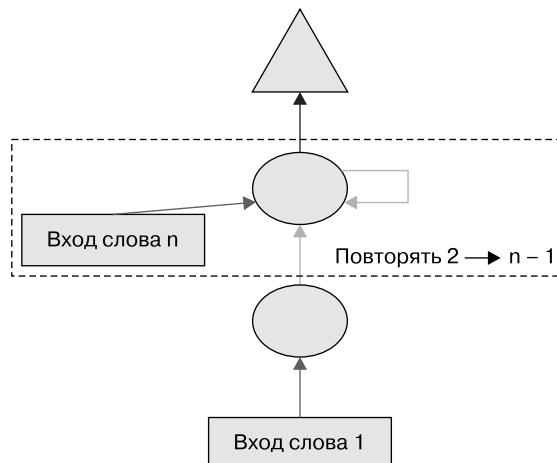


Рис. 12.4. Простая рекуррентная нейронная сеть

Вы увидите, что набор активаций обновляется каждый раз с помощью цикла и сохраняется в переменной h , которая называется *скрытым состоянием*.



ТЕРМИН: СКРЫТОЕ СОСТОЯНИЕ

Активации, обновляемые на каждом шаге рекуррентной нейронной сети.

Нейронная сеть, определенная с использованием подобного цикла, называется *рекуррентной нейронной сетью* (RNN). Важно понимать, что RNN представляет не новую усложненную архитектуру, а простой рефакторинг многослойной нейронной сети с помощью цикла `for`.



СЛОВО АЛЕКСИСУ

Если бы эти сети назывались «циклические нейронные сети», или LNN, это звучало бы на 50 % менее устрашающе.

Теперь, когда мы знаем, что такое RNN, пора приступить к ее улучшению.

Улучшение RNN

При внимательном рассмотрении нашего кода для RNN можно заметить проблему в том, что мы инициализируем скрытое состояние с нулем для каждой новой входной последовательности. Почему это проблема? Мы сделали наши последовательности образцов короткими, чтобы они могли легко вписаться в пакеты. Но если упорядочить эти образцы правильно, то модель, считывая их последовательности по порядку, будет вынуждена обрабатывать длинные отрезки исходной последовательности.

Помимо этого, стоит рассмотреть возможность получения дополнительного сигнала: почему мы должны прогнозировать только четвертое слово, когда можно использовать промежуточные прогнозы для второго и третьего слов? Рассмотрим возможность реализации этих изменений, начиная с добавления нового состояния.

Управление состоянием RNN

Инициализируя скрытое состояние модели для каждого нового образца как ноль, мы, по сути, выбрасываем всю информацию о просмотренных до этого предложениях, вследствие чего модель не знает, где именно она находится в общей последовательности счета. Это легко исправить простым перемещением инициализации скрытого состояния в `__init__`.

Но эта корректировка создаст свою небольшую, но важную проблему. В результате глубина нашей нейронной сети фактически станет равна общему числу токенов документа. Например, если в датасете будет 10 000 токенов, то мы так создадим сеть с 10 000 слоев. Чтобы понять, почему это происходит, рассмотрите исходное графическое представление нашей рекуррентной нейронной сети на рис. 12.3 до применения цикла `for`. Вы можете заметить, что каждый слой соотносится с одним входным токеном. Когда мы говорим о представлении рекуррентной нейронной сети до применения цикла `for`, мы называем это *неразвернутым представлением*, рассмотрение которого нередко помогает лучше понять RNN.

Проблема с нейронной сетью из 10 000 слоев в том, что если и когда вы достигнете 10 000-го слова датасета, то также будете вынуждены вычислить производные на всем пути обратно к первому слову. Это потребует очень много времени и ресурсов памяти. При этом вряд ли вам удастся уместить в GPU хотя бы один мини-пакет.

Решением данной проблемы будет сообщить PyTorch, что мы не хотим выполнять обратное распространение для вычисления производных через всю неявную нейронную сеть. Вместо этого мы сохраним только последние три слоя градиентов. Для удаления всей истории градиентов в PyTorch используется метод `detach`.

Ниже прописана новая версия нашей RNN. Теперь она поддерживает сохранение состояния, запоминая свои активации между разными вызовами `forward`, которые представляют ее использование для различных образцов из пакета:

```
class LMMModel3(Module):
    def __init__(self, vocab_sz, n_hidden):
        self.i_h = nn.Embedding(vocab_sz, n_hidden)
        self.h_h = nn.Linear(n_hidden, n_hidden)
        self.h_o = nn.Linear(n_hidden, vocab_sz)
        self.h = 0

    def forward(self, x):
        for i in range(3):
            self.h = self.h + self.i_h(x[:,i])
            self.h = F.relu(self.h_h(self.h))
        out = self.h_o(self.h)
        self.h = self.h.detach()
        return out

    def reset(self): self.h = 0
```

У этой модели будут одинаковые активации независимо от выбранной нами длины последовательности, потому что скрытое состояние будет запоминать последнюю активацию из предыдущего пакета. Единственное, что будет отличаться, — это вычисляемые на каждом этапе градиенты: они будут вычисляться только для длины последовательности токенов в прошлом, а не всего потока. Такой подход называется *алгоритмом обратного распространения во времени* (backpropagation through time — BPTT).



ТЕРМИН: ОБРАТНОЕ РАСПРОСТРАНЕНИЕ ВО ВРЕМЕНИ

Рассмотрение нейронной сети с фактически одним слоем в каждом шаге времени (обычно циклом `for`) как одной большой модели и вычисление для нее градиентов стандартным способом. Во избежание перерасхода памяти и времени мы обычно используем усеченный BPTT, который через каждые несколько шагов по времени «отделяет» историю этапов вычислений в скрытом состоянии.

Чтобы использовать `LMMModel3`, нам нужно гарантировать, что образцы будут просматриваться в определенном порядке. Как мы видели в главе 10, если первой строкой первого пакета является `dset[0]`, то во втором пакете первой строкой должна быть `dset[1]`, чтобы модель видела переход текста.

В главе 10 за нас это делал `LMDataloader`. Здесь же мы будем реализовывать процесс сами.

Для этого мы перераспределим наш датасет. Сначала мы разделим образцы на `m = len(dset) // bs` групп (это равнозначно разделению всего конкатенирован-

ногого датасета, к примеру, на 64 части одинакового размера, так как здесь мы используем `bs=64`). `m` — это длина каждой из этих частей. Например, если мы используем весь датасет (хотя мы вот-вот разделим его на обучающую и контрольную части), то у нас получится следующее:

```
m = len(seqs)//bs
m,bs,len(seqs)

(328, 64, 21031)
```

Первый пакет будет состоять из этих образцов:

```
(0, m, 2*m, ..., (bs-1)*m)
```

Второй из этих:

```
(1, m+1, 2*m+1, ..., (bs-1)*m+1)
```

И так далее. Таким образом, каждую эпоху модель будет на каждой строке пакета видеть фрагмент непрерывного текста размером $3*m$ (поскольку каждый текст имеет размер 3).

Данную переиндексацию выполняет следующая функция:

```
def group_chunks(ds, bs):
    m = len(ds) // bs
    new_ds = L()
    for i in range(m): new_ds += L(ds[i + m*j] for j in range(bs))
    return new_ds
```

Затем мы при создании `DataLoaders` просто передаем `drop_last=True`, чтобы отбросить последний пакет, не имеющий формы `bs`. Мы также передаем `shuffle=False`, обеспечивая чтение текстов по порядку:

```
cut = int(len(seqs) * 0.8)
dls = DataLoaders.from_dsets(
    group_chunks(seqs[:cut], bs),
    group_chunks(seqs[cut:], bs),
    bs=bs, drop_last=True, shuffle=False)
```

Последнее, что мы добавляем, — это небольшая доработка цикла обучения с помощью `Callback`. В главе 16 мы поговорим об обратных вызовах подробнее, что же касается этого, то он будет вызывать метод модели `reset` в начале каждой эпохи и перед каждой контрольной стадией. Поскольку мы реализовали этот метод, чтобы установить скрытое состояние модели на ноль, это гарантирует, что чтение тех непрерывных фрагментов текста мы будем начинать с чистым состоянием. В дополнение к этому можно продлить цикл обучения:

```
learn = Learner(dls, LMModel3(len(vocab), 64), loss_func=F.cross_entropy,
                 metrics=accuracy, cbs=ModelResetter)
learn.fit_one_cycle(10, 3e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	1.677074	1.827367	0.467548	00:02
1	1.282722	1.870913	0.388942	00:02
2	1.090705	1.651793	0.462500	00:02
3	1.005092	1.613794	0.516587	00:02
4	0.965975	1.560775	0.551202	00:02
5	0.916182	1.595857	0.560577	00:02
6	0.897657	1.539733	0.574279	00:02
7	0.836274	1.585141	0.583173	00:02
8	0.805877	1.629808	0.586779	00:02
9	0.795096	1.651267	0.588942	00:02

Уже лучше! Следующим шагом будет использование большего числа целей и их сравнение с промежуточными прогнозами.

Создание дополнительного сигнала

Еще одна проблема нашего текущего подхода в том, что мы прогнозируем только одно выходное слово на основе трех входных. В результате объем сигнала, отправляемого обратно для обновления весов, не настолько велик, насколько мог бы быть. В данном случае более эффективным будет делать прогноз следующего слова после каждого слова, а не после трех, как показано на рис. 12.5.

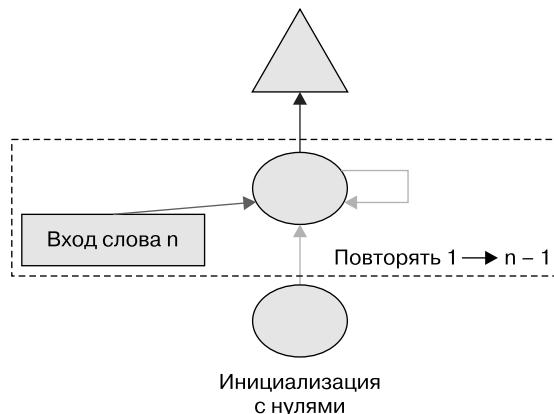


Рис. 12.5. RNN делает прогноз после каждого слова

Такое поведение легко добавить. Сначала нужно изменить данные так, чтобы в зависимой переменной было каждое из трех следующих слов после каждого из трех входных слов. Вместо `s1` мы используем атрибут `s1` (означающий длину последовательности) и несколько увеличим его:

```
s1 = 16
seqs = L((tensor(nums[i:i+s1]), tensor(nums[i+1:i+s1+1]))
          for i in range(0, len(nums)-s1-1, s1))
cut = int(len(seqs) * 0.8)
dls = DataLoaders.from_dsets(group_chunks(seqs[:cut], bs),
                               group_chunks(seqs[cut:], bs),
                               bs=bs, drop_last=True, shuffle=False)
```

Взгляните на первый элемент `seqs`, который содержит два списка одного размера. Второй список имеет тот же размер, что и первый, но смешен на один элемент:

```
[L(vocab[o] for o in s) for s in seqs[0]]
[(#16) ['one', '.', 'two', '.', 'three', '.', 'four', '.', 'five', '.',...
      (#16) ['.', 'two', '.', 'three', '.', 'four', '.', 'five', '.', 'six'...]]
```

Теперь нужно изменить модель так, чтобы она выводила прогноз после каждого слова:

```
class LMModel4(Module):
    def __init__(self, vocab_sz, n_hidden):
        self.i_h = nn.Embedding(vocab_sz, n_hidden)
        self.h_h = nn.Linear(n_hidden, n_hidden)
        self.h_o = nn.Linear(n_hidden, vocab_sz)
        self.h = 0

    def forward(self, x):
        outs = []
        for i in range(s1):
            self.h = self.h + self.i_h(x[:, i])
            self.h = F.relu(self.h_h(self.h))
            outs.append(self.h_o(self.h))
        self.h = self.h.detach()
        return torch.stack(outs, dim=1)

    def reset(self): self.h = 0
```

Эта модель будет возвращать вывод в форме `bs x s1 x vocab_sz` (поскольку мы сложили все в `dim=1`). Цели имеют форму `bs x s1`, поэтому нужно придать им более плоскую форму, прежде чем использовать в `F.cross_entropy`:

```
def loss_func(inp, targ):
    return F.cross_entropy(inp.view(-1, len(vocab)), targ.view(-1))
```

Теперь можно использовать эту функцию потерь для обучения модели:

```
learn = Learner(dls, LMMModel4(len(vocab), 64), loss_func=loss_func,
                 metrics=accuracy, cbs=ModelResetter)
learn.fit_one_cycle(15, 3e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	3.103298	2.874341	0.212565	00:01
1	2.231964	1.971280	0.462158	00:01
2	1.711358	1.813547	0.461182	00:01
3	1.448516	1.828176	0.483236	00:01
4	1.288630	1.659564	0.520671	00:01
5	1.161470	1.714023	0.554932	00:01
6	1.055568	1.660916	0.575033	00:01
7	0.960765	1.719624	0.591064	00:01
8	0.870153	1.839560	0.614665	00:01
9	0.808545	1.770278	0.624349	00:01
10	0.758084	1.842931	0.610758	00:01
11	0.719320	1.799527	0.646566	00:01
12	0.683439	1.917928	0.649821	00:01
13	0.660283	1.874712	0.628581	00:01
14	0.646154	1.877519	0.640055	00:01

Нужно продолжать обучение, поскольку задача изменилась и стала сложнее. Но в итоге мы приедем к хорошему результату... по крайней мере иногда. Если выполнить цикл обучения несколько раз, то вы заметите, что получаются достаточно разнящиеся результаты. Причина в том, что по факту у нас здесь очень глубокая нейронная сеть, что может приводить к очень большим или очень малым градиентам. Чуть позже в этой главе мы рассмотрим решение для этой ситуации.

Теперь же естественным способом улучшения модели будет ее углубление: сейчас в нашей базовой RNN между скрытым состоянием и выходными активациями есть всего один линейный слой, так что стоит попробовать добиться лучших результатов, добавив еще.

Многослойные RNN

Как показано на рис. 12.6, в многослойной RNN передача активаций происходит из одной рекуррентной нейронной сети во вторую рекуррентную нейронную сеть.

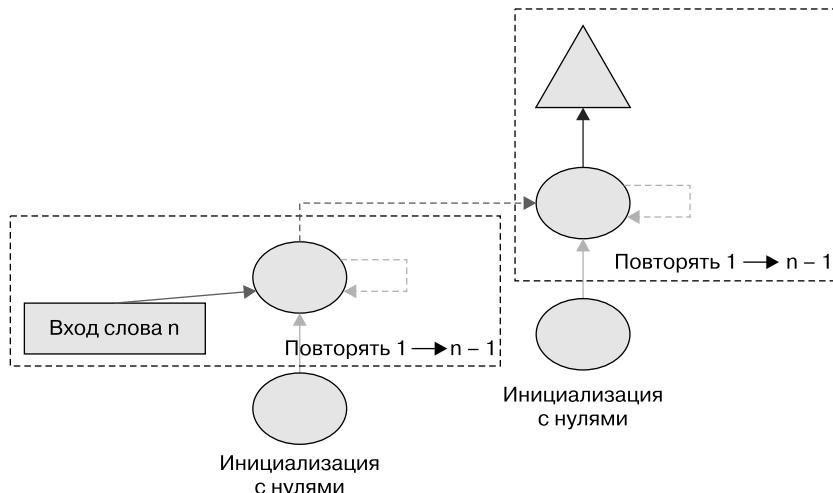


Рис. 12.6. Двухслойная RNN

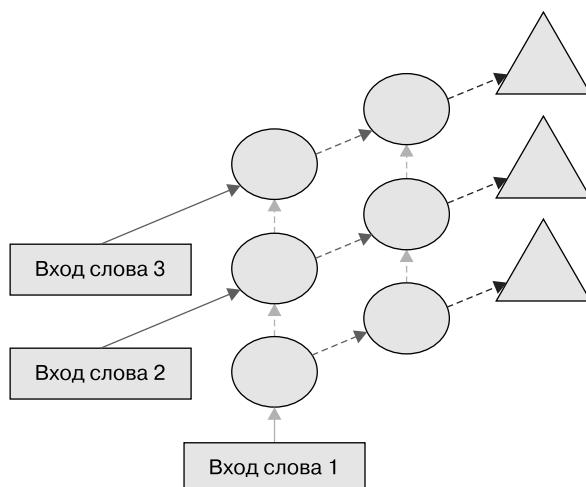


Рис. 12.7. Двухслойная развернутая RNN

Развернутое представление показано на рис. 12.7 (оно аналогично рис. 12.3).

Посмотрим, как это реализовать на практике.

Модель

Можно сэкономить время, использовав PyTorch-класс RNN, который реализует все то, что мы создали ранее, дополнительно позволяя объединить несколько RNN:

```
class LMModel5(Module):
    def __init__(self, vocab_sz, n_hidden, n_layers):
        self.i_h = nn.Embedding(vocab_sz, n_hidden)
        self.rnn = nn.RNN(n_hidden, n_hidden, n_layers, batch_first=True)
        self.h_o = nn.Linear(n_hidden, vocab_sz)
        self.h = torch.zeros(n_layers, bs, n_hidden)

    def forward(self, x):
        res,h = self.rnn(self.i_h(x),
                         self.h) self.h = h.detach()
        return self.h_o(res)

    def reset(self): self.h.zero_()

learn = Learner(dls, LMModel5(len(vocab), 64, 2),
                loss_func=CrossEntropyLossFlat(),
                metrics=accuracy, cbs=ModelResetter)
learn.fit_one_cycle(15, 3e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	3.055853	2.591640	0.437907	00:01
1	2.162359	1.787310	0.471598	00:01
2	1.710663	1.941807	0.321777	00:01
3	1.520783	1.999726	0.312012	00:01
4	1.330846	2.012902	0.413249	00:01
5	1.163297	1.896192	0.450684	00:01
6	1.033813	2.005209	0.434814	00:01
7	0.919090	2.047083	0.456706	00:01
8	0.822939	2.068031	0.468831	00:01
9	0.750180	2.136064	0.475098	00:01
10	0.695120	2.139140	0.485433	00:01
11	0.655752	2.155081	0.493652	00:01
12	0.629650	2.162583	0.498535	00:01
13	0.613583	2.171649	0.491048	00:01
14	0.604309	2.180355	0.487874	00:01

Итог несколько разочаровывает... Предыдущая однослойная RNN справлялась лучше. Почему же? Причина в том, что использование углубленной модели ведет к взрывающимся, или исчезающим, активациям.

Взрывающиеся, или исчезающие, активации

На практике создать точную модель на основе подобной RNN достаточно сложно. Результаты улучшатся, если мы будем реже вызывать `detach` и добавим больше слоев, — это даст нашей RNN больший временной горизонт для обучения и создания более богатых признаков. Но это также означает, что обучать придется более глубокую модель. Основная сложность в разработке глубокого обучения как раз и заключалась в нахождении эффективного способа обучения таких моделей.

Суть проблемы лежит в последствиях многократного умножения на матрицу. Чтобы легче это представить, подумайте, к чему может привести многократное умножение на число. Например, если вы будете умножать на 2, начав с единицы, то получите последовательность 1, 2, 4, 8... и после 32 повторений у вас уже будет 4 294 967 296. Аналогичная проблема возникает, если вы умножаете на 0,5: сначала получается 0,5, 0,25, 0,125... и через 32 операции уже 0,00000000023. Как видите, умножение на число даже немного меньшее или большее единицы после всего нескольких повторений приводит к взрыву или исчезновению исходного числа.

А так как матричное умножение подразумевает простое умножение чисел с последующим их сложением, то при повторении данного процесса происходит все то же самое. В этом и состоит нейронная сеть: каждый дополнительный слой добавляет очередную операцию матричного умножения. Это означает, что нейронной сети очень легко достигнуть чрезвычайно больших или малых чисел.

Здесь и кроется зерно проблемы, потому что способ, которым компьютеры хранят числа (используя *плавающую запятую*), подразумевает, что чем дальше числа уходят от нуля, тем больше утрачивается их точность. Схема на рис. 12.8, взятая из прекрасной статьи *What You Never Wanted to Know about Floating Point but Will Be Forced to Find Out* (https://oreil.ly/c_kG9) («Чего вы никогда не хотели знать о плавающей запятой, но будете вынуждены выяснить»), показывает, что по ходу числовой оси точность чисел с плавающей запятой изменяется.

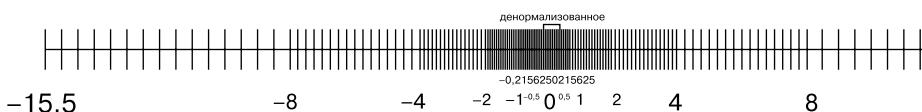


Рис. 12.8. Точность чисел с плавающей запятой

Эта неточность приводит к тому, что в глубоких сетях очень часто градиенты, вычисляемые для обновления весов, оказываются равны нулю или бесконечности. Такое явление принято называть *исчезающими, или взрывающимися, градиентами*. В итоге получается, что в SGD веса либо не обновляются совсем, либо подскакивают до бесконечности. В любом из этих сценариев их улучшения в ходе обучения не произойдет.

Исследователи разработали способы решения этой проблемы, которые мы будем рассматривать несколько позже. Один из них подразумевает изменение определения слоя так, чтобы снизить вероятность появления взрывающихся активаций. Подробно этот процесс мы рассмотрим в главе 13, когда будем говорить о нормализации пакетов, а также в главе 14 при обсуждении ResNets, хотя на практике эти подробности не играют существенной роли (если только вы не относитесь к исследователям, ищущим новые подходы для решения данной проблемы). Еще одна стратегия подразумевает повышенную осторожность при инициализации, что мы будем обсуждать в главе 17.

Для RNN с целью избегания взрывающихся активаций зачастую используются два типа слоев: слои *управляемых рекуррентных блоков* (GRU) и *долгой краткосрочной памяти* (LSTM). Оба этих вида доступны в PyTorch и выступают в качестве заменителей слоя RNN. В данной книге мы рассмотрим только LSTM. Тем не менее онлайн доступно немало обучающих материалов по GRU, которые представляют уменьшенный вариант структуры LSTM.

LSTM

LSTM — это архитектура, которая была представлена в 1997 году Юргеном Шмидхубером (Jürgen Schmidhuber) и Зеппом Хохрайтером (Sepp Hochreiter). Ее особенность — в использовании не одного, а двух скрытых состояний. В нашей базовой RNN скрытое состояние — это вывод RNN в предыдущем шаге времени. Отвечает это скрытое состояние за две вещи.

- Наличие верной информации, позволяющей выходным слоям правильно прогнозировать следующий токен.
- Удержание в памяти всего происходившего в предложении.

Рассмотрите, к примеру, предложения *Henry has a dog and he likes his dog very much* и *Sophie has a dog and she likes her dog very much*. Совершенно очевидно, что RNN нужно запомнить имя в начале предложения, чтобы успешно спрогнозировать *he/she* или *his/her*.

На практике RNN очень плохо справляются с удержанием информации о происходившем в предложении существенно раньше, что и побудило создать еще одно скрытое состояние (называемое *состоянием ячейки*) в LSTM. Состояние ячейки

будет отвечать за хранение *длгой краткосрочной памяти*, в то время как скрытое состояние будет сосредоточено на прогнозировании следующего токена. Давайте более подробно рассмотрим реализацию этого процесса и создадим LSTM с нуля.

Создание LSTM с нуля

Для построения LSTM сначала нам нужно понять ее архитектуру, которая схематично показана на рис. 12.9.

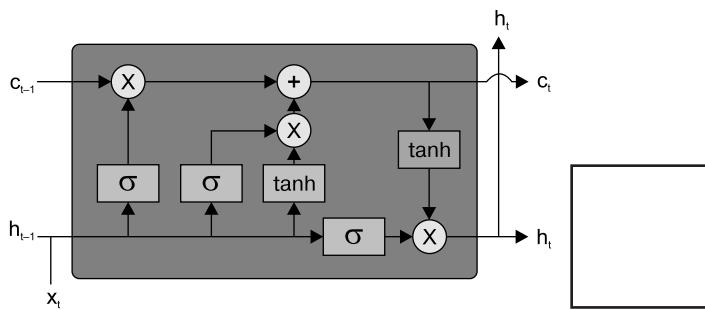


Рис. 12.9. Архитектура LSTM

На этой картинке вводный x_t входит слева вместе с предыдущим скрытым состоянием (h_{t-1}) и состоянием ячейки (c_{t-1}). Четыре оранжевых блока представляют четыре слоя (нейронные сети), где активация может быть либо сигмоидой (σ), либо \tanh (гиперболическим тангенсом). \tanh — это та же сигмоида, перемасштабированная в диапазон от -1 до 1 . Ее математическое выражение выглядит так:

$$\tanh(x) = \frac{e^x + e^{-x}}{e^x - e^{-x}} = 2\sigma(2x) - 1,$$

где является сигмоидной функцией. Зеленые круги на схеме представляют поэлементные операции. Справа же на выходе получается новое скрытое состояние (h_t) и новое состояние ячейки (c_t), готовые для следующего ввода. Новое скрытое состояние также используется в качестве вывода, о чем свидетельствует ответвление стрелки вверх.

Давайте поочередно рассмотрим четыре нейронные сети (называемые *фильтрами*) и объясним схему, но сначала обратите внимание, насколько незначительно изменяется состояние ячейки (вверху). Оно даже не проходит непосредственно через нейронную сеть, благодаря чему и переносит более долгосрочное состояние.

Сначала стрелки для ввода и предыдущего скрытого состояния объединяются. В RNN, которую мы написали ранее в этой главе, мы их складывали. В LSTM же мы составляем их в один большой тензор. Это означает, что размерность наших

вложений (представленная размерностью x_t) может отличаться от размерности скрытого состояния. Если мы вызовем эти `n_in` и `n_hid`, стрелка внизу будет иметь размер `n_in + n_hid`. Таким образом, все нейронные сети (оранжевые блоки) оказываются линейными слоями с вводами `n_in + n_hid` и выводами `n_hid`.

Первый фильтр слева называется *фильтром забывания*. Поскольку это линейный слой, сопровождаемый сигмоидой, его вывод будет состоять из скаляров со значениями между 0 и 1. Этот результат мы умножаем на состояние ячейки, чтобы определить, какую информацию сохранить, а какую отбросить. Значения ближе к 0 отклоняются, а значения ближе к 1 сохраняются. Это дает LSTM способность забывать свое долгосрочное состояние. Например, можно ожидать, что при пересечении точки или токена `<eos>` она (обучилась) будет сбрасывать свое состояние ячейки.

Второй фильтр называется *фильтром ввода*. Он работает совместно с третьим фильтром (у которого нет своего имени, но иногда его называют *фильтром ячейки*), обновляя состояние ячейки. Например, мы можем встретить новое местоимение пола, в случае чего понадобится заменить информацию о поле, которую фильтр забывания удалил. Аналогично фильтру забывания фильтр ввода решает, какие элементы состояния ячейки нужно обновить (значения ближе к 1), а какие — нет (значения ближе к 0). Третий фильтр определяет эти обновленные значения в диапазоне от -1 до 1 (на основе функции гиперболического тангенса), после чего полученный результат добавляется в состояние ячейки.

Последним идет *фильтр вывода*. Он определяет, какую информацию из ячейки состояния использовать для генерации вывода. Состояние ячейки проходит через \tanh до того, как совмещается с сигмоидальным выводом из фильтра вывода, в результате чего получается новое скрытое состояние. В коде мы можем написать все это так:

```
class LSTMCell(Module):
    def __init__(self, ni, nh):
        self.forget_gate = nn.Linear(ni + nh, nh)
        self.input_gate = nn.Linear(ni + nh, nh)
        self.cell_gate = nn.Linear(ni + nh, nh)
        self.output_gate = nn.Linear(ni + nh, nh)

    def forward(self, input, state):
        h,c = state
        h = torch.stack([h, input], dim=1)
        forget = torch.sigmoid(self.forget_gate(h))
        c = c * forget
        inp = torch.sigmoid(self.input_gate(h))
        cell = torch.tanh(self.cell_gate(h))
        c = c + inp * cell
        out = torch.sigmoid(self.output_gate(h))
        h = out * torch.tanh(c)
        return h, (h,c)
```

На практике после этого можно произвести рефакторинг. Помимо этого, с позиции производительности лучше будет выполнять одно большое матричное умножение, чем четыре меньших (потому что так мы запускаем специальное ускоренное ядро в GPU только один раз, более эффективно задействуя его параллельность). Составление в стек занимает некоторое время (поскольку нужно переместить один из тензоров в GPU, чтобы сформировать в итоге непрерывный массив), поэтому мы задействуем два раздельных слоя для ввода и скрытого состояния. После оптимизации и рефакторинга код выглядит так:

```
class LSTMCell(Module):
    def __init__(self, ni, nh):
        self.ih = nn.Linear(ni,4*nh)
        self.hh = nn.Linear(nh,4*nh)

    def forward(self, input, state):
        h,c = state
        # Одно большое умножение для всех фильтров лучше, чем четыре маленьких
        gates = (self.ih(input) + self.hh(h)).chunk(4, 1)
        ingate,forgetgate,outgate = map(torch.sigmoid, gates[:3])
        cellgate = gates[3].tanh()

        c = (forgetgate*c) + (ingate*cellgate)
        h = outgate * c.tanh()
        return h, (h,c)
```

Здесь мы используем PyTorch-метод `chunk` для разделения тензора на четыре части. Работает этот метод так:

```
t = torch.arange(0,10); t
tensor([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
t.chunk(2)
(tensor([0, 1, 2, 3, 4]), tensor([5, 6, 7, 8, 9]))
```

А теперь применим эту архитектуру для обучения языковой модели!

Обучение языковой модели с помощью LSTM

Вот та же сеть, что и `LModel15`, но использующая двухслойную LSTM. Обучив ее с повышенной скоростью и за более короткое время, мы улучшим точность:

```
class LModel16(Module):
    def __init__(self, vocab_sz, n_hidden, n_layers):
        self.i_h = nn.Embedding(vocab_sz, n_hidden)
        self.rnn = nn.LSTM(n_hidden, n_hidden, n_layers, batch_first=True)
```

```

        self.h_o = nn.Linear(n_hidden, vocab_sz)
        self.h = [torch.zeros(n_layers, bs, n_hidden) for _ in range(2)]

    def forward(self, x):
        res,h = self.rnn(self.i_h(x), self.h)
        self.h = [h_.detach() for h_ in h]
        return self.h_o(res)

    def reset(self):
        for h in self.h: h.zero_()

learn = Learner(dls, LMModel6(len(vocab), 64, 2),
                loss_func=CrossEntropyLossFlat(),
                metrics=accuracy, cbs=ModelResetter)
learn.fit_one_cycle(15, 1e-2)

```

epoch	train_loss	valid_loss	accuracy	time
0	3.000821	2.663942	0.438314	00:02
1	2.139642	2.184780	0.240479	00:02
2	1.607275	1.812682	0.439779	00:02
3	1.347711	1.830982	0.497477	00:02
4	1.123113	1.937766	0.594401	00:02
5	0.852042	2.012127	0.631592	00:02
6	0.565494	1.312742	0.725749	00:02
7	0.347445	1.297934	0.711263	00:02
8	0.208191	1.441269	0.731201	00:02
9	0.126335	1.569952	0.737305	00:02
10	0.079761	1.427187	0.754150	00:02
11	0.052990	1.494990	0.745117	00:02
12	0.039008	1.393731	0.757894	00:02
13	0.031502	1.373210	0.758464	00:02
14	0.028068	1.368083	0.758464	00:02

Вот теперь результат превосходит многослойную RNN. Хотя здесь по-прежнему наблюдается некоторое переобучение, разобраться с которым нам поможет регуляризация.

Регуляризация LSTM

Как правило, рекуррентные нейронные сети сложно поддаются обучению из-за проблемы исчезающих активаций и градиентов, которую мы недавно рассмотрели. Использование LSTM (или GRU) ячеек упрощает этот процесс, но при этом сохраняется склонность сетей к переобучению. При этом, несмотря на возможность использования, к текстам аугментация данных применяется существенно реже, чем к изображениям, потому что в большинстве случаев для этого нужно, чтобы случайные аугментации генерировала другая модель (например, переводя текст на другой язык и обратно). В целом техники аугментации в отношении текстовых данных пока недостаточно изучены.

Тем не менее для уменьшения переобучения доступны и другие подходы регуляризации, применение которых в LSTM было более тщательно изучено и изложено в работе Стивена Мерити и др. (Stephen Merity) *Regularizing and Optimizing LSTM Language Models* (<https://oreil.ly/Rf-OG>) («Регуляризация и оптимизация языковых моделей LSTM»). В этом исследовании показано, как эффективное использование дропаута, регуляризации активаций и регуляризации временных активаций позволяет LSTM превзойти эталонные результаты, для получения которых ранее требовалась гораздо более сложные модели. LSTM, использующие эти техники, авторы назвали *AWD-LSTM*. Мы рассмотрим каждую из этих техник по очереди.

Dropout

Дропаут (исключение) — это техника регуляризации, представленная Джеком Хинтоном (Geoffrey Hinton) и др. в работе *Improving Neural Networks by Preventing Co-Adaptation of Feature Detectors* (https://oreil.ly/-_xie) («Улучшение нейронных сетей с помощью предотвращения совместной адаптации детекторов признаков»). Основная ее идея — в случайном изменении во время обучения некоторых активаций на ноль. Это гарантирует активное участие всех нейронов в генерации вывода, как показано на рис. 12.10.

При объяснении в ходе интервью своего открытия дропаута Хинтон использовал интересную метафору:

Как-то раз я пошел в банк, обслуживавший мой счет. Там постоянно менялись операционисты, и я поинтересовался у одного из них, в чем причина. Он ответил, что причину не знает, но их постоянно перемещают с места на место. Я решил, что это мера безопасности — чтобы сотрудники не могли провернуть мошеннические схемы. Это натолкнуло меня на мысль, что случайное удаление разного подмножества нейронов в каждом примере ослабило бы их связность и уменьшило переобучение.

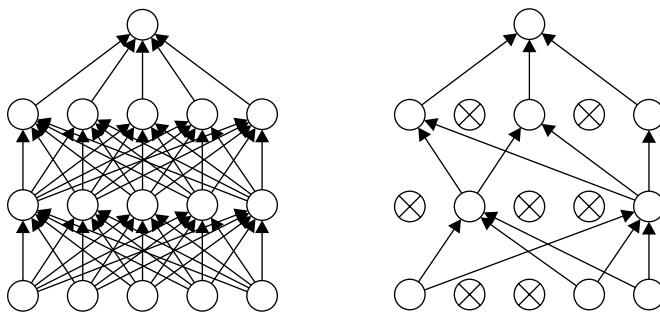


Рис. 12.10. Применение дропаута к нейронной сети (изображение любезно предоставлено Нитишем Шриваставой (Nitish Srivastava)). Слева: стандартная нейронная сеть с двумя скрытыми слоями. Справа: пример прореженной сети, созданной путем применения дропаута к сети слева. Перечеркнутые блоки были исключены

В том же интервью он сказал, что свою лепту в открытие внесла и нейронаука:

По сути, мы не знаем, почему нейроны передают импульсы. Есть теория, что они создают шум с целью регуляризации, так как у нас намного больше параметров, чем точек данных. Смысл дропаута в том, что при наличии шумных активаций вы можете позволить себе использовать гораздо большую модель.

Это объясняет, почему дропаут повышает обобщаемость: сначала он помогает нейронам наладить взаимодействие, а затем делает активации более шумными, повышая таким образом надежность модели.

Тем не менее стоит заметить, что если мы просто обнулим эти активации, не производя никаких других действий, то у модели возникнут сложности с обучением: если мы перейдем от суммы пяти активаций (все они являются положительными, так как мы применяем ReLu) к сумме всего двух, то масштаб уже будет другой. Следовательно, если мы применим дропаут с вероятностью p , то перемасштабируем все активации, разделив их на $1-p$ (в среднем p будет обнулено, поэтому останется $1-p$), как показано на рис. 12.11.

Вот полная реализация слоя дропаута в PyTorch (хотя родной слой PyTorch на самом деле написан на C, а не на Python):

```
class Dropout(Module):
    def __init__(self, p):
        self.p = p
    def forward(self, x):
        if not self.training: return x
        mask = x.new(*x.shape).bernoulli_(1-p)
        return x * mask.div_(1-p)
```

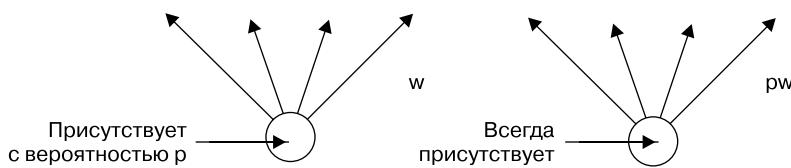


Рис. 12.11. Почему мы масштабируем активации при применении дропаута (изображение любезно предоставлено Нитишем Шриваставой (Nitish Srivastava) и др.). Слева: блок во время обучения, который присутствует с вероятностью p и связан с блоками в следующем слое с весами w . Справа: во время тестирования этот блок присутствует всегда, и веса умножаются на p . Вывод при тестировании такой же, что и ожидаемый вывод при обучении

Метод `bernoulli_` создает тензор случайных нулей (с вероятностью p) и единиц (с вероятностью $1-p$), который затем перед делением на $1-p$ умножается на ввод. Обратите внимание на использование атрибута `training`, который доступен в любом `nn.Module` PyTorch и указывает на то, выполняем мы обучение модели или ее вывод.



ЭКСПЕРИМЕНТИРУЙТЕ САМОСТОЯТЕЛЬНО

В предыдущих главах мы бы привели пример кода для `bernoulli_`, чтобы показать, как он работает. Но поскольку теперь вы уже многое умеете сами, мы будем давать все меньше и меньше примеров, предполагая, что вы будете выполнять собственные эксперименты для выяснения принципа работы тех или иных компонентов. В данном случае в вопроснике в конце главы мы предложим вам поэкспериментировать с `bernoulli_`, но это не означает, что всякий раз нужно ждать нашего призыва к действию, — лучше всего браться за изучение самим.

Использование дропаута до передачи вывода нашей LSTM в последний слой поможет уменьшить переобучение. Как мы увидим в главе 15, дропаут также применяется во многих других моделях, включая исходную вершину CNN, которая используется в `fastai.vision`, и может быть задействован в `fastai.tabular` путем передачи параметра `ps` (где каждый « p » передается в каждый добавленный слой `Dropout`).

Дропаут в режиме обучения и контроля ведет себя по-разному, что мы определяем в `Dropout` с помощью атрибута `training`. Вызов в `Module` метода `train` устанавливает `training` на `True` (как для модуля, в котором вызывается метод, так и для всех других модулей, которые в нем содержатся рекурсивно), а вызов `eval` устанавливает его на `False`. Это выполняется автоматически при вызове методов `Learner`, но если вы не используете этот класс, то не забывайте переключаться между режимами вручную.

Регуляризация активаций и регуляризация временных активаций

Методы *регуляризации активаций* (activation regularization – AR) и *регуляризации временных активаций* (temporal activation regularization – TAR) очень похожи на сокращение весов, о котором мы говорили в главе 8. Применяя сокращение весов, мы добавляем небольшой «штраф» к функции потерь, который способствует максимальному снижению весов. В случае регуляризации активаций мы будем стараться максимально уменьшить не веса, а последние активации, производимые LSTM.

Для регуляризации последних активаций нам нужно их где-то сохранить, а затем добавить средние их квадратов к потерям (вместе с множителем `alpha`, который аналогичен `wd` для сокращения весов):

```
loss += alpha * activations.pow(2).mean()
```

Регуляризация временных активаций связана с тем фактом, что мы прогнозируем токены в предложении. Это значит, что выходные данные наших LSTM при их прочтении по порядку должны иметь смысл. Чтобы это обеспечить, TAR добавляет «штраф» к потерям, максимально уменьшая разницу между двумя последовательными активациями: тензор наших активаций имеет форму `bs x s1 x n_hid`, и мы читаем последовательные активации по оси длины последовательности (измерение в середине). При этом TAR выражается так:

```
loss += beta * (activations[:,1:] - activations[:,:-1]).pow(2).mean()
```

`alpha` и `beta` представляют два настраиваемых гиперпараметра. Чтобы все заработало, наша модель с дропаутом должна возвращать три компонента: подходящий вывод, активации LSTM до дропаута и активации LSTM после дропаута. AR часто применяется к активациям после дропаута (чтобы не применять штраф к активациям, которые мы в итоге обратим в нули), а TAR – к активациям до дропаута (потому что между двумя последовательными временными шагами эти нули создают большие различия). После этого обратный вызов `RNNRegularizer` будет применять данную регуляризацию за нас.

Обучение регуляризованной LSTM со связанными весами

Далее мы совместим дропаут (примененный до перехода к выходному слою) с AR и TAR для обучения предыдущей LSTM. Нужно лишь вернуть три компонента вместо одного: обычный вывод LSTM, активации после дропаута и активации из всех LSTM. Последние две составляющие будут выбраны обратным вызовом `RNNRegularization` в качестве прибавки, которую он должен внести в потери.

В работе, посвященной AWD-LSTM (<https://oreil.ly/ETQ5X>), приводится еще один полезный прием, называемый *связыванием весов*. В языковой модели входные вложения представляют отображение английских слов в активации, а выходной скрытый слой представляет отображение из активаций в английские слова. Логично предположить, что эти отображения могут быть одинаковыми. В PyTorch это можно выразить, присвоив каждому из слоев одинаковую матрицу весов:

```
self.h_o.weight = self.i_h.weight
```

Эти заключительные доработки мы добавляем в `LMMModel7`:

```
class LMMModel7(Module):
    def __init__(self, vocab_sz, n_hidden, n_layers, p):
        self.i_h = nn.Embedding(vocab_sz, n_hidden)
        self.rnn = nn.LSTM(n_hidden, n_hidden, n_layers, batch_first=True)
        self.drop = nn.Dropout(p)
        self.h_o = nn.Linear(n_hidden, vocab_sz)
        self.h_o.weight = self.i_h.weight
        self.h = [torch.zeros(n_layers, bs, n_hidden) for _ in range(2)]

    def forward(self, x):
        raw,h = self.rnn(self.i_h(x), self.h)
        out = self.drop(raw)
        self.h = [h_.detach() for h_ in h]
        return self.h_o(out), raw, out

    def reset(self):
        for h in self.h: h.zero_()
```

Затем создаем регуляризованный `Learner`, используя обратный вызов `RNNRegularizer`:

```
learn = Learner(dls, LMMModel7(len(vocab), 64, 2, 0.5),
                 loss_func=CrossEntropyLossFlat(), metrics=accuracy,
                 cbs=[ModelResetter, RNNRegularizer(alpha=2, beta=1)])
```

`TextLearner` автоматически добавляет два этих обратных вызова (с этими значениями `alpha` и `beta` в качестве предустановленных), поэтому предыдущую строку можно упростить:

```
learn = TextLearner(dls, LMMModel7(len(vocab), 64, 2, 0.4),
                     loss_func=CrossEntropyLossFlat(), metrics=accuracy)
```

После этого мы обучаем модель и добавляем дополнительную регуляризацию, увеличивая сокращение весов до `0.1`:

```
learn.fit_one_cycle(15, 1e-2, wd=0.1)
```

epoch	train_loss	valid_loss	accuracy	time
0	2.693885	2.013484	0.466634	00:02
1	1.685549	1.187310	0.629313	00:02
2	0.973307	0.791398	0.745605	00:02
3	0.555823	0.640412	0.794108	00:02
4	0.351802	0.557247	0.836100	00:02
5	0.244986	0.594977	0.807292	00:02
6	0.192231	0.511690	0.846761	00:02
7	0.162456	0.520370	0.858073	00:02
8	0.142664	0.525918	0.842285	00:02
9	0.128493	0.495029	0.858073	00:02
10	0.117589	0.464236	0.867188	00:02
11	0.109808	0.466550	0.869303	00:02
12	0.104216	0.455151	0.871826	00:02
13	0.100271	0.452659	0.873617	00:02
14	0.098121	0.458372	0.869385	00:02

Вот этот результат уже намного превосходит нашу предыдущую модель!

Резюме

Теперь вы увидели все внутреннее наполнение архитектуры AWD-LSTM, которую мы использовали для классификации текста в главе 10. Дропаут же в ней используется и во многих других местах:

- дропаут вложений (внутри слоя вложений отбрасываются некоторые случайные строки);
- дропаут ввода (после слоя вложений);
- дропаут весов (применяется к весам LSTM на каждом этапе обучения);
- дропаут скрытого состояния (применяется к скрытому состоянию между двумя слоями).

Таким образом, выполняется еще большая регуляризация. Поскольку тонкая настройка этих пяти значений дропаута (включая дропаут перед выходным слоем) очень сложна, мы определили хорошие предустановленные варианты и позволили производить общую регулировку его величины с помощью параметра `drop_multi`, который встречался нам в этой главе (умножаемый на каждый дропаут).

Еще одной мощной архитектурой, особенно эффективной в задачах «последовательность — последовательность» (в которых зависимая переменная сама является последовательностью с длиной переменной, например, в переводе языков), является архитектура «Трансформер». Она описывается в бонусной главе книги на сайте (<https://book.fast.ai/>).

Вопросник

1. Что рекомендуется сделать, если работа с датасетом занимает очень много времени из-за его большого размера и сложности?
2. Почему перед созданием языковой модели мы конкатенируем документы датасета?
3. Какие две доработки модели нужно произвести, чтобы использовать стандартную полносвязную сеть для прогнозирования четвертого слова на основе трех предыдущих?
4. Как в PyTorch использовать общую матрицу для нескольких слоев?
5. Напишите модуль, прогнозирующий третье слово предложения на основе двух предшествующих. Не подглядывайте.
6. Что такое рекуррентная нейронная сеть?
7. Что такое скрытое состояние?
8. Каков эквивалент скрытого состояния в `LMMModel1`?
9. Почему для поддержания состояния в RNN важно передавать модели текст по порядку?
10. Что означает «развернутое» представление RNN?
11. Почему поддержание скрытого состояния в RNN ведет к проблемам с памятью и производительностью? Как это исправить?
12. Что такое ВРТТ?
13. Напишите код для вывода первых нескольких пакетов контрольной выборки, включая преобразование ID токенов в английские строки, как мы показывали для пакетов данных IMDb в главе 10.
14. Что делает обратный вызов `ModelResetter`? Зачем он нужен?

15. Каковы недостатки прогнозирования только одного выходного слова для каждого трех входных?
16. Почему нужна собственная функция потерь для `LMModel14`?
17. Почему обучение `LMModel14` нестабильно?
18. Из развернутого представления видно, что рекуррентная нейронная сеть содержит много слоев. Так почему же нам нужно соединять несколько RNN для улучшения результатов?
19. Нарисуйте представление составной (многослойной) RNN.
20. Почему менее частый вызов `detach` должен приводить к повышению эффективности RNN? Почему на практике этого может не происходить в случае с простой RNN?
21. Почему глубокая нейронная сеть может привести к генерации очень больших или малых активаций? Почему это важно?
22. Какие числа являются наиболее точными в компьютерном представлении чисел с плавающей запятой?
23. Почему исчезание градиентов мешает обучению?
24. Чем помогает наличие двух скрытых состояний в архитектуре LSTM? Какая у каждого из них задача?
25. Как называются эти два состояния в LSTM?
26. Что такое \tanh и как он связан с сигмоидой?
27. В чем задача этого кода в `LSTMCell`:
`h = torch.stack([h, input], dim=1)`
28. За что в PyTorch отвечает `chunk`?
29. Внимательно изучите версию `LSTMCell11` после рефакторинга, убедившись, что понимаете, как и почему она делает то же, что и версия до рефакторинга.
30. Почему для `LMModel16` можно использовать более высокую скорость обучения?
31. Какие три техники регуляризации применяются в модели AWD-LSTM?
32. Что такое дропаут?
33. Почему при использовании дропаута мы масштабируем активации? Когда это делается: при обучении модели, ее выводе или в обоих случаях?
34. Какова задача этой строки из `Dropout`:
`if not self.training: return x`
35. Поэкспериментируйте с `bernoulli_`, чтобы понять, как он работает.
36. Как в PyTorch перевести модель в режим обучения? А в режим оценки?

37. Напишите уравнение для регуляризации активаций (в коде или в математическом выражении). Чем оно отличается от сокращения весов?
38. Напишите уравнение для регуляризации временных активаций (в коде или в математическом выражении). Почему мы бы не стали использовать его для задач компьютерного зрения?
39. Что такое связывание весов в языковой модели?

Дополнительные задания

1. Почему `forward` в `LModel12` может начинаться с `h=0`? Почему нам не нужно сообщать `h=torch.zeros(...)`?
2. Напишите код для LSTM с нуля (можете обратиться к рис. 12.9).
3. Поиските в Сети информацию об архитектуре GRU и реализуйте ее с нуля, после чего попробуйте обучить модель. Посмотрите, удастся ли вам получить результаты, близкие к результатам этой главы. Сравните ваши успехи с результатами модуля GRU, встроенного в PyTorch.
4. Рассмотрите исходный код для AWD-LSTM в `fastai` и попытайтесь сопоставить каждую строку кода с понятиями из этой главы.

ГЛАВА 13

Сверточные нейронные сети

В главе 4 мы научились создавать нейронную сеть для распознавания изображений. При распознавании с ее помощью троек и семерок мы достигли точности 98 %, но при этом мы также видели, что встроенные классы fastai способны приблизиться и к 100 %. Постараемся сократить этот разрыв.

В этой главе мы начнем с углубления в понятие сверток и построения CNN (сверточной нейронной сети) с нуля. Затем перейдем к изучению техник по повышению стабильности обучения и узнаем все доработки, которые библиотека обычно выполняет за нас для получения высоких результатов.

Магия сверток

Один из наиболее мощных инструментов специалистов ML — это *инженерия признаков*. *Признак* представляет собой преобразование данных, упрощающее их восприятие моделью. Например, функция `add_datepart`, которую мы использовали для препроцессинга табличного датасета в главе 9, добавляла признаки дат в набор Bulldozers. А какие типы признаков мы можем получить из изображений?



ТЕРМИН: ИНЖЕНЕРИЯ ПРИЗНАКОВ

Создание новых преобразований входных данных для облегчения их восприятия моделью.

В случае с изображением признак, как правило, выражает его отличительный атрибут. Например, цифра 7 определяется двумя контурами: горизонтальным в верхней части и диагональным, идущим от его правого края вниз влево. С цифрой 3 уже сложнее: она определяется однонаправленными диагональными контурами в верхней левой и нижней правой частях, противоположными им контурами в левой нижней и правой верхней частях, а также горизонтальными

контурами, расположеными вверху, посередине, внизу и т. д. Здесь и кроется смысл: вместо анализа пикселей извлекать информацию о том, где именно на изображении располагаются контуры, и рассматривать эту информацию как признаки.

На деле поиск контуров на изображении является вполне стандартной задачей компьютерного зрения и сложности не представляет. Для этого используется так называемая *свертка*. Для реализации свертки требуется только умножение и сложение — две операции, отвечающие за большую часть работы, которую мы увидим в каждой модели глубокого обучения в данной книге.

Свертка применяется к изображению *ядру*. Ядро — это небольшая матрица, пример которой приведен на рис. 13.1, где она выделена красным (размер 3×3).

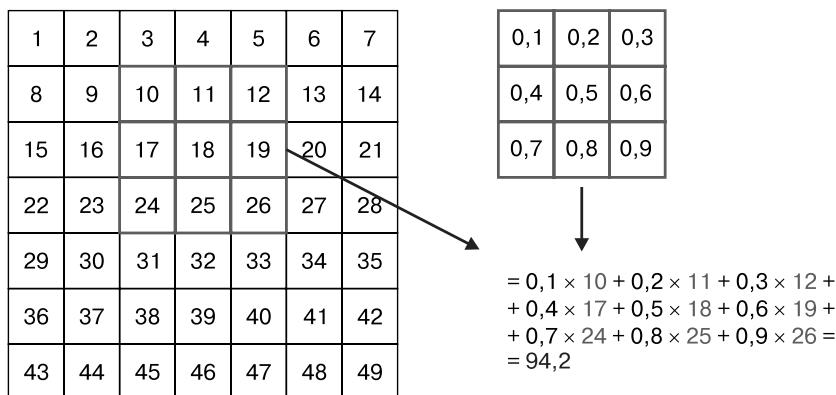


Рис. 13.1. Применение ядра к одному местоположению

Наблюдаемая сетка 7×7 представляет *изображение*, к которому мы будем применять ядро. Операция свертки умножает каждый элемент ядра на соответствующий элемент блока 3×3 изображения. После этого полученные результаты складываются. На рис. 13.1 показан пример применения ядра к одному участку изображения: блоку 3×3 вокруг ячейки 18.

Давайте выразим это кодом и начнем с создания небольшой матрицы 3×3 :

```
top_edge = tensor([[-1, -1, -1],
                   [ 0,  0,  0],
                   [ 1,  1,  1]]).float()
```

Это мы будем называть *ядром* (потому что его называют крутыми исследователи компьютерного зрения). При этом нам, естественно, понадобится изображение:

```
path = untar_data(URLs.MNIST_SAMPLE)

im3 = Image.open(path/'train'/'3'/'12.png')
show_image(im3);
```



Теперь мы возьмем на изображении верхний квадрат размером 3×3 пикселя и умножим значение каждого из этих пикселей на соответствующий ему элемент ядра, после чего сложим результаты так, как показано ниже:

```
im3_t = tensor(im3)
im3_t[0:3,0:3] * top_edge

tensor([[-0., -0., -0.],
       [0., 0., 0.],
       [0., 0., 0.]])
(im3_t[0:3,0:3] * top_edge).sum()

tensor(0.)
```

Пока ничего толкового — все пиксели в верхнем левом углу белые. Выберем пару более интересных участков:

```
df = pd.DataFrame(im3_t[:10,:20])
df.style.set_properties(**{'font-size':'6pt'}).background_gradient('Greys')
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
0	0	0	0	0	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	0	0	0	0	0	
2	0	0	0	0	0	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	0	0	0	0	0	
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
5	0	0	0	12	99	91	142	155	246	182	155	155	155	155	131	52	0	0	0	
6	0	0	0	138	254	254	254	254	254	254	254	254	254	254	252	210	122	33	0	
7	0	0	0	220	254	254	254	235	189	189	189	189	150	189	205	254	254	254	75	
8	0	0	0	35	74	35	35	25	0	0	0	0	0	0	13	224	254	254	153	
9	0	0	0	0	0	0	0	0	0	0	0	0	0	0	90	254	254	247	53	

В верхнем контуре есть ячейка 5, 7. Повторим вычисления для этой области:

```
(im3_t[4:7,6:9] * top_edge).sum()
tensor(762.)
```

А в правом сегменте есть ячейка 8, 18. Что мы получим здесь?

```
(im3_t[7:10,17:20] * top_edge).sum()
tensor(-29.)
```

Как вы видите, наши небольшие вычисления возвращают максимальное значение там, где квадрат 3×3 пикселя представляет верхний контур (то есть там, где в верхней части квадрата низкие значения, а сразу под ними высокие). Причина в том, что значения -1 в ядре оказывают в данном случае мало влияния в отличие от значений 1 .

Взглянем с позиции математики. Фильтр получает окно изображения размером 3×3 , и если мы укажем значения пикселей так:

$$\begin{matrix} a1 & a2 & a3 \\ a4 & a5 & a6 \\ a7 & a8 & a9 \end{matrix}$$

тогда фильтр вернет $a1 + a2 + a3 - a7 - a8 - a9$. Если при этом мы будем находиться в части изображения, где сумма $a1, a2$ и $a3$ совпадает с суммой $a7, a8$ и $a9$, то они друг друга уничтожат и мы получим 0. Если же $a1$ будет больше $a7, a2$ больше $a8$ и $a3$ больше $a9$, то и результат в итоге получится больше. Так фильтр обнаруживает горизонтальные края — говоря точнее, края, где наблюдается переход от светлых частей изображения вверху к более темным частям внизу.

Если изменить наш фильтр, указав в нем ряд 1 вверху и ряд -1 внизу, то он будет обнаруживать горизонтальные контуры перехода от темного к светлому. Поместив 1 и -1 в столбцы, а не ряды, мы получим фильтр, обнаруживающий вертикальные контуры. Каждый набор весов будет производить разный результат.

Создадим функцию, применяющую это к одному местоположению, и сверим результат с полученным ранее:

```
def apply_kernel(row, col, kernel):
    return (im3_t[row-1:row+2,col-1:col+2] * kernel).sum()

apply_kernel(5,7,top_edge)
tensor(762.)
```

Обратите внимание, что мы не можем применить ее к самому углу (например, участку $0,0$), поскольку там нет полноценного квадрата 3×3 .

Отображение ядра свертки

Отобразим `apply_kernel()` по всей сетке координат. Для этого мы возьмем ядро 3×3 и применим его к каждому участку 3×3 изображения. В качестве примера на рис. 13.2 показаны положения, в которых данное ядро может быть применено к первому ряду изображения 5×5 .

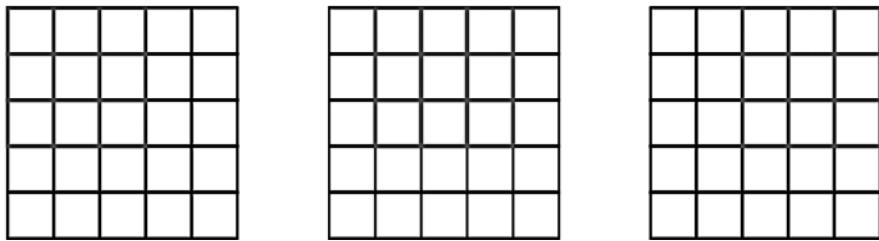


Рис. 13.2. Применение ядра ко всей сетке

Для получения сетки координат мы используем *генератор вложенных списков*:

```
[[[i,j) for j in range(1,5)] for i in range(1,5)]
[[[(1, 1), (1, 2), (1, 3), (1, 4)],
 [(2, 1), (2, 2), (2, 3), (2, 4)],
 [(3, 1), (3, 2), (3, 3), (3, 4)],
 [(4, 1), (4, 2), (4, 3), (4, 4)]]]
```



ГЕНЕРАТОРЫ ВЛОЖЕННЫХ СПИСКОВ

В Python такие генераторы используются часто, и если вам они раньше не встречались, уделите немного времени их изучению, чтобы понять описываемый здесь процесс. Рекомендуем поэкспериментировать с написанием собственных генераторов вложенных списков.

Вот результат применения ядра к сетке координат:

```
rng = range(1,27)
top_edge3 = tensor([[apply_kernel(i,j,top_edge) for j in rng] for i in rng])
show_image(top_edge3);
```



Выглядит неплохо! Верхние края темные, а нижние — светлые (поскольку *противоположны* верхним). Теперь, когда изображение содержит и отрицательные

числа, `matplotlib` автоматически изменила цвета так, что белый представляется малыми числами, черный — большими, а нули отображают серый фон.

То же самое мы сделаем с левыми контурами:

```
left_edge = tensor([[[-1,1,0],
                     [-1,1,0],
                     [-1,1,0]]]).float()

left_edge3 = tensor([[apply_kernel(i,j,left_edge) for j in rng] for i in rng])

show_image(left_edge3);
```



Как мы уже говорили, свертка — это операция применения ядра ко всей сетке. В работе Винсента Дюмулена (Vincent Dumoulin) и Франческо Визина (Francesco Visin) *A Guide to Convolution Arithmetic for Deep Learning* (<https://oreil.ly/les1R>) («Руководство по сверточной арифметике для глубокого обучения») приводится множество схем, показывающих применение ядер. На рис. 13.3 показан пример из этой работы, где к светло-синему изображению применяется темно-синее ядро 3×3 , в результате чего образуется зеленая карта активаций размером 2×2 .

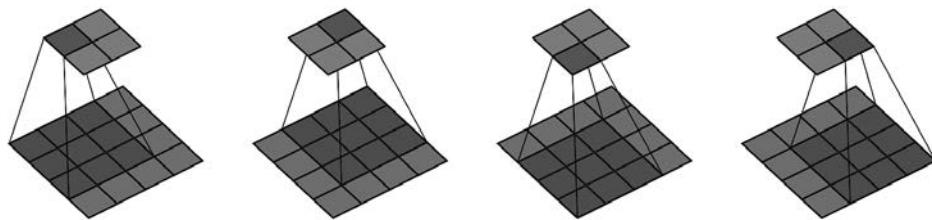


Рис. 13.3. Результат применения ядра 3×3 к изображению 4×4 (информация любезно предоставлена Винсентом Дюмуленом и Франческо Визином)

Взгляните на получившуюся в результате форму. Если исходное изображение имеет высоту h и ширину w , сколько окон 3×3 мы обнаружим? Из примера видно, что окон получается $h-2$ на $w-2$, в результате чего итоговое изображение имеет высоту $h-2$ и ширину $w-2$.

Мы не будем реализовывать эту функцию свертки с нуля, а используем готовую реализацию Python, так как это намного быстрее всего, что мы могли бы сделать в нем сами.

Свертки в PyTorch

Свертки относятся к очень важным и широко используемым операциям, поэтому в PyTorch они встроены по умолчанию и называются `F.conv2d` (напомним, что `F` означает импорт `fastai` из `torch.nn.functional`, как рекомендовано PyTorch). Согласно документации PyTorch, в ней используются следующие параметры:

- `input` — входной тензор формы (`minibatch, in_channels, iH, iW`);
- `weight` — фильтры формы (`out_channels, in_channels, kH, kW`).

Здесь `iH, iW` — это высота и ширина изображения (то есть 28, 28), а `kH, kW` — это высота и ширина ядра (3, 3). Очевидно, что PyTorch ожидает для обоих этих аргументов тензоры ранга 4, но у нас сейчас тензоры ранга 2 (то есть матрицы или массивы с двумя осями).

Необходимость в двух дополнительных осях обусловлена рядом специальных приемов, используемых PyTorch. Первый — это возможность PyTorch применять свертку к нескольким изображениям одновременно, что позволяет вызвать ее для каждого элемента пакета сразу.

Второй прием заключается в способности одновременно использовать несколько ядер. Дополнительно создадим ядра для диагональных контуров, после чего составим все четыре в один тензор:

```
diag1_edge = tensor([[ 0,-1,  1],
                     [-1,  1,  0],
                     [ 1,  0,  0]]).float()
diag2_edge = tensor([[ 1,-1,  0],
                     [ 0,  1,-1],
                     [ 0,  0,  1]]).float()

edge_kernels = torch.stack([left_edge, top_edge, diag1_edge, diag2_edge])
edge_kernels.shape

torch.Size([4, 3, 3])
```

Для проверки нам потребуется `DataLoader` и образец мини-пакета. Задействуем API блока данных:

```
mnist = DataBlock((ImageBlock(cls=PILImageBW), CategoryBlock),
                  get_items=get_image_files,
                  splitter=GrandparentSplitter(),
                  get_y=parent_label)

dls = mnist.dataloaders(path)
xb,yb = first(dls.valid)
xb.shape

torch.Size([64, 1, 28, 28])
```

По умолчанию `fastai` при использовании блоков данных помещает данные в GPU. Мы же для наших примеров переместим их в CPU:

```
xb,yb = to_cpu(xb),to_cpu(yb)
```

В одном пакете находится 64 изображения, каждое состоит из одного канала размером 28×28 пикселей. `F.conv2d` способна обрабатывать и многоканальные (цветные) изображения. Канал представляет один базовый цвет — в стандартных полноцветных изображениях используется три канала: красный, зеленый и синий. PyTorch представляет изображение как тензор ранга 3, используя следующие измерения:

[каналы, строки, столбцы]

Обработку нескольких каналов мы рассмотрим в этой главе чуть позже. Передаваемые в `F.conv2d` ядра должны быть тензорами ранга 4:

[признаки_вых, каналы_вх, строки, столбцы]

В `edge_kernels` сейчас не хватает одного из измерений: нам нужно сообщить PyTorch, что число входных каналов ядра равно единице, для чего мы вставим ось размером 1 (называемую *единичной осью*) в первую часть, где согласно документации PyTorch ожидается `in_channels`.

```
edge_kernels.shape,edge_kernels.unsqueeze(1).shape (torch.Size([4, 3, 3]),
torch.Size([4, 1, 3, 3]))
```

Так мы корректируем форму `edge_kernels`. Теперь передадим все это в `conv2d`:

```
edge_kernels = edge_kernels.unsqueeze(1)
batch_features = F.conv2d(xb, edge_kernels)
batch_features.shape

torch.Size([64, 4, 26, 26])
```

Выходная форма показывает, что в мини-пакете у нас 64 изображения, 4 ядра и 26×26 карт контуров (мы начали с 28×28 изображений, но утратили по одному пикслю с каждого края, о чем говорилось ранее). Получаем мы те же результаты, что и в случае, когда делали это вручную:

```
show_image(batch_features[0,0]);
```



Самый важный из спецприемов PyTorch заключается в том, что он может заставить GPU, чтобы проделать всю эту работу в параллельном режиме,

а именно применить несколько ядер к множеству изображений по нескольким каналам. Подобное задействование параллельной обработки GPU является решающим для его эффективного использования. Если делать все это по-очередно, то зачастую процесс в целом будет замедляться в сотни раз (а если задействовать ручной цикл свертки из предыдущего раздела, то работа станет в миллионы раз медленнее!). Поэтому чтобы стать успешным специалистом глубокого обучения, необходимо научиться единовременно загружать GPU большим объемом работы.

Было бы здорово не терять те самые два пикселя на каждой оси. Для этого мы добавляем так называемое *заполнение* (padding), в качестве которого вдоль внешней стороны изображения размещаются дополнительные пиксели. Чаще всего эти пиксели представляют нули.

Штрихи и заполнение

Применив заполнение, мы можем гарантировать, что карта выходных активаций будет иметь те же размеры, что и исходное изображение. Это существенно упростит процесс построения архитектур. На рис. 13.4 показано, как заполнение позволяет применить ядро в углах изображения, чего ранее мы сделать не могли.

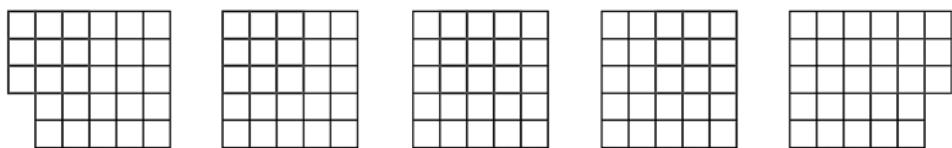


Рис. 13.4. Свертка с заполнением

На рис. 13.5 показано, что при входном изображении 5×5 , ядре 4×4 и заполнении шириной два пикселя мы получаем карту активаций 6×6 .

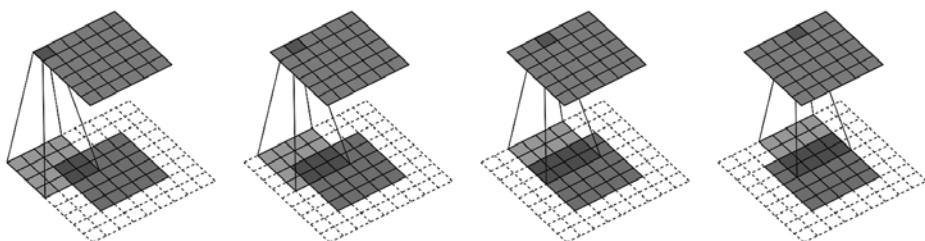


Рис. 13.5. Ядро 4×4 с входным изображением 5×5 и заполнением 2 пикселя (изображение любезно предоставлено Винсентом Дюмулином и Франческо Визином)

Если мы добавляем ядро размером ks на ks (где ks является нечетным числом), то необходимое для сохранения формы заполнение каждой стороны будет равно $ks//2$. В случае четного значения ks для верха/низа и левой/правой стороны потребуется разная величина заполнения, но на практике четный размер фильтра мы используем крайне редко.

До сих пор, применяя ядро к сетке, мы перемещались на один пиксель за шаг, но можно перемещаться и дальше. Например, после каждого применения ядра мы можем шагать на два пикселя, как показано на рис. 13.6. Это называется сверткой с *шагом 2*. На практике чаще всего используется ядро размером 3×3 и заполнение единицами. Как вы позже увидите, свертки с шагом 2 пригождаются для уменьшения размера выходов, а свертки с шагом 1 — для добавления слоев без изменения выходного размера.

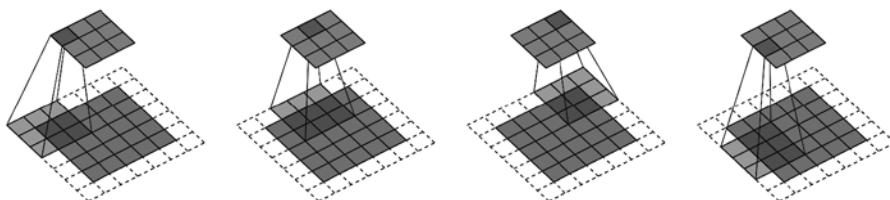


Рис. 13.6. Ядро 3×3 с входом 5×5 , сверткой с шагом 2 и заполнением 1 пиксель
(изображение любезно предоставлено Винсентом Дюмулином и Франческо Визином)

В изображении размеров h на w при использовании заполнения 1 и шага 2 в результате мы получим размер $(h+1)//2$ на $(w+1)//2$. Общая формула для каждого из измерений следующая:

$$(n + 2*pad - ks) // stride + 1$$

Здесь pad означает заполнение, ks — размер ядра, а $stride$ — размер шага.

Теперь рассмотрим процесс вычисления значений пикселей наших сверток.

Понимание сверточных уравнений

Для объяснения математики сверток студент fast.ai Мэтт Кляйнсмит (Matt Kleinsmith) предложил очень крутую идею показывать CNN с разных точек зрения. Эта его идея оказалась настолько полезна, что мы используем ее и здесь.

Вот наше пиксельное изображение 3×3 , где пиксели обозначены буквами:

A	B	C
D	E	F
G	H	J

А вот ядро, где каждый вес тоже обозначен буквой, но уже греческой:

α	β
γ	δ

Поскольку фильтр сопоставляется с изображением четыре раза, мы получаем четыре результата:

P	Q
R	S

На рис. 13.7 показано, как мы применили ядро к каждому участку изображения для получения каждого значения результата.

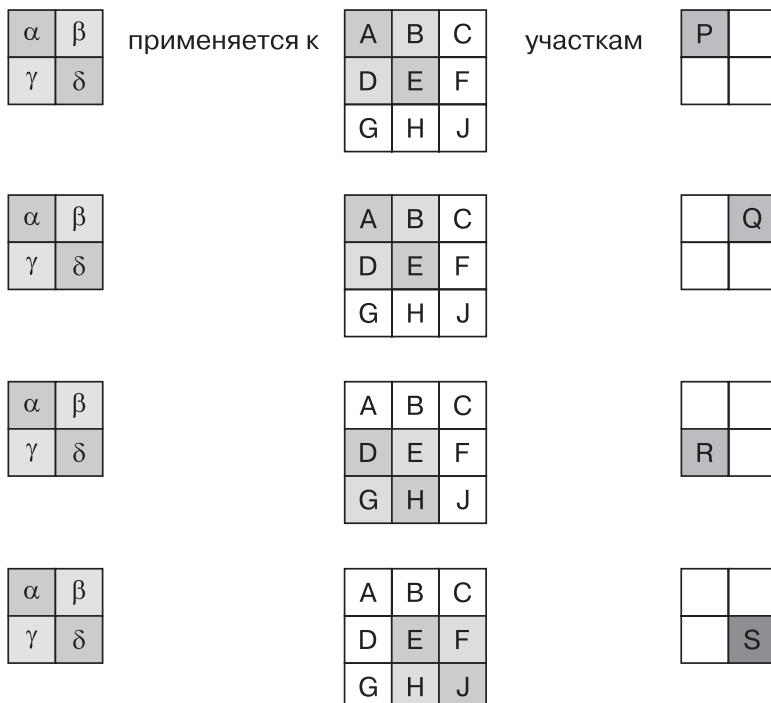


Рис. 13.7. Применение ядра

А вот и само уравнение:

$$\begin{array}{l} \alpha * A + \beta * B + \gamma * D + \delta * E + b = P \\ \alpha * B + \beta * C + \gamma * E + \delta * F + b = Q \\ \alpha * D + \beta * E + \gamma * G + \delta * H + b = R \\ \alpha * E + \beta * B + \gamma * H + \delta * J + b = S \end{array}$$

Рис. 13.8. Уравнение

Обратите внимание, что выражение смещения, b , одинаково для каждого участка изображения. Теперь смещение стоит рассматривать как часть фильтра, так же как и веса ($\alpha, \beta, \gamma, \delta$).

А вот интересное открытие: свертку можно представить как особый вид матричного умножения, что показано на рис. 13.9. Здесь матрица весов аналогична весовой, которая используется в обычных нейронных сетях, но в данном случае у нее есть два особых свойства.

1. Нули, выделенные серым, не обучаются. Это означает, что они останутся нулями в течение всего процесса оптимизации.
2. Некоторые веса равны, и, несмотря на свою обучаемость (то есть способность к изменению), они должны оставаться равными и называются *общими весами*.

Нули соответствуют пикселям, которые фильтр не затрагивает. Каждая строка матрицы весов отражает одно применение фильтра.

Теперь, когда мы понимаем, что такое свертки, построим с их помощью нейронную сеть.

$$\begin{array}{c|ccccccccc} \alpha & \beta & 0 & \gamma & \delta & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & \alpha & \beta & 0 & \gamma & \delta & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \alpha & \beta & 0 & \gamma & \delta & 0 & 0 \\ 0 & 0 & 0 & 0 & \alpha & \beta & 0 & \gamma & \delta & 0 \end{array} * \begin{array}{c|c} A & b \\ \hline B & b \\ C & b \\ D & b \\ E & b \\ F & b \\ G & b \\ H & b \\ J & b \end{array} = \begin{array}{c|c} \alpha A + \beta B + 0C + \gamma D + \delta E + 0F + 0G + 0H + 0J + b & P \\ \hline 0A + \alpha B + \beta C + 0D + \gamma E + \delta F + 0G + 0H + 0J + b & Q \\ 0A + 0B + 0C + \alpha D + \beta E + 0F + \gamma G + \delta H + 0J + b & R \\ 0A + 0B + 0C + 0D + \alpha E + \beta F + 0G + \gamma H + \delta J + b & S \end{array}$$

Рис. 13.9. Свертка как матричное умножение

Первая сверточная нейронная сеть

Нет оснований полагать, что наилучшими ядрами для распознавания изображений окажутся какие-то определенные краевые фильтры. Более того, мы видели, что в более поздних слоях ядра сверток становятся сложными преобразованиями признаков, полученных из более ранних слоев, но для построения их вручную хорошего способа у нас нет.

Вместо этого будет лучше обучить значения ядер, как мы это уже умеем: с помощью SGD. В результате модель изучит полезные для классификации признаки. При использовании свертки вместо (или в дополнение к ним) стандартных линейных слоев мы создаем *сверточную нейронную сеть* (convolutional neural network – CNN).

Создание CNN

Вернемся к простой нейронной сети из главы 4. Вот ее определение:

```
simple_net = nn.Sequential(
    nn.Linear(28*28,30),
    nn.ReLU(),
    nn.Linear(30,1)
)
```

А вот определение модели:

```
simple_net
Sequential(
  (0): Linear(in_features=784, out_features=30, bias=True)
  (1): ReLU()
  (2): Linear(in_features=30, out_features=1, bias=True)
)
```

Теперь нам нужно создать аналогичную этой линейной модели архитектуру, но вместо линейных слоев задействовать сверточные. `nn.Conv2d` – это модульный эквивалент `F.conv2d`, который более удобен при создании архитектуры, так как генерирует матрицу весов автоматически при инстанцировании.

Вот вариант архитектуры:

```
broken_cnn = sequential(
    nn.Conv2d(1,30, kernel_size=3, padding=1),
    nn.ReLU(),
    nn.Conv2d(30,1, kernel_size=3, padding=1)
)
```

Обратите внимание, что нам не потребовалось указывать входной размер `28*28`. Причина в том, что линейному слою для каждого пикселя нужен вес из матрицы

весов, для чего ему необходимо знать общее число пикселей, а свертка применяется для каждого пикселя автоматически. Как мы видели в предыдущем разделе, веса зависят только от числа входных и выходных каналов, а также размера ядра.

Подумайте, какой должна быть форма выходного сигнала, а затем давайте проверим ваши предположения:

```
broken_cnn(xb).shape
torch.Size([64, 1, 28, 28])
```

Это мы не сможем использовать для классификации, поскольку нам нужна одна выходная активация для каждого изображения, а не карта активаций 28×28 . Решить это можно, задействовав достаточное число сверток с шагом 2 так, чтобы размер финального слоя получился равным 1. После одной свертки с шагом 2 размер станет 14×14 , после двух — 7×7 , затем 4×4 , 2×2 и в конце 1.

Давайте пробовать. Начнем с определения функции с базовыми параметрами, которые будем использовать в каждой свертке:

```
def conv(ni, nf, ks=3, act=True):
    res = nn.Conv2d(ni, nf, stride=2, kernel_size=ks, padding=ks//2)
    if act: res = nn.Sequential(res, nn.ReLU())
    return res
```

РЕФАКТОРИНГ



Рефакторинг элементов нейронной сети, подобных этому, существенно снижает риск получения ошибок из-за несогласованности внутри архитектур. Помимо этого, он делает более очевидным, какие части слоев фактически меняются.

Используя свертку с шагом 2, мы обычно увеличиваем количество признаков, потому что число активаций в карте активаций уменьшается в четыре раза. При этом слишком сильно уменьшать емкость слоя за один раз нам не нужно.

ТЕРМИН: КАНАЛЫ И ПРИЗНАКИ



Эти выражения широко используются одно вместо другого и означают размер второй оси матрицы весов, которая представляет количество активаций в каждой ячейке сетки после свертки. Термин *признаки* никогда не применяется к входным данным, в то время как *каналы* могут означать либо входные данные (обычно каналы — это цвета), либо активации внутри сети.

Вот пример построения простой CNN:

```
simple_cnn = sequential(
    conv(1, 4), # 14 x 14
    conv(4, 8), # 7 x 7
```

```

    conv(8 ,16),           # 4 × 4
    conv(16,32),          # 2 × 2
    conv(32,2, act=False), # 1 × 1
    Flatten(),
)

```



СЛОВО ДЖЕРЕМИ

Я люблю добавлять комментарии вроде тех, что приведены в этом коде после каждой свертки, чтобы показать размер карты признаков после каждого слоя. Приведенные здесь комментарии предполагают, что размер входного изображения был 28×28 .

Теперь на выходе сети получается две активации, которые сопоставляются с двумя возможными уровнями из наших меток:

```

simple_cnn(xb).shape
torch.Size([64, 2])

```

Далее можно создать Learner:

```
learn = Learner(dls, simple_cnn, loss_func=F.cross_entropy, metrics=accuracy)
```

Для просмотра подробностей происходящего в модели процесса используется метод `summary`:

```

learn.summary()
Sequential (Input shape: ['64 x 1 x 28 x 28'])
=====
Layer (type)      Output Shape       Param #  Trainable
=====
Conv2d           64 x 4 x 14 x 14   40        True
ReLU            64 x 4 x 14 x 14   0         False
Conv2d           64 x 8 x 7 x 7    296       True
ReLU            64 x 8 x 7 x 7    0         False
Conv2d           64 x 16 x 4 x 4   1,168     True
ReLU            64 x 16 x 4 x 4   0         False
Conv2d           64 x 32 x 2 x 2   4,640     True
ReLU            64 x 32 x 2 x 2   0         False
Conv2d           64 x 2 x 1 x 1    578       True
Flatten          64 x 2

```

```
Total params: 6,722
Total trainable params: 6,722
Total non-trainable params: 0

Optimizer used: <function Adam at 0x7fb9c258cb0>
Loss function: <function cross_entropy at 0x7fbca9ba0170>

Callbacks:
    TrainEvalCallback
    Recorder
    ProgressCallback
```

Заметьте, что выход последнего слоя `Conv2d` имеет форму $64 \times 2 \times 1 \times 1$. Нам нужно удалить эти лишние оси 1×1 , что делает `Flatten`. По сути, он аналогичен PyTorch-методу `squeeze`, но в виде модуля.

Посмотрим, как пойдет обучение. Поскольку это первая углубленная сеть, построенная нами с нуля, мы проведем его с меньшей скоростью и в течение большего числа эпох:

```
learn.fit_one_cycle(2, 0.01)
```

epoch	train_loss	valid_loss	accuracy	time
0	0.072684	0.045110	0.990186	00:05
1	0.022580	0.030775	0.990186	00:05

Успех! Мы приблизились к ранее полученному результату `resnet18`, хоть это и не совсем то, что нужно, к тому же потребовалось большее число эпох и пониженная скорость обучения. Нам еще предстоит изучить несколько приемов, но мы уже почти готовы к созданию современной CNN с нуля.

Разъяснение арифметики сверток

Из сводки видно, что входные данные имеют размер $64 \times 1 \times 28 \times 28$, обозначающий оси `batch`, `channel`, `height`, `width`. Обычно они обозначаются `NCHW` (где N – это размер пакета). Для сравнения TensorFlow, например, упорядочивает оси как `NHWC`. Вот первый слой:

```
m = learn.model[0]
m

Sequential(
  (0): Conv2d(1, 4, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1))
  (1): ReLU()
)
```

Итак, здесь у нас один входной канал, четыре выходных и ядро 3×3 . Проверим веса первой свертки:

```
m[0].weight.shape
torch.Size([4, 1, 3, 3])
```

Сводка показывает, что у нас 40 параметров, но $4*1*3*3$ будет 36, откуда же еще 4 параметра? Взглянем на содержимое смещения:

```
m[0].bias.shape
torch.Size([4])
```

Теперь с помощью этой информации можно прояснить сказанное в предыдущем разделе: «Используя свертку с шагом, 2 мы зачастую увеличиваем количество признаков, потому что уменьшаем число активаций в карте активаций в 4 раза. При этом мы не хотим слишком сильно уменьшать емкость слоя за один раз».

Для каждого канала есть одно смещение. (Иногда каналы называются *признаками*, или *фильтрами*, но только когда они не являются входными.) На выходе получается форма $64 \times 4 \times 14 \times 14$. Она и станет входной формой для следующего слоя, который, согласно *summary*, содержит 296 параметров. Проигнорируем ось размера пакета, чтобы не усложнять. Итак, каждую из $14 \times 14 = 196$ точек мы умножаем на $296 - 8 = 288$ весов (для простоты игнорируя смещение), получая на данном слое $196 * 288 = 56_448$ умножений. А на следующем слое будет $7 * 7 * (1168 - 16) = 56_448$ умножений.

Здесь *размер сетки* нашей свертки с шагом 2 был уменьшен вдвое с 14×14 в 7×7 , и при этом мы удвоили *количество фильтров* с 8 до 16, что привело к изменению общего количества вычислений. Если оставлять число каналов в каждом слое с шагом 2 без изменений, то по мере углубления сети объем производимых ею вычислений будет постепенно сокращаться. Но нам известно, что более глубокие слои должны вычислять семантически богатые признаки (например, глаза или мех), поэтому вряд ли уменьшение вычислений будет иметь смысла.

По-другому на это можно посмотреть через призму рецептивных полей.

Рецептивные поля

Рецептивное поле — это область изображения, задействованная в вычислении слоя. На сайте книги (<https://book.fast.ai/>) есть электронная таблица Excel под названием *conv-example.xlsx*, которая демонстрирует вычисления двух сверточных слоев с шагом 2 на примере цифры из MNIST. У каждого слоя есть одно ядро. На рис. 13.10 показано, что мы видим, когда щелкаем на одной из ячеек в разделе *conv2*, отражающей выход второго сверточного слоя, и затем выбираем *Отследить прецеденты*.

Ячейка с зеленой рамкой — это та, по которой мы щелкнули, а обведенные синим — это ее *прецеденты* — ячейки, используемые для вычисления ее значения.

Они соответствуют области ячеек 3×3 из входного слоя слева и ячейкам из фильтра справа. Теперь давайте еще раз щелкнем на *Отследить прецеденты*, чтобы увидеть, какие ячейки задействуются для вычисления этих выходных значений. На рис. 13.11 показано, что именно происходит.

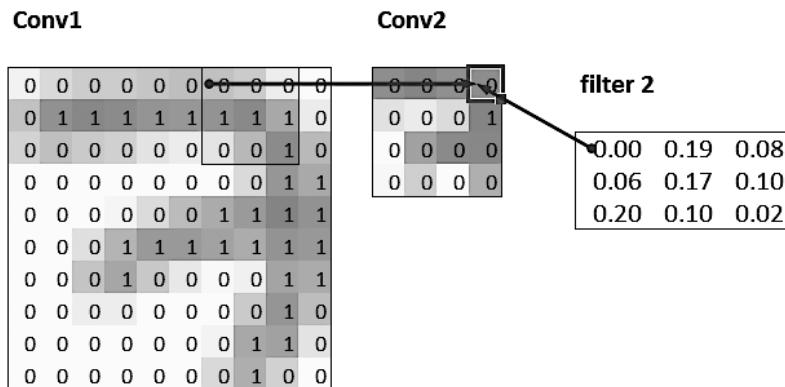


Рис. 13.10. Ближайшие прецеденты слоя Conv2

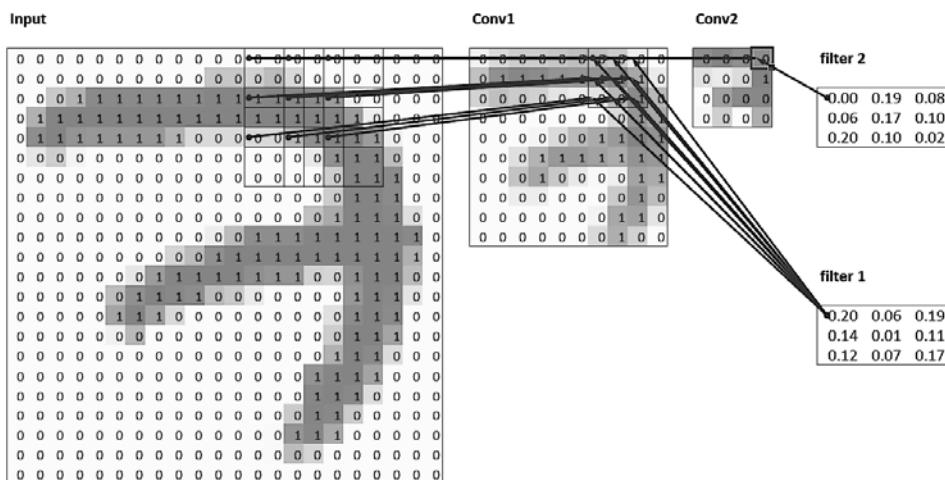


Рис. 13.11. Вторичные прецеденты слоя Conv2

В этом примере у нас всего два сверточных слоя, каждый с шагом 2, и мы прослеживаем по ним прецеденты обратно до входного изображения. Его область размером 7×7 используется для вычисления одной зеленою ячейки в слое Conv2 и называется *рецептивным полем* входа обведенной зеленым активации

в Conv2. Помимо этого, теперь нам требуется второе ядро фильтра, поскольку сейчас у нас два слоя.

Из этого примера видно, что чем глубже мы уходим в сеть (а именно, чем больше позади слоя остается сверток с шагом 2), тем больше становится рецептивное поле активации в этом слое. Большое рецептивное поле означает, что для вычисления каждой активации в этом слое задействуется большая часть входного изображения. Теперь мы знаем, что появляющиеся в углубленных слоях сети семантически богатые признаки соответствуют более крупным рецептивным полям. Следовательно, можно ожидать, что для обработки такой возрастающей сложности нам потребуется больше весов для каждого из признаков. Это еще один способ выразить то, о чём мы говорили в предыдущем разделе: при внедрении в сети свертки с шагом 2 нам также нужно увеличивать количество каналов.

В процессе написания текущей главы у нас было много непроясненных вопросов, на которые требовалось ответить, чтобы как можно лучше объяснить вам особенности CNN. Как ни странно, но многие из ответов мы нашли в Twitter. По этому поводу сейчас мы сделаем короткий перерыв, чтобы поговорить с вами об этом, а затем уже перейдем к рассмотрению цветных изображений.

0 Twitter

Вообще, мы не активные пользователи соцсетей. Но цель этой книги — помочь вам стать преуспевающим специалистом в области глубокого обучения, и было бы упущением не упомянуть, насколько важен оказался Twitter в нашем собственном путешествии по миру ML.

Здесь просто нужно учесть, что у этой социальной сети есть и другая сторона, отдаленная от Дональда Трампа и семейства Кардашьян, в которой исследующие и практикующие глубокое обучение специалисты ведут постоянный диалог. В ходе написания текущего раздела Джереми хотел перепроверить все сказанное нами о свертках с шагом 2, поэтому задал вопрос в Twitter:



The screenshot shows a tweet from Jeremy Howard (@jeremyphoward). The profile picture is a stylized face. The tweet text is:

I forget: why did we move from stride 1 convs with maxpool to stride 2 convs? And why don't we use stride 1 convs with avgpool (or do some modern nets do that)? Is it just an empirical thing, or is there some deeper reason? Has someone done the ablation studies?

At the bottom of the tweet, it says "11:21 AM · Feb 23, 2020 · Twitter Web App".

Спустя несколько минут был получен ответ:

 **Christian Szegedy**
@ChrSzegedy

Replies to @jeremyphoward

This depends on a lot of factors: the overall network architecture, the accelerator (CPU vs GPU vs TPU) etc. Some Inception models used concatenation of (conv + max/avg/l2 pooling). The quality differences were marginal. Some pooling methods are more residual friendly.

11:39 AM · Feb 23, 2020 · Twitter Web App

Кристиан Сегеди (Christian Szegedy) — это первый автор Inception (https://oreil.ly/hGE_Y) (Зарождение), победитель соревнования ImageNet 2014 года. Он сделал многие ключевые открытия, используемые в современных нейронных сетях. Спустя два часа появилось следующее сообщение:

 **Yann LeCun**
@ylecun

Replies to @jeremyphoward

My original early 1989 NeurComp paper on ConvNet used stride 2, no pooling, simply because the computation was fast.
The 2nd paper (NIPS 1989) used stride +average pooling/tanh.
It worked better on zipcode digits. But it could have been due to many reasons.

1:35 PM · Feb 23, 2020 · Twitter for Android

Узнаете имя? Оно встречалось нам в главе 2, когда мы обсуждали лауреатов премии Тьюринга, заложивших основы современного глубокого обучения. Джереми также обратился в Twitter с просьбой проверить точность нашего описания сглаживания меток, о котором шла речь в главе 7. Ответил ему Кристиан Сегеди (изначально сглаживание меток было представлено в работе Inception):

 **Christian Szegedy**
@ChrSzegedy

Replies to @jeremyphoward

It was mostly written by Sergey, so he might be the best person to ask. IMO, yours is a fair motivation of label-smoothing. It interprets the passage correctly.

5:24 PM · Feb 21, 2020 · Twitter Web App

Многие из значимых в мире глубокого обучения людей — активные пользователи Twitter и охотно идут на контакт с сообществом. Ради интереса можно начать с просмотра последних лайков Джереми или Сильвейна. Так вы увидите список пользователей Twitter, которым есть что сказать полезного и интересного.

Twitter является основным средством, помогающим нам оставаться в курсе интересных научных работ, релизов ПО и других новостей из области глубокого обучения. Для установления связи с сообществом мы рекомендуем стать участником как форумов fast.ai (<https://forums.fast.ai/>), так и Twitter.

Ну а теперь вернемся к сути главы. До сих пор мы показывали примеры картинок только в черно-белом цвете, у которых для каждого пикселя было по одному значению. На практике же большинство полноцветных изображений содержат по три определяющих цвет значения для каждого пикселя. Так что далее мы переходим к примерам именно цветных изображений.

Цветные изображения

Цветное изображение — это тензор ранга 3:

```
im = image2tensor(Image.open('images/grizzly.jpg'))
im.shape

torch.Size([3, 1000, 846])

show_image(im);
```



Первая ось содержит каналы красного, зеленого и синего:

```
_,axs = subplots(1,3)
for bear,ax,color in zip(im,axs,('Reds','Greens','Blues')):
    show_image(255-bear, ax=ax, cmap=color)
```



Мы видели операцию свертки для одного фильтра на одном канале изображения (примеры выполнялись для квадрата). Сверточный слой получает изображение с определенным числом каналов (три для первого слоя в стандартных RGB-изображениях), а выводит изображение уже с другим числом каналов. Как и в случае со скрытым размером, который представлял количество нейронов в линейном слое, мы имеем возможность использовать столько фильтров, сколько захотим. При этом каждый из них сможет специализироваться (одни для обнаружения горизонтальных краев, другие — для вертикальных и т. д.), формируя что-то наподобие тех примеров, которые мы рассматривали в главе 2.

В одном скользящем окне у нас есть определенное число каналов, и нам нужно столько же фильтров (мы не задействуем одинаковое ядро для всех каналов). Поэтому размер ядра будет не 3×3 , а ch_in на 3×3 . На каждом канале мы умножаем элементы окна на элементы соответствующего фильтра, суммируем результаты (как делали ранее), после чего складываем эти результаты по всем фильтрам. В примере на рис. 13.12 итогом сверточного слоя для рассматриваемого окна будет красный + зеленый + синий.

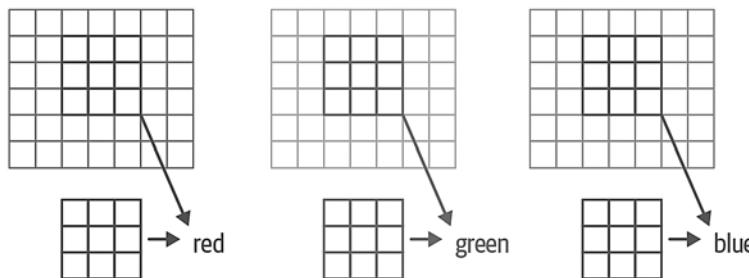


Рис. 13.12. Свертка изображения RGB

Поэтому чтобы применить свертку к цветной картинке, нам нужен тензор ядер, размер которого будет соответствовать первой оси. В каждой области соответствующие части ядра и участка изображения перемножаются.

После этого полученные результаты суммируются, образуя одно число для каждой области сетки каждого выходного признака, как показано на рис. 13.13.

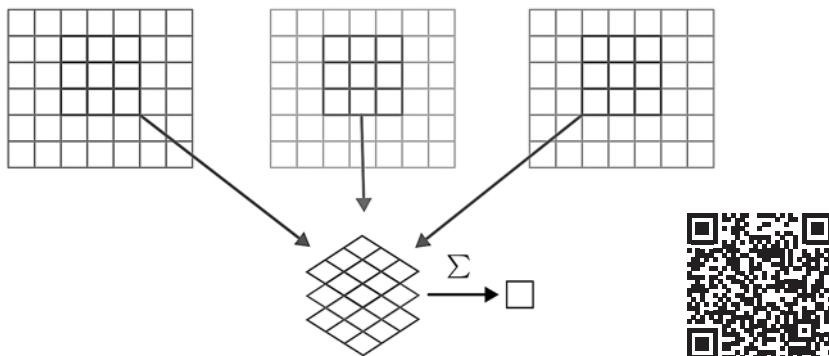


Рис. 13.13. Сложение RGB-фильтров

Так мы получаем `ch_out` фильтров, и в итоге результатом сверточного слоя будет пакет изображений с `ch_out` каналов, а также высотой и шириной, заданными обозначенной ранее формулой. Это дает нам `ch_out` тензоров размером `ch_in x ks x ks`, которые мы представляем в виде одного большого четырехмерного тензора. В PyTorch порядок измерений весов будет таким: `ch_out x ch_in x ks x ks`.

Дополнительно к этому для каждого фильтра нам может потребоваться смещение. В предыдущем примере итоговый результат нашего сверточного слоя был равен $\gamma_R + \gamma_G + \gamma_B + b$. Как и в линейном слое, смещений будет столько же, сколько имеется ядер, поэтому размер их вектора — `ch_out`.

При настройке CNN для обучения цветных изображений никакой особый механизм не задействуется. Просто убедитесь, что у первого слоя три входа.

Для обработки цветных изображений есть много способов. Например, можно перевести их в черно-белые, из цветового пространства RGB в HSV (тон, насыщенность, значение) и т. п. Эксперименты, как правило, доказывают, что изменение цветовой кодировки не внесет большой разницы в результаты модели, если только вы не утратите в ходе преобразования часть информации. Поэтому перевод изображений в черно-белые не будет удачной идеей, поскольку полностью утратится информация цвета (что может оказаться очень важным, так как, например, порода животного может отличаться конкретным окрасом). При этом перевод в HSV обычно ничего не меняет.

Теперь вы знаете, что означают те картинки из работы Зейлера и Фергуса (<https://oreil.ly/Y6dzZ>), о которой мы писали в главе 1. В качестве напоминания приведем одну из этих картинок, отражающую веса первого слоя:



Здесь берутся три среза сверточного ядра для каждого выходного признака и показываются в виде изображений. Можно заметить, что хотя создатели нейронной сети и не создавали напрямую ядра для нахождения контуров, она автоматически обнаружила эти признаки с помощью SGD.

Далее мы рассмотрим обучение этих CNN и покажем вам все внутренние техники, используемые fastai для повышения эффективности обучения.

Повышение стабильности обучения

Мы уже хорошо можем отличать тройки от семерок, так что перейдем к задаче посложнее: распознаванию всех десяти цифр. Это подразумевает использование MNIST вместо MNIST_SAMPLE:

```
path = untar_data(URLs.MNIST)

path.ls()

(#2) [Path('testing'), Path('training')]
```

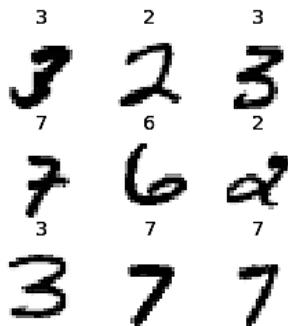
Данные находятся в двух каталогах: *training* и *testing*, о чем необходимо сообщить `GrandparentSplitter` (по умолчанию он задействует `train` и `valid`). Это мы делаем с помощью функции `get_dls`, которая в дальнейшем упростит для нас изменение размера пакета:

```
def get_dls(bs=64):
    return DataBlock(
        blocks=(ImageBlock(cls=PILImageBW), CategoryBlock),
        get_items=get_image_files,
        splitter=GrandparentSplitter('training', 'testing'),
        get_y=parent_label,
        batch_tfms=Normalize()
    ).dataloaders(path, bs=bs)

dls = get_dls()
```

Не забывайте, что перед использованием данных всегда желательно их сначала изучить:

```
dls.show_batch(max_n=9, figsize=(4,4))
```



Закончив подготовку данных, мы можем обучить на них простую модель.

Базовая модель

Ранее в этой главе мы создали модель на основе функции `conv`:

```
def conv(ni, nf, ks=3, act=True):
    res = nn.Conv2d(ni, nf, stride=2, kernel_size=ks, padding=ks//2)
    if act: res = nn.Sequential(res, nn.ReLU())
    return res
```

Для начала задействуем в качестве основы простую CNN. Мы используем ту же, что и ранее, но с одной доработкой: задействуем больше активаций. Поскольку нам нужно дифференцировать большее число цифр, то и фильтров для этого понадобится больше.

Как мы уже говорили, при использовании сверточного слоя с шагом 2 мы каждый раз удваиваем количество фильтров. Один из способов увеличения их числа по ходу сети — это удвоивание активаций в первом слое, в результате чего и каждый последующий слой будет получаться вдвое больше предыдущего.

Но здесь возникает небольшая проблема. Рассмотрим ядро, применяемое к каждому пикселью. По умолчанию мы используем размер 3×3 . Следовательно, в каждой области это ядро применяется к 9 пикселям. Ранее у нашего первого слоя было четыре выходных фильтра, в результате чего из девяти пикселей в каждой области вычислялось четыре значения. Подумайте, что произойдет, если мы удвоим количество выходных фильтров до восьми. Тогда при применении ядра мы будем вычислять на основе девяти пикселей уже девять чисел. Это означает, что обучения практически не происходит: размер выхода почти равен

размеру входа. Нейронные сети будут создавать полезные признаки, только если их к этому подтолкнуть, то есть сделать число выходов операции значительно меньше числа входов.

Чтобы это исправить, можно задействовать в первом слое более крупное ядро. Если взять размер 5×5 , то при его применении будет рассматриваться уже 25 пикселей. Создание на их основе восьми фильтров будет означать, что нейронной сети придется найти ряд полезных признаков:

```
def simple_cnn():
    return sequential(
        conv(1, 8, ks=5),           # 14 × 14
        conv(8, 16),               # 7 × 7
        conv(16, 32),              # 4 × 4
        conv(32, 64),              # 2 × 2
        conv(64, 10, act=False),   # 1 × 1
        Flatten(),
    )
```

Как вы вскоре увидите, мы можем заглянуть внутрь моделей в момент их обучения с целью найти способы улучшения этого обучения. Для этого мы используем обратный вызов `ActivationStats`, который записывает среднее значение, стандартное отклонение и гистограмму активаций каждого обучаемого слоя (как мы видели, обратные вызовы задействуются для добавления в обучающий цикл нужного поведения; принцип их работы мы рассмотрим подробнее в главе 16):

```
from fastai.callback.hook import *
```

Нам нужно быстрое обучение, что подразумевает установку более высокой скорости. Давайте посмотрим, как мы будем работать при 0,06:

```
def fit(epochs=1):
    learn = Learner(dls, simple_cnn(), loss_func=F.cross_entropy,
                    metrics=accuracy, cbs=ActivationStats(with_hist=True))
    learn.fit(epochs, 0.06)
    return learn

learn = fit()
```

epoch	train_loss	valid_loss	accuracy	time
0	2.307071	2.305865	0.113500	00:16

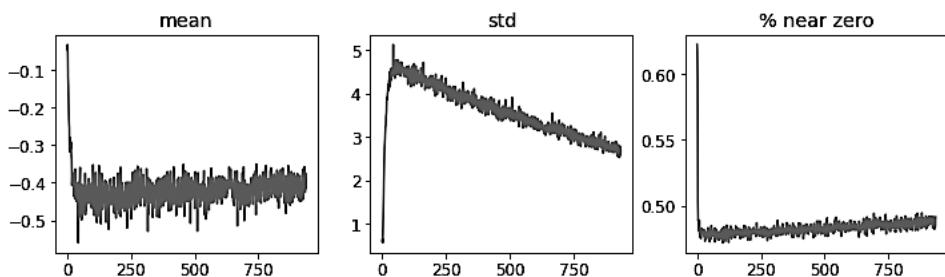
Результат абсолютно никчемный. Давайте выясним почему.

Одна из удобных особенностей передаваемых в `Learner` обратных вызовов заключается в том, что они становятся доступны автоматически с тем же именем, что и их класс, только не в `snake_case`. Так что к нашему обратному вызову

`ActivationStats` можно обратиться через `activation_stats`. Уверены, что вы помните `learn.recorder...` можете догадаться, как он реализуется? Правильно, это обратный вызов `Recorder`.

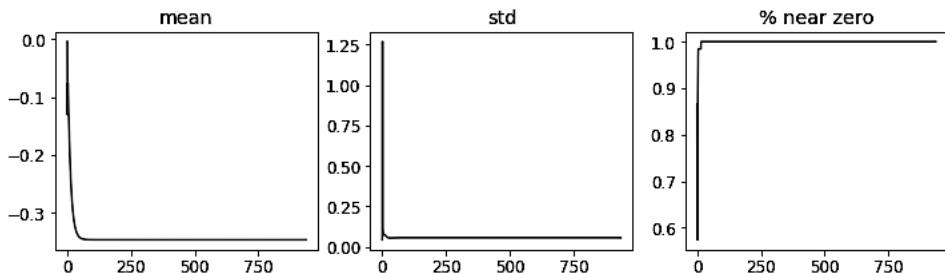
`ActivationStats` содержит ряд удобных утилит для графического отображения активаций в процессе обучения. `plot_layer_stats(idx)` отображает среднее значение и стандартное отклонение активаций слоя номер `idx` наряду с процентом активаций, приближенных к нулю. Вот график для первого слоя:

```
learn.activation_stats.plot_layer_stats(0)
```



Как правило, у модели в процессе обучения должно быть единообразное или по меньшей мере гладкое среднее и стандартное отклонение активаций слоя. Близкие к нулю активации особенно проблематичны, так как это означает, что выполняемые в модели вычисления практически ничего не дают (поскольку умножение на ноль дает ноль). Когда у вас в слое присутствуют нули, они будут переноситься в последующий слой, который, в свою очередь, будет создавать дополнительные нули, и т. д. Вот предпоследний слой нашей сети:

```
learn.activation_stats.plot_layer_stats(-2)
```



Как и ожидалось, проблемы усугубляются к концу сети, так как нестабильность и нулевые активации по ходу слоев нарастают. Давайте посмотрим, что в таком случае можно предпринять для стабилизации обучения.

Увеличение размера пакета

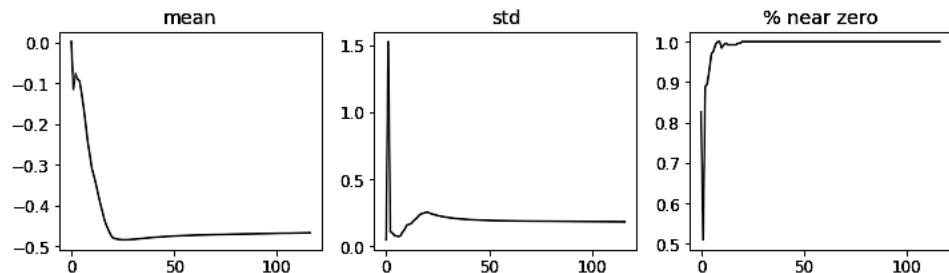
Один из способов стабилизации обучения — это увеличение размера пакета. Более крупные пакеты имеют более точные градиенты, поскольку вычисляются они на основе большего объема данных. Хотя, с другой стороны, увеличение размера пакетов приведет к уменьшению их количества в эпоху, что снизит возможности модели для обновления весов. Посмотрим, поможет ли нам изменение размера пакета до 512:

```
dls = get_dls(512)
learn = fit()
```

epoch	train_loss	valid_loss	accuracy	time
0	2.309385	2.302744	0.113500	00:08

Вот как выглядит предпоследний слой теперь:

```
learn.activation_stats.plot_layer_stats(-2)
```



Большинство активаций опять скатились к нулю. Посмотрим, что еще можно сделать в этом отношении.

Обучение 1cycle

Наши изначальные веса недостаточно подходят для решаемой задачи. Поэтому начинать обучение с высокой скоростью опасно, так как мы видели, что это может привести к его хаотичному расхождению. При этом также не стоит заканчивать обучение на высокой скорости, чтобы не проскочить искомый минимум. Тем не менее в остальной период ее можно повысить, что позволит ускорить обучение в целом. Следует изменять скорость в процессе, начиная с низкой, а затем перейти на высокую и закончить опять же низкой.

Лесли Смит (Leslie Smith) (да, тот самый, который придумал искатель скорости обучения) развивает эту идею в статье *Super-Convergence: Very Fast Training of*

Neural Networks Using Large Learning Rates (<https://oreil.ly/EB8NU>) («Суперконвергенция: очень быстрое обучение нейронных сетей на высоких скоростях»). Он разработал порядок изменения скорости обучения, представляющий две фазы: в первой скорость возрастает с минимальной до максимальной (*разогрев*), а во второй — уменьшается обратно к минимуму (*отжиг, annealing*). Смит назвал эту комбинацию подходов *обучением 1cycle*.

Данный метод позволяет использовать гораздо более высокую скорость обучения по сравнению с другими подходами, предоставляя при этом два преимущества.

- Повышение общей скорости обучения — явление, которое Смит назвал *суперконвергенцией*.
- При более высокой скорости обучения снижается вероятность переобучения, так как мы пропускаем острые локальные минимумы, оказываясь в более гладких (следовательно, и более обобщаемых) участках функции потерь.

Далее следует интересная и утонченная деталь, которая основана на том, что хорошо обобщающаяся модель — это та, чьи потери почти не меняются при небольшом изменении входных данных. Если модель обучается с высокой скоростью достаточно долго и может при этом найти хороший показатель потерь, то скорее всего, она находит хорошо обобщающуюся область благодаря частому перескакиванию от пакета к пакету (что, по сути, и определяет высокую скорость обучения). Проблема же в том, что, как мы и говорили, резкий переход на высокую скорость приводит к отклонению показателей потерь, а не их улучшению. Поэтому мы не переходим резко на высокую скорость. Вместо этого мы начинаем с низкой, где показатели потерь не расходятся, и позволяем оптимизатору постепенно найти все более гладкие области параметров путем постепенного увеличения скорости.

Затем, когда мы нашли достаточно гладкую область для параметров, мы переходим к поиску наилучшего участка этой области, для чего скорость обучения нужно снова снизить. Именно поэтому метод 1cycle задействует постепенный «разогрев» с последующим «охлаждением». Многие исследователи обнаружили, что на практике этот подход ведет к формированию более точных моделей и ускоряет обучение в целом. В связи с этим он по умолчанию используется для `fine_tune` в fastai.

В главе 16 мы подробно разберем понятие *импульса* в SGD. Говоря коротко, импульс — это техника, в которой оптимизатор не только совершает шаг в направлении градиентов, но также продолжает движение по направлению предыдущих шагов. Лесли Смит предложил в работе *A Disciplined Approach to Neural Network Hyper-Parameters: Part 1* (<https://oreil.ly/oL7GT>) («Дисциплинированный подход к гиперпараметрам нейронной сети: часть 1») идею *циклического импульса*.

Он предполагает, что импульс изменяется в противоположном к скорости обучения направлении: при высокой скорости мы задействуем меньший импульс, увеличивая его в фазе отжига.

Использовать обучение 1cycle в fastai можно с помощью вызова `fit_one_cycle`:

```
def fit(epochs=1, lr=0.06):
    learn = Learner(dls, simple_cnn(), loss_func=F.cross_entropy,
                    metrics=accuracy, cbs=ActivationStats(with_hist=True))
    learn.fit_one_cycle(epochs, lr)
    return learn

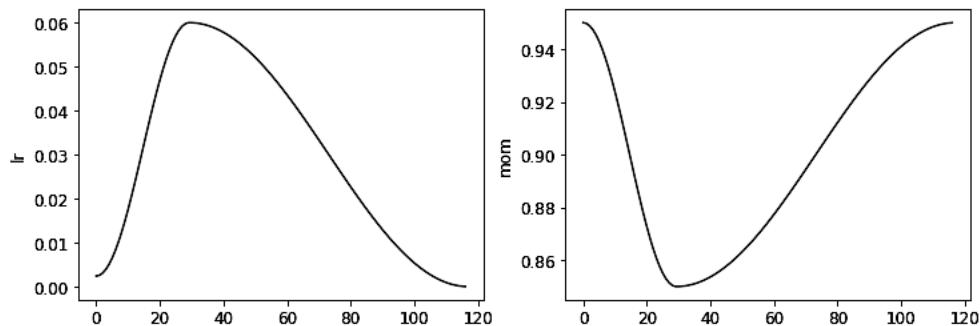
learn = fit()
```

epoch	train_loss	valid_loss	accuracy	time
0	0.210838	0.084827	0.974300	00:08

Наконец-то мы добились успеха, достигнув вполне адекватной точности!

Просмотреть значения скорости обучения и импульса в процессе обучения можно, вызвав для `learn.recorder` метод `plot_sched`. `learn.recorder` записывает все происходящее в ходе обучения, включая потери, метрики и такие гиперпараметры, как скорость обучения и импульс:

```
learn.recorder.plot_sched()
```



В первой своей работе, посвященной обучению 1cycle, Смит использовал линейный разогрев и линейный отжиг. Вы можете заметить, что мы адаптировали этот подход в fastai, совместив его с другим популярным подходом: косинусным отжигом. `fit_one_cycle` предоставляет следующие настраиваемые параметры.

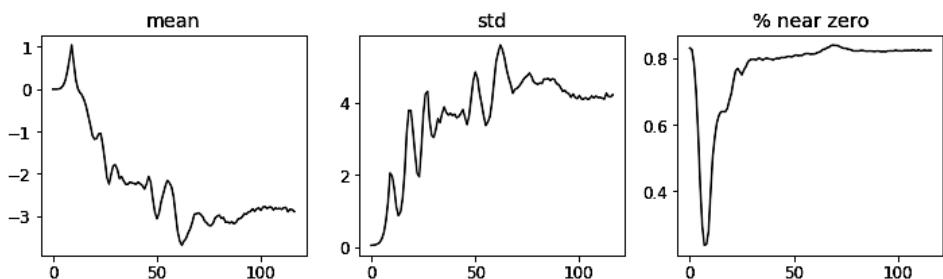
- `lr_max` — максимальная скорость обучения, которая будет использована (она может быть представлена списком скоростей для каждой группы сло-

ев либо Python-объекта `slice`, содержащим скорости первой и последней группы слоев).

- `div` — на сколько разделить `lr_max` для получения стартовой скорости обучения.
- `div_final` — на сколько разделить `lr_max` для получения конечной скорости обучения.
- `pct_start` — какой процент пакетов использовать для разогрева.
- `mom` — кортеж (`mom1, mom2, mom3`), где `mom1` — это начальный импульс, `mom2` — минимальный импульс, а `mom3` — финальный.

Еще раз взглянем на состояние слоев:

```
learn.activation_stats.plot_layer_stats(-2)
```



Показатель оклонулевых весов существенно улучшается, хотя еще недостаточно. С помощью `color_dim` можно увидеть дополнительные подробности о происходящем в процессе обучения, передав в него индекс слоя:

```
learn.activation_stats.color_dim(-2)
```



`color_dim` был разработан fast.ai совместно с нашим студентом Стефано Джомо (Stefano Giomo). Джомо определяет его как *цветовое измерение*, давая подробное объяснение (<https://oreil.ly/bPXGw>) истории и деталей, стоящих за этим методом. Основная идея — в создании гистограммы активаций слоя, которая,

как мы надеемся, будет следовать плавному шаблону аналогично нормальному распространению (рис. 13.14).

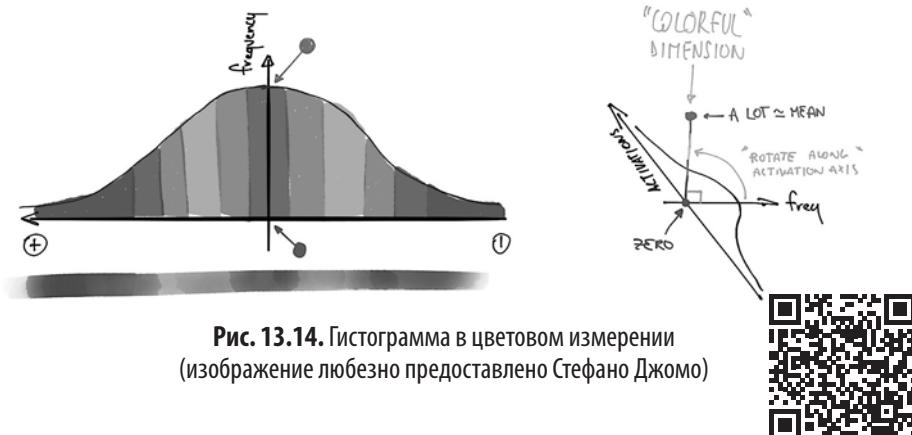


Рис. 13.14. Гистограмма в цветовом измерении
(изображение любезно предоставлено Стефано Джомо)

Чтобы создать `color_dim`, мы берем гистограмму, показанную слева, и преобразуем ее в цветовое представление, приведенное под ней. Затем мы поворачиваем ее на бок, как показано справа. Мы выяснили, что распределение будет более отчетливым, если получить логарифм значений гистограммы. Далее Джомо говорит:

Последний график для каждого слоя получается путем составления гистограммы активаций из каждого пакета вдоль горизонтальной оси. Таким образом, каждый вертикальный срез визуализации представляет гистограмму активаций одного пакета. Интенсивность цвета соответствует высоте гистограммы; иначе говоря, числу активаций в каждой ячейке гистограммы.

На рис. 13.15 показано, как все это соотносится.

Это показывает, почему $\log(f)$ более «цветастая», чем f , когда f следует нормальному распределению, потому что получение логарифма изменяет гауссову кривую на аппроксимированную квадратичную функцию.

С учетом этого еще раз взглянем на результат предпоследнего слоя:

```
learn.activation_stats.color_dim(-2)
```



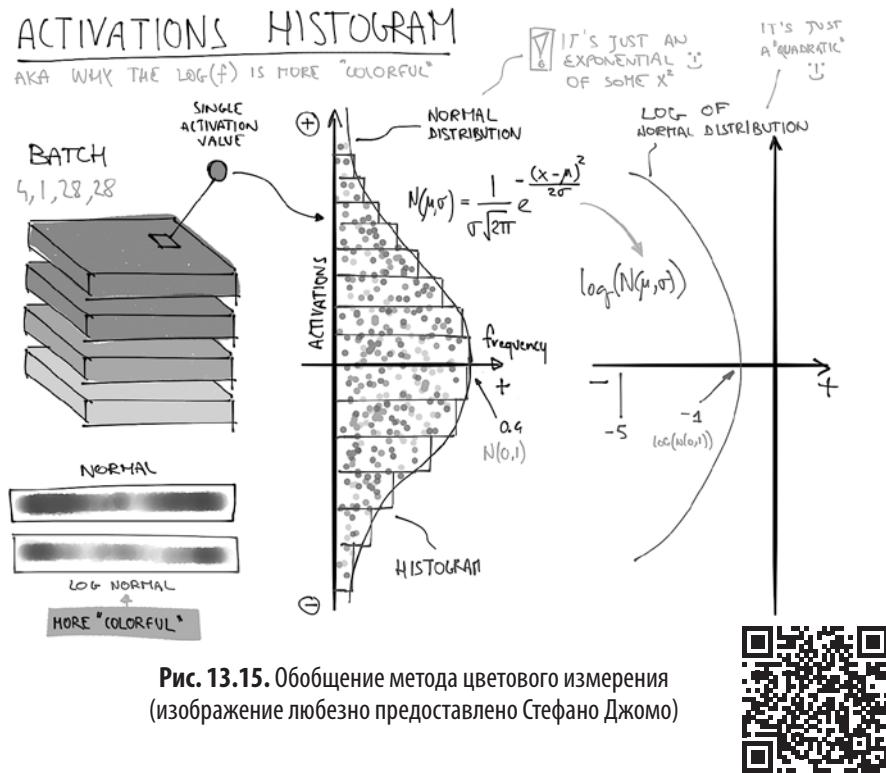


Рис. 13.15. Обобщение метода цветового измерения
(изображение любезно предоставлено Стефано Джомо)



Здесь мы видим классическую картину «плохого обучения». Начинаем мы при почти всех активациях, близких к нулю: это видно слева с краю, где большая часть области заполнена темно-синим. Ярко-желтый цвет в нижней части графика представляет околонулевые активации. Затем, спустя пару пакетов, мы наблюдаем экспоненциальный рост ненулевых активаций. Но он заходит слишком далеко, и происходит разрушение. Здесь снова появляется темно-синий, и нижняя часть обретает ярко-желтый цвет. Похоже, как будто обучение начинается заново. Затем активации вновь возрастают и разрушаются. Спустя несколько таких циклов наблюдается распространение активаций по всему диапазону.

Гораздо лучше, если обучение будет идти гладко изначально. Циклы экспоненциального роста и последующего разрушения приводят ко множеству околонулевых активаций, замедляя обучение и давая плохие итоговые результаты. Для решения этой проблемы можно задействовать пакетную нормализацию.

Пакетная нормализация

Чтобы исправить проблему медленного обучения и плохих итоговых результатов, важно разобраться с изначально большим процентом околонулевых акти-

ваций и постараться поддерживать их хорошее распределение на протяжении всего обучения.

Сергей Иоффе (Sergey Ioffe) и Кристиан Сегеди (Christian Szegedy) представили решение этой проблемы в 2015 году в своей работе под названием *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift* (<https://oreil.ly/MTZJL>) («Нормализация пакетов: ускорение обучения глубокой нейронной сети путем уменьшения внутреннего ковариантного сдвига»). В аннотации они описывают проблему, с которой мы только что столкнулись:

Обучение глубоких нейронных сетей усложняется тем фактом, что распределение входных данных каждого слоя в течение этого процесса изменяется, так как меняются параметры предшествующих слоев. В связи с этим требуется снижение скорости обучения и более осторожная инициализация параметров, что и замедляет весь итоговый процесс. Мы называем это явление внутренним ковариантным сдвигом и решаем его проблему нормализацией входных данных слоев.

Свое решение они описывают так:

Внедрить нормализацию как часть архитектуры модели и выполнить ее для каждого обучающего мини-пакета. Пакетная нормализация позволяет использовать гораздо более высокую скорость обучения и не утруждаться тщательностью инициализации.

Публикация этой научной работы вызвала бурное восхищение, так как включала график, приведенный на рис. 13.16, который показывает, что с применением нормализации пакетов можно обучить модель даже более точную, чем ее эталонный аналог (архитектура *Inception*), и примерно в пять раз быстрее.

Пакетная нормализация (зачастую называемая *batchnorm*) работает путем получения среднего математического ожидания и стандартных отклонений активаций слоя, используя их для нормализации активаций. Тем не менее это грозит проблемами, так как сети для точного прогнозирования может потребоваться, чтобы некоторые активации были очень высокими. В связи с этим авторы также добавили два обучаемых параметра (которые будут обновляться на этапе SGD), как правило, называемых *gamma* и *beta*. После нормализации активаций для получения нового вектора активаций слой *batchnorm* возвращает $\text{gamma} * y + \text{beta}$.

Благодаря этому активации могут иметь любое среднее или дисперсию, которые не будут зависеть от среднего или стандартного отклонения результатов предыдущего слоя. Обучаются эти статистики отдельно, что упрощает обучение модели. При этом процесс отличается для этапа обучения и контроля: в процессе обучения мы используем среднее значение и стандартное отклонение пакета для нормализации данных, а во время контроля мы, наоборот, используем скользящее среднее статистик, вычисленное в ходе обучения.

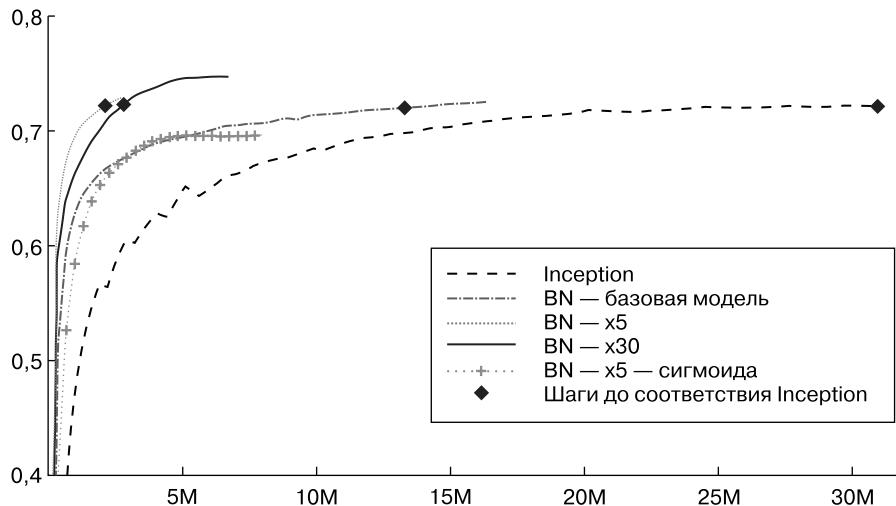


Рис. 13.16. Влияние нормализации пакетов (изображение любезно предоставлено Сергеем Иоффе и Кристианом Сегеди)

Добавим слой пакетной нормализации в conv:

```
def conv(ni, nf, ks=3, act=True):
    layers = [nn.Conv2d(ni, nf, stride=2, kernel_size=ks, padding=ks//2)]
    layers.append(nn.BatchNorm2d(nf))
    if act: layers.append(nn.ReLU())
    return nn.Sequential(*layers)
```

И подстроим модель:

```
learn = fit()
```

epoch	train_loss	valid_loss	accuracy	time
0	0.130036	0.055021	0.986400	00:10

Отличный результат! Давайте взглянем на `color_dim`:

```
learn.activation_stats.color_dim(-4)
```



Это мы и хотим видеть: плавное развитие активаций без провалов. Пакетная нормализация реально выполнила обещанное. В действительности эта техника оказалась настолько успешной, что мы встречаем ее (или нечто очень похожее) практически во всех современных нейронных сетях.

Интересное наблюдение относительно моделей, содержащих слои пакетной нормализации, в том, что они лучше обобщаются, чем их аналоги без нормализации. Хотя мы пока и не встречали тщательного анализа этого процесса, большинство исследователей видят причину в том, что пакетная нормализация привносит в процесс обучения дополнительную случайность. Каждый мини-пакет будет иметь математическое ожидание (среднее) и стандартное отклонение, несколько отличные от всех других пакетов. Следовательно, активации будут каждый раз нормализоваться разными значениями. Чтобы делать точные прогнозы, модели потребуется обучиться устойчивости к этим вариациям. Внесение же дополнительной случайности в процесс обучения зачастую помогает.

Раз уж все у нас идет хорошо, продолжим обучение в течение еще нескольких эпох и посмотрим, что получится. Давайте даже *увеличим* скорость обучения, поскольку авторы пакетной нормализации заявили об «обучении на гораздо более высоких скоростях»:

```
learn = fit(5, lr=0.1)
```

epoch	train_loss	valid_loss	accuracy	time
0	0.191731	0.121738	0.960900	00:11
1	0.083739	0.055808	0.981800	00:10
2	0.053161	0.044485	0.987100	00:10
3	0.034433	0.030233	0.990200	00:10
4	0.017646	0.025407	0.991200	00:10

```
learn = fit(5, lr=0.1)
```

epoch	train_loss	valid_loss	accuracy	time
0	0.183244	0.084025	0.975800	00:13
1	0.080774	0.067060	0.978800	00:12
2	0.050215	0.062595	0.981300	00:12
3	0.030020	0.030315	0.990700	00:12
4	0.015131	0.025148	0.992100	00:12

Вот теперь можно с уверенностью сказать, что мы умеем распознавать цифры. Пришло время переходить к задачам посложнее...

Резюме

Мы видели, что свертки представляют собой просто вид матричного умножения с двумя ограничениями в матрице весов: некоторые элементы всегда равны нулю, а некоторые связаны (вынуждены всегда иметь одинаковые значения). В главе 1 мы приводили восемь условий из книги 1986 года «Параллельная распределенная обработка». Одним из них было условие наличия «образца связи между нейронами». Именно это и реализуют данные ограничения: обусловливают определенный шаблон связей.

Эти ограничения позволяют использовать в модели намного меньше параметров, не жертвуя возможностью представлять сложные визуальные признаки. Это значит, что мы можем обучать более глубокие модели быстрее с меньшим риском переобучения. И хотя теорема об универсальной аппроксимации показывает, что *теоретически* в полносвязной сети *возможно* представить что угодно в одном скрытом слое, теперь мы видели, что *на практике* можно обучать более эффективные модели, просто продумав архитектуру сети.

Свертки являются наиболее распространенным шаблоном связей в нейронных сетях (наряду со стандартными линейными слоями, которые мы называем *полносвязными*), но наверняка скоро будут разработаны и другие.

Помимо этого, мы научились интерпретировать активации слоев сети, чтобы понять, успешно ли проходит обучение, и узнали, как пакетная нормализация помогает регуляризовать обучение, делая его более плавным. В следующей главе мы используем оба вида этих слоев для построения самой популярной архитектуры в компьютерном зрении: остаточной нейронной сети.

Вопросник

1. Что такое признак?
2. Напишите матрицу сверточного ядра для обнаружения верхних контуров.
3. Напишите математическую операцию, применяемую ядром 3×3 к одному пикселю изображения.
4. Каково будет значение сверточного ядра, примененного к матрице нулей размером 3×3 ?
5. Что такое заполнение?
6. Что такое шаг?

7. Создайте генератор вложенных списков для выполнения задачи по вашему выбору.
8. Каковы формы параметров `input` и `weight` для 2D-свертки PyTorch?
9. Что такое канал?
10. Что общего между сверткой и матричным умножением?
11. Что такое сверточная нейронная сеть?
12. В чем польза рефакторинга компонентов определения нейронной сети?
13. Что такое `Flatten`? Где требуется его добавление в MNIST CNN? Почему?
14. Что означает NCHW?
15. Почему третий слой MNIST CNN содержит $7*7*(1168-16)$ умножений?
16. Что такое рецептивное поле?
17. Каков будет размер рецептивного поля активации после двух сверток с шагом 2? Почему?
18. Выполните самостоятельно `conv-example.xlsx` и поэкспериментируйте с отслеживанием прецедентов.
19. Загляните в список недавних «лайков» Джереми и Сильвейна — возможно, там вы найдете интересную полезную информацию.
20. Как цветное изображение представляется в виде тензора?
21. Как работает свертка с цветным входным изображением?
22. Какой метод можно использовать для просмотра этих данных в `DataLoaders`?
23. Почему мы удваиваем количество фильтров после каждой свертки с шагом 2?
24. Почему при работе с MNIST (в случае с `simple_cnn`) мы используем в первой свертке большее ядро?
25. Какую информацию `ActivationStats` сохраняет для каждого слоя?
26. Как обратиться к обратному вызову `learner` после обучения?
27. Какие три статистики отражает `plot_layer_stats`? Что представляет ось x?
28. Почему оклонулевые активации вызывают проблемы?
29. Каковы плюсы и минусы обучения с увеличенным размером пакета?
30. Почему стоит избегать использования высокой скорости обучения в начале этого процесса?
31. Что такое обучение 1cycle?
32. Каковы преимущества обучения с высокой скоростью?
33. Почему мы предпочитаем заканчивать обучение на низкой скорости?
34. Что такое циклический импульс?

35. Какой обратный вызов отслеживает в процессе обучения значения гиперпараметров (наряду с остальной информацией)?
36. Что представляет столбец пикселей на графике `color_dim`?
37. Как на графике `color_dim` выглядит «плохое обучение»? Почему?
38. Какие обучаемые параметры содержит слой пакетной нормализации?
39. Какие статистики используются для нормализации в процессе обучения?
А какие в процессе контроля?
40. Почему модели, использующие пакетную нормализацию, лучше обобщаются?

Дополнительные задания

1. Какие признаки помимо детекторов контуров использовались в компьютерном зрении (особенно до роста популярности глубокого обучения)?
2. В PyTorch доступны и другие слои нормализации. Попробуйте их и определите наиболее эффективные. Узнайте, почему были разработаны другие слои нормализации и чем они отличаются от пакетной нормализации.
3. Попробуйте перенести функцию активации после слоя пакетной нормализации в `conv`. Видите ли вы разницу? Поиските информацию о том, какой порядок рекомендуется и почему.

ГЛАВА 14

ResNet

В этой главе мы будем делать надстройку над CNN, которую освоили в предыдущей главе, и объясним суть архитектуры ResNet (остаточная сеть). Эта архитектура была впервые представлена в 2015 году Каймингом Хе (Kaiming He) и др. в статье *Deep Residual Learning for Image Recognition* (<https://oreil.ly/b68K8>) («Глубокое остаточное обучение для распознавания изображений») и пока что является наиболее широко используемой. Более современные разработки в моделях распознавания изображений почти всегда задействуют тот же прием остаточных связей и в большинстве случаев просто являются вариациями оригинальной ResNet.

Сначала мы покажем вам базовую ResNet в том виде, в каком она была спроектирована изначально, а затем объясним современные доработки, повышающие ее эффективность. Но сперва нам понадобится определить несколько более сложную задачу, чем датасет MNIST, поскольку в нем мы уже приблизились к точности 100 % с помощью стандартной CNN.

Возвращение к Imagenette

Будет непросто оценивать любые улучшения, вносимые нами в модели, когда мы уже достигли высокой точности, как с датасетом MNIST из предыдущей главы. Поэтому мы возьмемся за более сложную задачу классификации, вернувшись к Imagenette. В данном случае с целью ускорения процесса мы задействуем небольшие изображения.

Соберем данные — мы будем использовать их версию с уже измененным размером до 160 пикселей и сделаем случайную обрезку до 128 пикселей:

```
def get_data(url, presize, resize):
    path = untar_data(url)
    return DataBlock(
        blocks=(ImageBlock, CategoryBlock), get_items=get_image_files,
        splitter=GrandparentSplitter(valid_name='val'),
        get_y=parent_label, item_tfms=Resize(presize),
```

```

batch_tfms=[*aug_transforms(min_scale=0.5, size=resize),
    Normalize.from_stats(*imagenet_stats)],
).dataloaders(path, bs=128)

dls = get_data(URLs.IMAGENETTE_160, 160, 128)
dls.show_batch(max_n=4)

```



При работе с MNIST мы имели дело с изображениями 28×28 пикселей. В случае с Imagenette мы будем производить обучение на изображениях 128×128 пикселей. Позже нам бы хотелось иметь возможность использовать изображения большего размера: по меньшей мере 224×224 пикселя, что является стандартом ImageNet. Вспомните ли вы, как нам удавалось получить по одному вектору активаций для каждого изображения сверточной нейронной сети на датасете MNIST?

Использованный нами подход подразумевал такое количество сверток с шагом 2, чтобы размер сетки последнего слоя получался равным 1. Затем мы просто сглаживали полученные единичные оси, создавая вектор для каждого изображения (матрицу активаций для мини-пакета). То же самое можно сделать и для Imagenette, но это вызовет две проблемы.

- Для получения сетки 1×1 нам потребуется множество слоев с шагом 2 — возможно, больше, чем мы выбрали бы в противном случае.
- Модель не будет работать на изображениях, отличающихся по размеру от изображений, на которых она обучалась.

Для решения первой проблемы можно сгладить финальный сверточный слой так, чтобы он обрабатывал сетку с размером, превышающим 1×1 . Для этого можно

просто преобразовать матрицу в вектор, как мы уже делали ранее, выстроив каждую строку после ей предшествующей. В действительности такой подход нейронные сети задействовали практически всегда вплоть до 2013 года. Самый известный пример — сеть VGG, победившая в соревновании ImageNet в тот самый 2013 год, все еще иногда используется и по сей день. Но у этой архитектуры есть другая проблема: она не только не работала с изображениями, чей размер отличался от обучающих, но также требовала большого объема памяти, так как преобразование сверточного слоя приводило к передаче в финальные слои большого числа активаций. Следовательно, матрицы весов этих слоев получались огромными.

Эта проблема была решена путем создания *полностью сверточных нейронных сетей*. Хитрость таких сетей заключается в получении среднего значения активаций по всей сверточной сетке. Другими словами, можно просто использовать эту функцию:

```
def avg_pool(x): return x.mean((2,3))
```

Она получает среднее по осям X и Y . Эта функция будет всегда преобразовывать сетку активаций в одну активацию для каждого изображения. PyTorch предоставляет более гибкий модуль `nn.AdaptiveAvgPool2d`, который усредняет сетку активаций в любой определяемый вами размер (хотя мы почти всегда используем размер 1).

Таким образом, полностью сверточная сеть содержит ряд сверточных слоев, некоторые из них будут иметь шаг 2. В конце этого ряда идет адаптивный слой среднего пулинга (выборки средних значений), сглаженный слой для удаления единичных осей и, наконец, линейный слой. Вот наша первая полностью сверточная нейронная сеть:

```
def block(ni, nf): return ConvLayer(ni, nf, stride=2)
def get_model():
    return nn.Sequential(
        block(3, 16),
        block(16, 32),
        block(32, 64),
        block(64, 128),
        block(128, 256),
        nn.AdaptiveAvgPool2d(1),
        Flatten(),
        nn.Linear(256, dls.c))
```

Вскоре реализации `block` будут заменены на другие варианты, поэтому мы больше не называем ее `conv`. Кроме того, мы экономим время, задействовав преимущества `ConvLayer`, который уже предоставляет функциональность `conv` из предыдущей главы и не только.

Разобравшись со сверточными слоями, мы получим активации размера $bs \times ch \times h \times w$ (размер пакета, число каналов, высота, ширина). Нам нужно преобразовать все это

в тензор размера $bs \times ch$, поэтому мы берем среднее по двум последним измерениям и сглаживаем это хвостовое измерение 1×1 , как уже делали в предыдущей модели.



ОСТАНОВИТЕСЬ И ПОДУМАЙТЕ

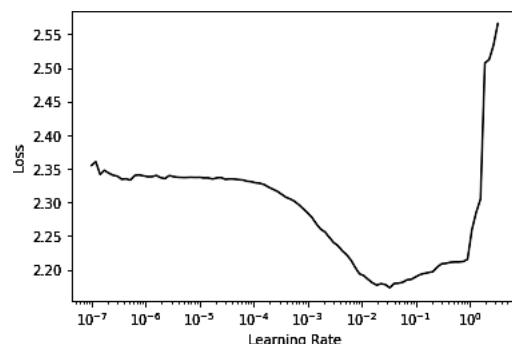
Задумайтесь над следующим вопросом: имел бы этот подход смысл для задачи оптического распознавания символов (optical character recognition — OCR), такой как датасет MNIST? Подавляющее большинство специалистов, работающих с задачами OCR и аналогичными им, склонны использовать полностью сверточные сети, потому что так на сегодня делают практически все. Но, по сути, в этом нет смысла. К примеру, вы не сможете определить, является цифра тройкой или восьмеркой, если разделите ее на небольшие части, перемешаете их и будете решать, похожа ли каждая из них в среднем на тройку или на восьмерку. Но именно это и делает адаптивная выборка средних значений! Полностью сверточные сети отлично подходят только для объектов, не имеющих одного верного расположения или размера (например, большинство естественных фотографий).

Это отличается от стандартной выборки в том смысле, что те слои будут получать среднее (при выборке средних значений) или максимальное (при выборке максимальных значений) из окна заданного размера. Например, слои макс-пулинга размером 2, которые были очень популярны в старых CNN, уменьшают размер изображения наполовину по каждому измерению путем выделения максимального значения из каждого окна 2×2 (с шагом 2).

Как и ранее, мы определяем `Learner` с нашей настраиваемой моделью и затем обучаем ее на предварительно собранных данных:

```
def get_learner(m):
    return Learner(dls, m, loss_func=nn.CrossEntropyLoss(), metrics=accuracy
                  ).to_fp16()

learn = get_learner(get_model())
learn.lr_find()
(0.47863011360168456, 3.981071710586548)
```



Зе-3 зачастую оказывается удачным значением скорости обучения для CNN. И поскольку здесь как раз такой случай, проверим это:

```
learn.fit_one_cycle(5, 3e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	1.901582	2.155090	0.325350	00:07
1	1.559855	1.586795	0.507771	00:07
2	1.296350	1.295499	0.571720	00:07
3	1.144139	1.139257	0.639236	00:07
4	1.049770	1.092619	0.659108	00:07

Неплохое начало, учитывая, что нам нужно выбирать правильный вариант из десяти категорий, и мы начинаем обучение с нуля, выполняя его на протяжении всего пяти эпох. Можно добиться гораздо лучших результатов, задействовав более глубокую модель, но простое нанизывание слоев не особо улучшит итог (можете проверить сами). Чтобы обойти эту проблему, в ResNet был введен принцип *пропускающих соединений*. В следующем разделе мы познакомимся с этим и другими аспектами данной архитектуры.

Построение современной CNN: ResNet

Теперь у нас есть все элементы, необходимые для построения моделей, которые мы использовали в задачах компьютерного зрения с самого начала книги: ResNet. Мы поясним основную лежащую в их основе идею и покажем, как данная архитектура повышает точность, на примере Imagenette по сравнению с нашей предыдущей моделью, а потом перейдем к созданию версии со всеми последними доработками.

Пропускающие соединения

В 2015 году авторы работы, посвященной ResNet, заметили, кое-что, как им показалось, любопытное. Они обнаружили, что даже после использования пакетной нормализации сеть, использующая больше слоев, справлялась хуже, чем ее альтернатива с их меньшим количеством, при этом ничем другим они между собой не отличались. Более же интересно то, что разницу заметили не только на контрольной выборке, но также и на обучающей. Так что это была не просто проблема обобщаемости, а проблема обучения. Вот как это отражено в самой работе:

Вопреки ожиданиям, это ухудшение не явилось следствием переобучения, и добавление большего числа слоев для подходящих глубоких моделей действительно ведет к увеличению ошибки обучения, как [ранее сообщалось] и тщательно проверялось рядом экспериментов.

Этот феномен был проиллюстрирован графиком, приведенным на рис. 14.1, где ошибка обучения показана слева, а ошибка тестовой выборки — справа.

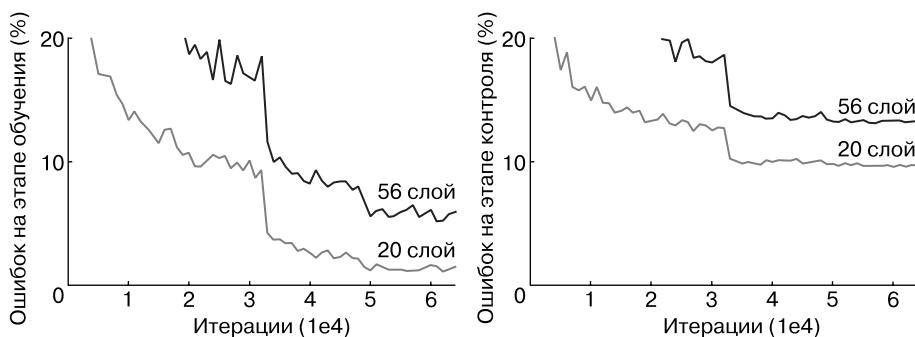


Рис. 14.1. Обучение сетей разной глубины
(изображение любезно предоставлено Каймингом Хе и др.)

Как отмечают авторы, они не первые заметили этот занятный факт, но первыми совершили очень важный прорыв:

Рассмотрим малослойную архитектуру и ее углубленный аналог, добавляющий несколько слоев. Существует решение для построения более глубокой модели: добавляемые слои являются тождественным отображением, а остальные копируются из обученной малослойной модели.

Поскольку это академическая работа, то процесс описан весьма непонятным языком, но принцип здесь, по сути, прост: начать с хорошо обученной нейронной сети из 20 слоев, добавить поверх них еще 36 слоев, которые вообще ничего не делают (например, это могут быть линейные слои с одним весом, равным 1, и смещением, равным 0). В результате мы получим сеть из 56 слоев, которая делает в точности то, что и изначальная сеть из двадцати. Это доказывает, что всегда существуют глубокие нейронные сети, которые должны быть *по меньшей мере так же хороши*, как любая малослойная сеть, но по какой-то причине SGD не может их обнаружить.



ТЕРМИН: ТОЖДЕСТВЕННОЕ ОТБРАЖЕНИЕ

Возвращение входных данных без изменений. Выполняет этот процесс *функция тождества*.

Фактически есть и другой, гораздо более интересный, способ создания этих дополнительных 36 слоев. Что, если мы заменим все вхождения $\text{conv}(x)$ на $x + \text{conv}(x)$, где conv — это функция из предыдущей главы, добавляющая вторую свертку, затем слой пакетной нормализации, а после этого ReLU? Более того, вспомните, что пакетная нормализация выполняет операцию $\gamma y + \beta$. Что, если инициализировать γ нулем для каждого из финальных batchnorm-слоев? Поскольку β уже инициализирована нулем, то наша $\text{conv}(x)$ для тех дополнительных 36 слоев всегда будет равна 0, поэтому $x + \text{conv}(x)$ всегда будет равняться x .

Что это нам дало? Суть в том, что те 36 дополнительных слоев сами по себе представляют *тождественное отображение*, но у них есть *параметры*, а значит, они *обучаемы*. Итак, мы можем начать с нашей обученной 20-слойной модели, добавить эти 36 слоев, которые изначально вообще ничего не делают, а затем *тонко настроить всю 56-слойную модель*. После этого добавленные 36 слоев могут обучить параметры, которые сделают их максимально полезными.

В научной работе по ResNet предлагается вариация этого подхода, подразумевающая «пропуск» каждой второй свертки, в результате чего мы получаем $x + \text{conv2}(\text{conv1}(x))$. Это показано на диаграмме на рис. 14.2 (взятой из данной работы).

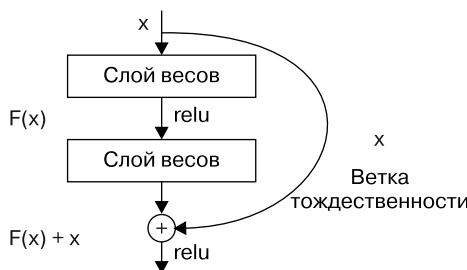


Рис. 14.2. Простой блок ResNet (изображение любезно предоставлено Каймингом Хе и др.)

Стрелка справа — это просто часть x из $x + \text{conv2}(\text{conv1}(x))$, которая называется *веткой тождественности*, или *пропускающим соединением*. Путь слева от нее — это часть $\text{conv2}(\text{conv1}(x))$. Можете рассматривать путь тождественности как предоставляющий прямой маршрут от входа к выходу.

В ResNet мы не начинаем с обучения небольшого числа слоев, добавляя в итоге дополнительные и производя тонкую настройку в конце. Вместо этого мы на протяжении всей CNN используем блоки ResNet, подобные показанному на рис. 14.2, инициализируя их с нуля обычным образом и также стандартным способом обучая их через SGD. Облегчения обучения посредством SGD мы добиваемся с помощью пропускающих соединений.

Блоки ResNet можно рассмотреть и с другой (по сути, эквивалентной) стороны. Вот как об этом говорится в работе:

Вместо расчета на то, что каждая пара соседних слоев напрямую согласуется с желаемым базовым отображением, мы явно позволяем этим слоям согласоваться с остаточным отображением. Формально обозначая желаемое базовое отображение как $H(x)$, мы позволяем стеку нелинейных слоев согласоваться с другим отображением $F(x) := H(x) - x$. Исходное отображение приводится к $F(x) + x$. Предполагается, что легче оптимизировать остаточное отображение, чем исходное, не имеющее ссылок. В крайнем случае, если бы тождественное отображение являлось оптимальным, то было бы проще привести остаточное к нулю, чем согласовывать тождественное стеком нелинейных слоев.

Опять же, это слишком сложно для понимания — попробуем перевести на понятный язык. Если выход заданного слоя — это x и мы используем блок ResNet, возвращающий $y = x + \text{block}(x)$, то мы не просим этот блок спрогнозировать y . Мы просим его дать прогноз разности между y и x . Поэтому работа данных блоков не в прогнозировании конкретных признаков, а в минимизации ошибки между x и желаемым y . Таким образом, ResNet хорошо изучает небольшие отличия между бездействием и прохождением через блок из двух сверточных слоев (с обучаемыми весами). Именно здесь и кроется суть наименования этих моделей: они прогнозируют остатки (напоминание: «остаток» — это прогноз минус цель).

В обоих случаях подразумевается один общий ключевой принцип — простота обучения. Это важная мысль. Вспомним теорему об универсальной аппроксимации, утверждающую, что достаточно крупная сеть может выучить что угодно. Это по-прежнему так, но оказывается, что есть очень существенное отличие между тем, что сеть *может выучить* в принципе, и тем, что *ей будет легко выучить* на основе реалистичных данных и режимов обучения. Многие достижения в области нейронных сетей за последние десять лет, аналогично блоку ResNet, стали следствием обнаружения способа реально осуществить то, что всегда считалось только возможным.



ИСТИННЫЙ ПУТЬ ТОЖДЕСТВЕННОСТИ

В оригинальной работе не действовался прием с использованием нуля для начального значения гамма в последнем слое пакетной нормализации каждого блока. Это решение было разработано несколькими годами позже. Поэтому в начальной версии ResNet обучение не начиналось с истинного пути тождественности, пролегающего через блоки ResNet, но при этом наличие возможности «пройти через» пропускающие соединения позволяла улучшить обучение. Добавление приема с инициализацией гамма дало возможность обучать модели с еще большими скоростями.

Вот определение простого блока ResNet (fastai инициализирует веса `gamma` последнего слоя пакетной нормализации с нулем из-за `norm_type=NormType.BatchZero`):

```
class ResBlock(Module):
    def __init__(self, ni, nf):
        self.convs = nn.Sequential(
            ConvLayer(ni,nf),
            ConvLayer(nf,nf, norm_type=NormType.BatchZero))

    def forward(self, x): return x + self.convs(x)
```

Но здесь есть две проблемы: он не может обработать шаг больше 1 и требует, чтобы `ni=nf`. Приостановитесь и подумайте, почему это так.

Сложность в том, что при шаге, например, 2 в одной из сверток размер сетки выходных активаций будет вдвое меньше по каждой оси входа. Поэтому мы не сможем добавить ее обратно к `x` в `forward`, так как `x` и выходные активации будут иметь разные измерения. Та же проблема возникает, если `ni!=nf`, — разные формы входных и выходных связей не позволят сложить их вместе.

Чтобы это исправить, нам потребуется изменить форму `x` для соответствия результату `self.convs`. Уменьшение размера сетки вдвое реализуется с помощью слоя среднего пулинга с шагом 2, то есть слоя, который получает участки 2×2 из входа и заменяет их выбранными средними значениями.

Изменение числа каналов осуществляется с помощью свертки. Тем не менее нам нужно, чтобы это пропускающее соединение было максимально близко к карте тождественности, что подразумевает максимальное упрощение данной свертки. Самый простой вариант — это свертка, размер ядра которой равен 1, то есть будет записан как $ni \times nf \times 1 \times 1$, в результате чего оно будет получать только скалярное произведение по каналам каждого входного пикселя, не совмещая полученные для пикселей значения. Этот вид свертки 1×1 широко используется в современных CNN, так что не пожалейте времени на осмысление принципа ее действия.



ТЕРМИН: СВЕРТКА 1×1

Свертка, размер ядра которой равен 1.

Бот ResBlock, использующий эти приемы для обработки изменения формы в пропускающем соединении:

```
def _conv_block(ni,nf,stride):
    return nn.Sequential(
        ConvLayer(ni, nf, stride=stride),
```

```

    ConvLayer(nf, nf, act_cls=None, norm_type=NormType.BatchZero))

class ResBlock(Module):
    def __init__(self, ni, nf, stride=1):
        self.convs = _conv_block(ni, nf, stride)
        self.idconv = noop if ni==nf else ConvLayer(ni, nf, 1, act_cls=None)
        self.pool = noop if stride==1 else nn.AvgPool2d(2, ceil_mode=True)

    def forward(self, x):
        return F.relu(self.convs(x) + self.idconv(self.pool(x)))

```

Обратите внимание, что здесь мы задействуем функцию `noop`, которая просто возвращает свой вход в неизменном виде (`noop` — это термин информатики, означающий «отсутствие операций»). В данном случае `idconv` ничего не делает, если `ni=nf`, а `pool` ничего не делает, если `stride=1`, что нам и требовалось для пропускающего соединения.

Вы также увидите, что мы удалили ReLU (`act_cls=None`) из финальной свертки в `convs` и из `idconv`, переместив ее в область *после* добавления пропускающего соединения. Смысл в том, что весь блок ResNet подобен слою, а вам нужно, чтобы активации шли после вашего слоя.

Заменим `block` на `ResBlock` и посмотрим:

```

def block(ni,nf): return ResBlock(ni, nf, stride=2)
learn = get_learner(get_model())

learn.fit_one_cycle(5, 3e-3)

```

epoch	train_loss	valid_loss	accuracy	time
0	1.973174	1.845491	0.373248	00:08
1	1.678627	1.778713	0.439236	00:08
2	1.386163	1.596503	0.507261	00:08
3	1.177839	1.102993	0.644841	00:09
4	1.052435	1.038013	0.667771	00:09

Не намного лучше. Но весь смысл был в том, чтобы получить возможность обувать *более глубокие* модели, преимуществом которых пока что мы не пользуемся. Для создания модели, которая, к примеру, вдвое глубже, нам нужно лишь заменить `block` на два `ResBlock` подряд:

```

def block(ni, nf):
    return nn.Sequential(ResBlock(ni, nf, stride=2), ResBlock(nf, nf))

learn = get_learner(get_model())
learn.fit_one_cycle(5, 3e-3)

```

epoch	train_loss	valid_loss	accuracy	time
0	1.964076	1.864578	0.355159	00:12
1	1.636880	1.596789	0.502675	00:12
2	1.335378	1.304472	0.588535	00:12
3	1.089160	1.065063	0.663185	00:12
4	0.942904	0.963589	0.692739	00:12

Вот теперь прогресс очевиден!

Авторы работы, посвященной ResNet, победили в соревновании ImageNet 2015 года. В то время это было важнейшее ежегодное событие в области компьютерного зрения. Мы с вами уже знакомились с другими победителями ImageNet 2013 года — Зейлером и Фергусом. Интересно заметить, что в обоих случаях отправными точками научных прорывов послужили экспериментальные наблюдения. В случае Фергуса и Зейлера это были наблюдения за тем, что по факту изучают слои, а в случае авторов ResNet — за тем, какие виды сетей можно обучить. Эта способность разрабатывать и анализировать эксперименты, с интересом реагируя даже просто на неожиданные результаты, а самое главное, старательно и упорно докапываться после этого до сути происходящего, лежит в основе многих научных открытий. Глубокое обучение не похоже на чистую математику. Это по большей части область экспериментирования, поэтому здесь важно уделять много времени именно практике, а не теоретизированию.

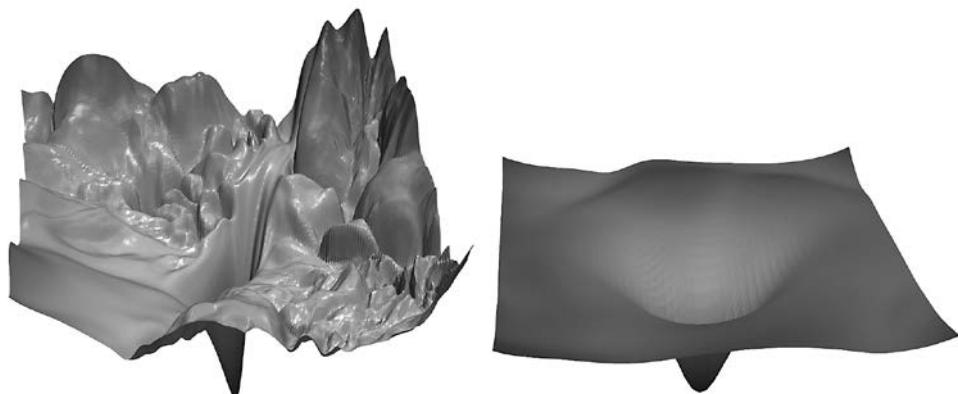


Рис. 14.3. Влияние ResNet на ландшафт потерь
(изображение любезно предоставлено Хао Ли и др.)

С момента своего появления ResNet широко изучалась и применялась во многих направлениях. Одна из наиболее интересных работ по этой части — *Visualizing the Loss Landscape of Neural Nets* (<https://oreil.ly/C9cFi>) («Визуализация ландшафта потерь нейронных сетей») была представлена Хао Ли (Naoh Li) и др. в 2018 году. В ней показывается, что использование пропускающих соединений помогает сгладить функцию потерь, что упрощает обучение и помогает избежать провала в области резких углублений. На рис. 14.3 приводится завораживающая картина из этой работы, иллюстрирующая разницу между выпуклым рельефом, по которому SGD необходимо пройти для оптимизации стандартной CNN (слева) и гладкой поверхностью ResNet (справа).

Наша первая модель уже эффективна, но дополнительные исследования обнаружили возможность применения и других приемов, которые сделают ее еще лучше. Рассмотрим теперь их.

Актуальная ResNet

В работе *Bag of Tricks for Image Classification with Convolutional Neural Networks* (<https://oreil.ly/n-qhd>) («Набор приемов для классификации изображений с помощью сверточных нейронных сетей») Тонг Хе (Tong He) и др. рассматривают вариации архитектуры ResNet, реализация которых не требует лишних затрат в плане количества параметров и объема вычислений. Используя доработанную ResNet-50 совместно с Mixup, авторам удалось добиться на датасете ImageNet точности топ-5 94,6 %, что можно сравнить с 92,2 %, полученными посредством стандартной ResNet-50 без Mixup. Этот результат превосходит получаемый обычными моделями ResNet, которые вдвое глубже (а также вдвое медленнее и более склонны к переобучению).



ТЕРМИН: ТОЧНОСТЬ ТОП-5

Метрика, оценивающая, как часто нужная нам метка попадает в топ-5 прогнозов модели. Она использовалась в соревновании ImageNet, так как многие изображения содержали по несколько объектов либо такие объекты, которые можно было легко спутать или которые имели ошибочно присвоенные похожие метки. В таких ситуациях рассмотрение точности топ-1 может быть неподходящим вариантом. Хотя с недавних пор CNN настолько усовершенствовались, что стали достигать точности топ-5 около 100 %, поэтому некоторые исследователи стали использовать точность топ-1 и для ImageNet.

Мы будем применять эту доработанную версию в процессе масштабирования до полной ResNet, потому что она существенно лучше. Ее отличие от предыдущей реализации в том, что начинается она не с блоков ResNet, а с нескольких сверточных слоев, сопровождаемых слоем макс-пулинга. Вот как выглядят первые слои, называемые *опорой* нейронной сети:

```

def _resnet_stem(*sizes):
    return [
        ConvLayer(sizes[i], sizes[i+1], 3, stride = 2 if i==0 else 1)
        for i in range(len(sizes)-1)
    ] + [nn.MaxPool2d(kernel_size=3, stride=2, padding=1)]
_resnet_stem(3,32,32,64)

[ConvLayer(
    (0): Conv2d(3, 32, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1))
    (1): BatchNorm2d(32, eps=1e-05, momentum=0.1)
    (2): ReLU()
), ConvLayer(
    (0): Conv2d(32, 32, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (1): BatchNorm2d(32, eps=1e-05, momentum=0.1)
    (2): ReLU()
), ConvLayer(
    (0): Conv2d(32, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (1): BatchNorm2d(64, eps=1e-05, momentum=0.1)
    (2): ReLU()
), MaxPool2d(kernel_size=3, stride=2, padding=1, ceil_mode=False)]

```



ТЕРМИН: ОПОРА

Первые несколько слоев CNN. Обычно структура опоры отличается от основного тела CNN.

Использование для опоры простых сверточных слоев вместо блоков ResNet обусловлено одной важной чертой всех глубоких сверточных сетей, а именно тем, что большая часть вычислений производится именно в ранних слоях. Поэтому следует делать эти слои максимально быстрыми и простыми.

Чтобы увидеть, почему на ранних слоях происходит столько вычислений, рассмотрите самую первую свертку для входного изображения размером 128 пикселей. Если это будет свертка с шагом 1, то она применит ядро к каждому из 128 пикселей, а это немало работы! В более поздних слоях размер сетки должен стать 4×4 или даже 2×2 , а значит, и операций применения ядра будет намного меньше.

С другой стороны, у свертки первого слоя есть только три входных признака и 32 выходных. Поскольку это ядро 3×3 , то параметров в весах получится $3 \times 32 \times 3 \times 3 = 864$. Но в последней свертке будет 256 входных признаков и 512 выходных, что означает 1 179 648 весов! Поэтому первые слои содержат основную часть вычислений, а последние — основную часть параметров.

Блок ResNet производит больше вычислений, чем простой сверточный блок, поскольку (в случае с шагом 2) блок ResNet содержит три свертки и слой пулинга. Именно поэтому нам нужно начинать ResNet с простых сверток.

Теперь мы готовы показать вам реализацию современной ResNet с «набором приемов». Она задействует четыре группы блоков ResNet с 64, 128, 256, а затем

512 фильтрами. Каждая группа начинается с блока с шагом 2, что не относится к первой, так как она идет после слоя MaxPooling:

```
class ResNet(nn.Sequential):
    def __init__(self, n_out, layers, expansion=1):
        stem = _resnet_stem(3,32,32,64)
        self.block_szs = [64, 64, 128, 256, 512]
        for i in range(1,5): self.block_szs[i] *= expansion
        blocks = [self._make_layer(*o) for o in enumerate(layers)]
        super().__init__(*stem, *blocks,
                         nn.AdaptiveAvgPool2d(1), Flatten(),
                         nn.Linear(self.block_szs[-1], n_out))

    def _make_layer(self, idx, n_layers):
        stride = 1 if idx==0 else 2
        ch_in, ch_out = self.block_szs[idx:idx+2]
        return nn.Sequential([
            ResBlock(ch_in if i==0 else ch_out, ch_out, stride if i==0 else 1)
            for i in range(n_layers)
        ])
```

Функция `_make_layer` отвечает за создание серии блоков `n_layers`. Первый охватывает промежуток от `ch_in` до `ch_out` с указанным `stride`, а все остальные применяют шаг 1 с тензорами от `ch_out` до `ch_out`. После определения блоков модель является чисто последовательной, поэтому мы определяем ее как подкласс `nn.Sequential`. (Пока не обращайте внимания на параметр `expansion`; о нем мы будем говорить в следующем разделе. Сейчас он равен 1, то есть влияния оказывать не будет.)

Различные версии моделей (ResNet-18, -34, -50 и т. д.) просто изменяют количество блоков в каждой из этих групп. Вот определение ResNet-18:

```
rn = ResNet(dls.c, [2,2,2,2])
```

Немного ее обучим и посмотрим, как она покажет себя в сравнении с предыдущей моделью:

```
learn = get_learner(rn)
learn.fit_one_cycle(5, 3e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	1.673882	1.828394	0.413758	00:13
1	1.331675	1.572685	0.518217	00:13
2	1.087224	1.086102	0.650701	00:13
3	0.900428	0.968219	0.684331	00:12
4	0.760280	0.782558	0.757197	00:12

Несмотря на то что здесь у нас больше каналов (следовательно, модель более точна), благодаря оптимизированной опоре обучение выполняется так же быстро, как и раньше.

Чтобы углубить модель, не задействуя слишком много вычислений или памяти, можно использовать другой вид слоя, предложенный в работе ResNet для сетей с глубиной в 50 и более: зауженный слой.

Зауженные слои

Вместо составления двух сверточных слоев с ядром размером 3 зауженные слои используют три свертки: две 1×1 (в начале и в конце) и одну 3×3 . Схематично это отражено на рис. 14.4.

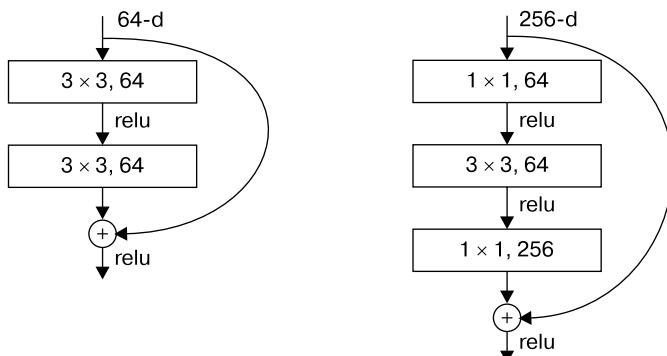


Рис. 14.4. Сравнение стандартного и зауженного блоков ResNet (изображение любезно предоставлено Каймингом Хе и др.)

В чем же его польза? Свертки размером 1×1 намного быстрее, поэтому даже если эта структура выглядит сложнее, выполняется этот блок быстрее, чем первый блок ResNet, который мы видели. Это, в свою очередь, дает возможность использовать больше фильтров: как видно из картинки, количество фильтров входа и выхода в четыре раза больше (256 вместо 64). Свертки 1×1 сначала уменьшают, а затем восстанавливают количество каналов (отсюда и название *зауженный*). Их преимущество в том, что мы за тот же промежуток времени можем использовать больше фильтров.

Заменим `ResBlock` на такую зауженную альтернативу:

```
def _conv_block(ni, nf, stride):
    return nn.Sequential(
```

```
ConvLayer(ni, nf//4, 1),
ConvLayer(nf//4, nf//4, stride=stride),
ConvLayer(nf//4, nf, 1, act_cls=None, norm_type=NormType.BatchZero))
```

Мы задействуем это для создания ResNet-50 с размерами групп (3, 4, 6, 3). Теперь нужно передать в параметр `expansion` нашей ResNet значение 4, поскольку вначале число каналов необходимо уменьшить в четыре раза, после чего в конце их количество будет обратно увеличено в те же четыре раза.

Глубокие сети, подобные этой, обычно не демонстрируют улучшений при обучении в течение всего пяти эпох, поэтому мы увеличим их количество до двадцати, чтобы добиться от нашей увеличенной модели максимума. Для повышения эффективности используем также и увеличенные изображения:

```
dls = get_data(URLs.IMAGENETTE_320, presize=320, resize=224)
```

Нам не нужно ничего делать, чтобы учесть более крупные изображения в 224 пикселя. Благодаря полностью сверточной сети все уже работает. Это также позволило нам использовать *постепенное изменение размера* ранее — задействованные нами тогда модели были сверточными, поэтому мы могли тонко настраивать даже те из них, которые обучались на других размерах изображений. Теперь можно обучить модель и оценить эффект:

```
rn = ResNet(dls.c, [3,4,6,3], 4)
```

```
learn = get_learner(rn)
learn.fit_one_cycle(20, 3e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	1.613448	1.473355	0.514140	00:31
1	1.359604	2.050794	0.397452	00:31
2	1.253112	4.511735	0.387006	00:31
3	1.133450	2.575221	0.396178	00:31
4	1.054752	1.264525	0.613758	00:32
5	0.927930	2.670484	0.422675	00:32
6	0.838268	1.724588	0.528662	00:32
7	0.748289	1.180668	0.666497	00:31
8	0.688637	1.245039	0.650446	00:32
9	0.645530	1.053691	0.674904	00:31

epoch	train_loss	valid_loss	accuracy	time
10	0.593401	1.180786	0.676433	00:32
11	0.536634	0.879937	0.713885	00:32
12	0.479208	0.798356	0.741656	00:32
13	0.440071	0.600644	0.806879	00:32
14	0.402952	0.450296	0.858599	00:32
15	0.359117	0.486126	0.846369	00:32
16	0.313642	0.442215	0.861911	00:32
17	0.294050	0.485967	0.853503	00:32
18	0.270583	0.408566	0.875924	00:32
19	0.266003	0.411752	0.872611	00:33

Вот это уже очень хороший результат! Попробуйте добавить к этому Mixup и обучить модель на протяжении ста эпох, пока будете обедать. Таким образом вы получите очень точный классификатор изображений, обученный с нуля.

Показанная здесь структура с заужением обычно используется только в моделях ResNet-50, -101 и -152. В моделях ResNet-18 и -34 чаще задействуют структуру без заужения, как мы видели в предыдущем разделе. Тем не менее мы заметили, что зауженный слой обычно работает лучше даже для малослойных сетей. Это просто показывает, что мелкие детали из научных работ могут применяться на протяжении многих лет, не будучи доведенными до ума. Ставить под сомнение предположения и понятия, которые «все знают», — правильно, потому что эта область все еще плохо изучена и множество деталей нередко оказываются недоработанными.

Резюме

Теперь вы узнали, как строятся модели компьютерного зрения, которые мы изучали с самой первой главы. Вам также известно, что применение пропускающих соединений позволяет обучать их более глубокие аналоги. Несмотря на то что в поиске наилучших архитектур было произведено немало исследований, они все так или иначе задействуют этот прием построения прямого пути от входа к концу сети. При использовании переноса обучения ResNet выступает в роли предварительно обученной модели. В следующей главе мы рассмотрим заключительные детали создания на ее основе использованных нами моделей.

Вопросник

1. Как мы получили один вектор активаций в CNN, использованной для MNIST в предыдущих главах? Почему этот вариант не подходит для Imagenette?
2. Что мы вместо этого делаем для Imagenette?
3. Что такое адаптивный пулинг?
4. Что такое средний пулинг?
5. Зачем нужен `Flatten` после адаптивного слоя среднего пулинга?
6. Что такое пропускающее соединение?
7. Почему пропускающие соединения позволяют обучать более глубокие модели?
8. Что показано на рис. 14.1? Как это привело к идее пропускающих соединений?
9. Что такое тождественное отображение?
10. Каково базовое уравнение для блока ResNet (без учета слоев пакетной нормализации и ReLU)?
11. Какое отношение ResNet имеют к остаткам?
12. Как мы поступаем с пропускающим соединением при использовании свертки с шагом 2? А как действуем в случае изменения количества фильтров?
13. Как можно выразить свертку 1×1 в виде скалярного произведения векторов?
14. Создайте свертку 1×1 с `F.conv2d` или `nn.Conv2d` и примените ее к изображению. Как изменится форма этого изображения?
15. Что возвращает функция `poop`?
16. Поясните, что показано на рис. 14.3.
17. Когда метрика точности топ-5 лучше, чем точность топ-1?
18. Что такое «опора» CNN?
19. Почему в опоре CNN мы используем простые свертки вместо блоков ResNet?
20. Чем отличается зауженный блок от простого блока ResNet?
21. Почему зауженный блок быстрее?
22. Каким образом полностью сверточные сети (и сети с адаптивным пулингом в общем) допускают постепенное изменение размера?

Дополнительные задания

1. Попробуйте создать полностью сверточную сеть с аддитивным средним пулингом для MNIST (обратите внимание, что вам потребуется меньше слоев с шагом 2). Как она сопоставляется с сетью без слоя пулинга?
2. В главе 17 мы поясняем *нотацию суммирования Эйнштейна*. Забегите вперед и ознакомьтесь с этой темой, чтобы понять, как эта нотация работает, после чего напишите реализацию операции свертки 1×1 , используя `torch.einsum`. Сравните ее с такой же операцией, использующей `torch.conv2d`.
3. Напишите функцию точности топ-5 с помощью только PyTorch или только Python.
4. Обучите модель на Imagenette на протяжении большего числа эпох со сглаживанием меток и без него. Загляните в таблицу лидеров по этому датасету и оцените, насколько близко вы смогли подобраться к лучшим результатам. Прочтите привязанные к лидирующим результатам страницы, где описываются ведущие подходы.

ГЛАВА 15

Архитектуры приложений

Сейчас мы находимся в очень удачном положении, так как можем полноценно понимать архитектуры, которые использовали для эталонных моделей компьютерного зрения, обработки естественного языка и табличного анализа. В текущей главе мы заполним все недостающие детали, касающиеся функционирования моделей, а также научим вас их создавать.

Помимо этого, мы вернемся к настраиваемому конвейеру предварительной обработки, который видели в главе 11 при работе с сиамскими сетями, и покажем, как использовать его компоненты в библиотеке fastai, чтобы создавать собственные предварительно обученные модели для новых задач.

Начнем с компьютерного зрения.

Компьютерное зрение

Для построения моделей компьютерного зрения мы задействуем функции `cnn_learner` и `unet_learner` в зависимости от задачи. В этом разделе мы узнаем, как создавать объекты `Learner`, которые использовали в частях I и II.

`cnn_learner`

Рассмотрим, что происходит, когда мы применяем функцию `cnn_learner`. Начнем с передачи в эту функцию архитектуры для использования в качестве *тела* сети. В большинстве случаев мы задействуем ResNet, которую вы уже умеете создавать, поэтому нам нет необходимости дополнительного в нее углубляться. Предварительно обученные веса скачиваются согласно требованиям и загружаются в ResNet.

Затем для использования в переносе обучения данную сеть необходимо *обрезать*. Этот процесс подразумевает отделение последнего слоя, отвечающего за категоризацию, актуальную только для ImageNet. В действительности нам

нужно отрезать не только этот слой, но и все, начиная с адаптивного слоя пулинга и далее. Причина этого вскоре станет понятна. Поскольку разные архитектуры могут задействовать разные типы слоев пулинга или даже совершенно другие виды *вершин*, мы не просто ищем адаптивный слой подвыборки, чтобы решить, какую часть модели отделить. Вместо этого у нас есть словарь информации, с помощью которого мы для каждой модели определяем, где заканчивается ее тело и начинается вершина. Этот словарь мы называем `model_meta`, и вот его вариант для `resnet50`:

```
model_meta[resnet50]
{'cut': -2,
 'split': <function fastai.vision.learner._resnet_split(m)>,
 'stats': ([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])}
```



ТЕРМИНЫ: ТЕЛО И ВЕРШИНА

Вершина нейронной сети — это часть, специализирующаяся на конкретной задаче. В случае CNN эта часть обычно идет после адаптивного слоя среднего пулинга. Тело, в свою очередь, представляет всю оставшуюся часть и включает опору (о которой мы узнали в главе 14).

Если мы возьмем все слои, предшествующие точке отделения `-2`, то получим часть модели, которую `fastai` сохраняет для переноса обучения. Далее идет этап прикрепления новой вершины, которая создается с помощью функции `create_head`:

```
create_head(20, 2)
Sequential(
    (0): AdaptiveConcatPool2d(
        (ap): AdaptiveAvgPool2d(output_size=1)
        (mp): AdaptiveMaxPool2d(output_size=1)
    )
    (1): Flatten()
    (2): BatchNorm1d(20, eps=1e-05, momentum=0.1, affine=True)
    (3): Dropout(p=0.25, inplace=False)
    (4): Linear(in_features=20, out_features=512, bias=False)
    (5): ReLU(inplace=True)
    (6): BatchNorm1d(512, eps=1e-05, momentum=0.1, affine=True)
    (7): Dropout(p=0.5, inplace=False)
    (8): Linear(in_features=512, out_features=2, bias=False)
)
```

С помощью этой функции можно выбирать, сколько дополнительных слоев добавлять в конце, величину дропаута после каждого из них, а также тип используемого пулинга. По умолчанию `fastai` применяет и средний пулинг, и макс-пулинг, конкatenируя их (слой `AdaptiveConcatPool2d`). Это не особо распространенный подход, но он был недавно разработан независимо в `fastai` и других

исследовательских лабораториях с целью предоставить некоторые улучшения по сравнению с использованием только среднего пулинга.

Fastai отличается от большинства библиотек тем, что по умолчанию добавляет в вершину CNN два линейных слоя, а не один. Причина в том, что перенос обучения может сохранять свою эффективность даже, как мы видели, при адаптации предварительно обученной модели под совершенно другую область задачи. Тем не менее простого использования одного линейного слоя в таких случаях будет недостаточно. Мы обнаружили, что применение двух таких слоев позволит задействовать перенос обучения быстрее и проще в большем спектре ситуаций.



ОДНА ПОСЛЕДНЯЯ ПАКЕТНАЯ НОРМАЛИЗАЦИЯ

Стоит рассмотреть параметр для `create_head`, называемый `bn_final`. Его определение как `True` приведет к добавлению слоя пакетной нормализации в качестве финального. Это может пригодиться для правильного масштабирования модели под выходные активации. Мы пока что не встречали этот подход в публикациях, но на практике во всех испробованных нами ситуациях он работал отлично.

Рассмотрим, какую роль выполняла `unet_learner` в задаче сегментации, приведенной в главе 1.

unet_learner

В главе 1 для сегментации изображений мы задействовали одну из самых интересных архитектур в глубоком обучении. Сегментация — непростая задача, так как нужный выход является изображением или сеткой пикселей, содержащей спрогнозированную метку для каждого пикселя. Другие задачи применяют аналогичную структуру, например, повышая разрешение изображения (*сверхвысокое разрешение*), добавляя цвет в черно-белое изображение (*колоризация*) или преобразуя фото в синтетический рисунок (*перенос стиля*), — эти задачи рассматриваются в онлайн-главе книги (<https://book.fast.ai>), так что обязательно ознакомьтесь с ней после текущего материала. В каждом случае мы начинаем с изображения и преобразуем его в другое изображение с теми же измерениями или соотношением сторон, но с уже измененными определенным образом пикселями. Такие модели мы называем *генеративными зрительными моделями*.

Реализовать это можно, начав с того же подхода по разработке вершины CNN, который мы рассматривали в предыдущей главе. В качестве примера мы начнем с ResNet и отрежем аддитивный слой пулинга вместе со всем последующим содержимым. Затем удаленные слои мы заменим на подготовленную вершину, выполняющую генеративную задачу.

В последнем предложении многое непонятно. Как вообще создать вершину CNN, генерирующую изображение? Если мы начнем, например, с входного изображения размером 224 пикселя, тогда в конце тела ResNet у нас будет сетка сверточных активаций размером 7×7 . Как мы можем преобразовать ее в маску для сегментации размером 224 пикселя?

Конечно же, мы это делаем с помощью нейронной сети! Итак, нам нужен слой, способный увеличивать размер сетки в CNN. Простым подходом будет заменить каждый пиксель сетки 7×7 на четыре пикселя квадрата 2×2 . Каждый из этих четырех пикселей будет иметь одно значение — это называется *интерполяцией методом ближайшего соседа* (ступенчатой интерполяцией). PyTorch предоставляет слой, который делает это за нас, поэтому одним из вариантов будет создать вершину, содержащую сверточные слои с шагом 1 (вместе со слоями пакетной нормализации и ReLU) вперемежку со слоями ступенчатой интерполяции. Вообще-то, вы можете попробовать это прямо сейчас. Попробуйте создать собственную вершину, спроектированную подобным образом, и опробовать ее в задаче по сегментации CamVid. У вас должны получиться вполне успешные результаты, хотя и не такие замечательные, как полученные нами в главе 1.

Второй вариант — заменить комбинацию ступенчатой интерполяции и свертки на *транспонированную свертку*. Она идентична обычной свертке, но начальное заполнение нулями происходит между всеми пикселями входного изображения. Взгляните на рис. 15.1, где приводится пример из замечательной работы, посвященной арифметике сверток (<https://oreil.ly/hu06c>), о которой мы говорили в главе 13. Здесь вы видите транспонированную свертку размером 3×3 , примененную к изображению 3×3 .

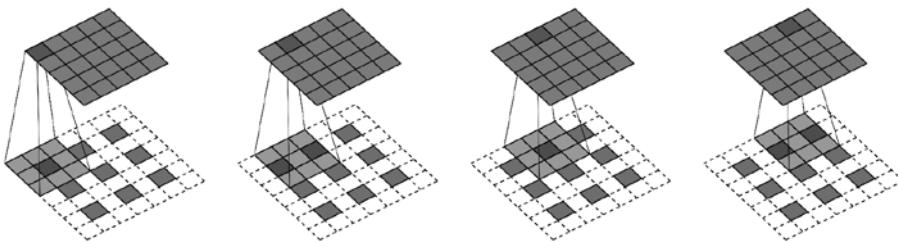


Рис. 15.1. Транспонированная свертка (изображение любезно предоставлено Винсентом Дюмулином и Франческо Визином)

Здесь мы видим, что результатом становится увеличение размера входного изображения. Можете попробовать это прямо сейчас, используя fastai-класс `ConvLayer`. Передайте параметр `transpose=True`, чтобы создать в вершине транспонированную свертку вместо стандартной.

Как бы то ни было, ни один из перечисленных подходов не идеален. Проблема в том, что наша сетка 7×7 элементарно не содержит достаточно информации

для создания вывода размером 224×224 пикселя. Чтобы получить необходимый объем информации для полноценной регенерации каждого пикселя на выходе, требуется огромное множество активаций каждой ячейки этой сетки.

Решением будет использовать *пропускающие соединения*, как и в ResNet, но начинать пропуск с активаций в теле ResNet и следовать обратно к активациям транспонированной свертки на противоположном конце архитектуры. Этот подход, показанный на рис. 15.2, был представлен в 2015 году Олафом Роннебергером (Olaf Ronneberger) и др. в работе *U-Net: Convolutional Networks for Biomedical Image Segmentation* (<https://oreil.ly/6ely4>) («U-Net: сверточные сети для сегментации биомедицинских изображений»). Несмотря на то что эта работа ориентирована на медицинскую сферу, U-Net произвела революцию в моделях генеративного зрения.

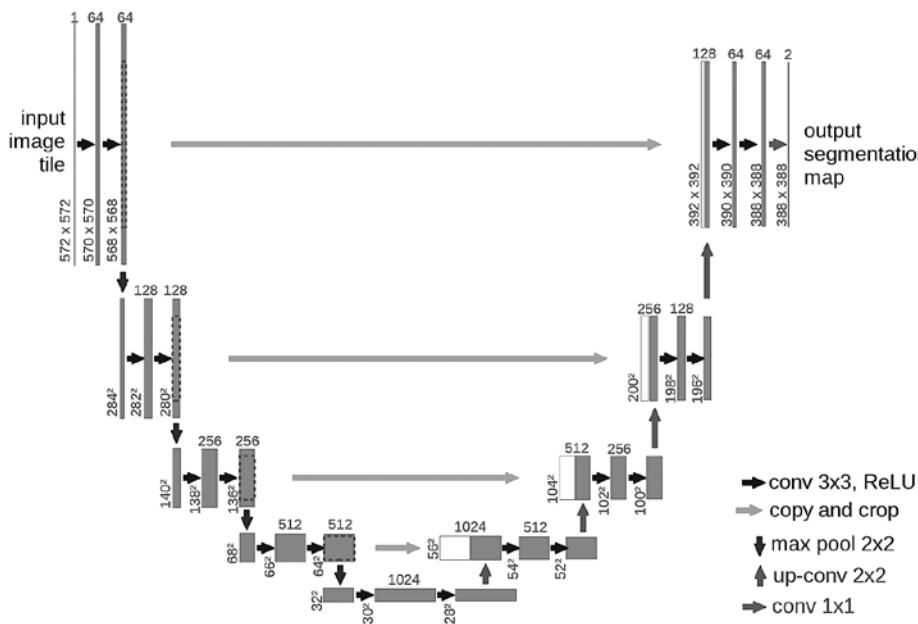


Рис. 15.2. Архитектура U-Net (изображение любезно предоставлено Олафом Роннебергером, Филиппом Фишером (Philipp Fischer) и Томасом Броксом (Thomas Brox))

На этой схеме слева показано тело CNN (в данном случае это стандартная CNN, а не ResNet, и авторы вместо сверток с шагом 2 используют макс-пулинг 2×2 , поскольку рассматриваемая работа была написана до появления ResNet), а транспонированные сверточные («up-conv») слои — справа. Дополнительные пропускающие соединения показаны в виде серых стрелок, направленных слева направо (иногда их называют *кросс-соединениями*). Теперь становится очевидным, почему архитектура называется *U-Net* (*сеть U*).

При ее использовании на вход транспонированных сверток передается не только сетка низкого разрешения из предыдущего слоя, но и сетка высокого разрешения из вершины ResNet. Это позволяет U-Net задействовать всю необходимую информацию исходного изображения. Проблема использования U-Net в том, что конкретная архитектура зависит от размера изображения. В fastai есть уникальный класс `DynamicUnet`, который автоматически генерирует архитектуру нужного размера на основе предоставленных данных.

Разберем пример, в котором мы задействуем библиотеку fastai для написания собственной модели.

Сиамская сеть

Вернемся к входному конвейеру, который настроили в главе 11 для сиамской сети. Вы можете помнить, что он состоял из пары изображений, с меткой либо `True`, либо `False`, что определялось их возможным отношением к одному классу.

Создадим на основе только что разобранного материала модель для этой задачи и обучим ее. Как мы это сделаем? Мы задействуем предварительно обученную архитектуру и передадим через нее два наших изображения. После мы конкатенируем результаты и отправим их в настроенную вершину, которая вернет два прогноза. В виде модулей все это будет выглядеть так:

```
class SiameseModel(Module):
    def __init__(self, encoder, head):
        self.encoder, self.head = encoder, head

    def forward(self, x1, x2):
        ftrs = torch.cat([self.encoder(x1), self.encoder(x2)], dim=1)
        return self.head(ftrs)
```

Для создания энкодера нам нужно просто взять предварительно обученную модель и обрезать ее, и это мы поручим функции `create_body`. Потребуется лишь передать ей информацию об участке, который нужно обрезать. Как мы уже видели, согласно словарю метаданных для предварительно обученных моделей, значение среза для ResNet будет -2:

```
encoder = create_body(resnet34, cut=-2)
```

Затем мы создадим вершину. Судя по энкодеру, последний слой содержит 512 признаков, значит, этой вершине нужно будет получить $512*4$. Почему 4? Сначала нам нужно умножить на 2, потому что у нас два изображения. Затем в связи с приемом конкатенации пулингов нужно выполнить второе умножение на 2. Следовательно, вершину мы создаем так:

```
head = create_head(512*4, 2, ps=0.5)
```

Теперь с помощью энкодера и вершины можно построить модель:

```
model = SiameseModel(encoder, head)
```

Прежде чем использовать `Learner`, нужно определить еще два компонента. Сначала необходимо определить функцию потерь, которой в данном случае выступит перекрестная энтропия, но так как наши цели являются логическими значениями, нужно преобразовать их в целые числа, иначе PyTorch выбросит ошибку:

```
def loss_func(out, targ):
    return nn.CrossEntropyLoss()(out, targ.long())
```

При этом чтобы полноценно задействовать преимущества переноса обучения, необходимо определить настраиваемый *разделитель*. Разделитель — это функция, сообщающая fastai, как разделять модель на группы параметров. При переносе обучения эти группы используются фоном только для обучения вершины.

Здесь нам нужно две группы параметров: одна для энкодера и вторая для вершины. Исходя из этого, мы определяем следующий разделитель (`params` — это просто функция, возвращающая все параметры заданного модуля):

```
def siamese_splitter(model):
    return [params(model.encoder), params(model.head)]
```

Затем мы определяем `Learner`, передавая данные, модель, функцию потерь, разделитель и любую нужную метрику. Поскольку мы не задействуем вспомогательную функцию fastai для переноса обучения (например, `cnn_learner`), потребуется вызвать `learn.freeze` вручную. Это гарантирует обучение только последней группы параметров (в данном случае вершины):

```
learn = Learner(dls, model, loss_func=loss_func,
                 splitter=siamese_splitter, metrics=accuracy)
learn.freeze()
```

После этого мы непосредственно обучим модель обычным способом:

```
learn.fit_one_cycle(4, 3e-3)
```

epoch	train_loss	valid_loss	accuracy	time
0	0.367015	0.281242	0.885656	00:26
1	0.307688	0.214721	0.915426	00:26
2	0.275221	0.170615	0.936401	00:26
3	0.223771	0.159633	0.943843	00:26

Далее разморозим и дополнительно тонко настроим всю модель с помощью дискриминативных скоростей обучения (то есть пониженной скорости для тела и повышенной — для вершины):

```
learn.unfreeze()  
learn.fit_one_cycle(4, slice(1e-6,1e-4))
```

epoch	train_loss	valid_loss	accuracy	time
0	0.212744	0.159033	0.944520	00:35
1	0.201893	0.159615	0.942490	00:35
2	0.204606	0.152338	0.945196	00:36
3	0.213203	0.148346	0.947903	00:36

94,8 % — это очень хороший показатель, если вспомнить, что так же обученный классификатор (без аугментации данных) допускал 7 % ошибок.

Теперь, рассмотрев создание полноценных эталонных моделей компьютерного зрения, пора переходить к NLP.

Обработка естественного языка

Преобразование языковой модели AWD-LSTM в классификатор путем переноса обучения, как мы делали в главе 10, похоже на то, что мы делали с `cnn_learner` в первом разделе главы. В этом случае нам не требуется метасловарь, так как в теле нет такого разнообразия поддерживаемых архитектур. Нам нужно лишь выбрать сформированную RNN для энкодера в языковой модели, который представлен одним модулем PyTorch. Этот энкодер предоставит активацию для каждого слова входных данных, потому что языковой модели нужно выводить прогноз для каждого следующего слова.

Чтобы создать на основе этого классификатор, мы используем подход, описанный в работе ULMFiT (<https://oreil.ly/3hdSj>) как «BPTT для классификации текста (BPT3C)»:

Мы разделяем документ на пакеты фиксированной длины с размером b . В начале каждого пакета модель инициализируется с конечным состоянием из предыдущего пакета. Мы отслеживаем скрытые состояния для среднего и макс-пулинга. Градиенты вычисляются для пакетов, чьи скрытые состояния повлияли на итоговый прогноз. На практике мы используем последовательности обратного распространения переменной длины.

Говоря иначе, классификатор содержит цикл `for`, который перебирает каждый пакет последовательности. Состояние поддерживается во всех пакетах, и активации каждого пакета сохраняются. В конце мы используем тот же прием конкатенации среднего и макс-пулинга, который применяем в моделях компьютерного зрения, но на этот раз выполняем выборку не из ячеек сетки CNN, а из последовательностей RNN.

Для цикла `for` нужно собрать данные в пакеты, но каждый текст должен рассматриваться отдельно, так как у каждого из них собственная метка. Тем не менее очень вероятно, что не все тексты будут одинаковой длины, а это означает, что мы не сможем поместить их в один массив, как делали с языковой моделью.

Здесь нам на выручку приходит заполнение: при извлечении набора текстов мы определяем один из них с наибольшей длиной. Затем мы заполняем более короткие специальным токеном `xpad`. Чтобы избежать крайних случаев, в которых текст с 2000 токенов оказывается в одном пакете с текстом, имеющим десять токенов (получится много заполнения и много вычислительных затрат), мы изменяем критерий случайности, обеспечивая совмещение в пакете текстов сопоставимого размера. Эти тексты по-прежнему будут иметь некий случайный порядок для обучающей выборки (для контрольной можно просто упорядочить их по длине), но не полностью.

Все это библиотека `fastai` выполняет фоном автоматически при создании `DataLoaders`.

Табличные модели

Напоследок рассмотрим модели `fastai.tabular`. (Нам не потребуется отдельно разбирать коллаборативную фильтрацию, так как мы уже видели, что эти модели являются просто табличными моделями или используют подход скалярного произведения, который мы реализовывали ранее.)

Вот метод `forward` для `TabularModel`:

```
if self.n_emb != 0:
    x = [e(x_cat[:,i]) for i,e in enumerate(self.embeds)]
    x = torch.cat(x, 1)
    x = self.emb_drop(x)

if self.n_cont != 0:
    x_cont = self.bn_cont(x_cont)
    x = torch.cat([x, x_cont], 1) if self.n_emb != 0 else x_cont
return self.layers(x)
```

Здесь мы не будем показывать `__init__`, поскольку он не столь интересен, но при этом рассмотрим каждую строку кода `forward`. Первая просто проверяет, есть

ли вложения для обработки, эту часть можно пропустить, если используются непрерывные переменные:

```
if self.n_emb != 0:
```

`self.embeds` содержит матрицы вложений, следовательно, эта строка получает активации каждого

```
x = [e(x_cat[:, i]) for i, e in enumerate(self.embeds)]
```

и конкатенирует их в один тензор:

```
x = torch.cat(x, 1)
```

После этого применяется дропаут. Для изменения этого значения можно передать `emb_drop` в `__init__`:

```
x = self.emb_drop(x)
```

Теперь мы проверяем, есть ли непрерывные переменные для обработки:

```
if self.n_cont != 0:
```

Они передаются через слой пакетной нормализации

```
x_cont = self.bn_cont(x_cont)
```

и конкатенируются с активациями вложений, если таковые присутствуют:

```
x = torch.cat([x, x_cont], 1) if self.n_emb != 0 else x_cont
```

Наконец, все это передается через линейные слои (каждый из которых включает пакетную нормализацию, если `use_bn` установлен как `True`, и дропаут, если в `ps` установлено значение или список значений):

```
return self.layers(x)
```

Поздравляем! Вот вы и познакомились с каждой деталью архитектур, используемых в библиотеке fastai!

Резюме

Теперь подробности архитектур глубокого обучения не должны вас пугать. Можете смело заглянуть в код fastai и PyTorch, проанализировав, что именно в нем происходит. При этом будет важнее не просто понять, что происходит, а то, почему. Почитайте научные работы, на которые идут ссылки в коде, и постараитесь понять, как этот код сопоставляется с описанными в них алгоритмами.

Разобравшись со всеми элементами модели и передаваемыми ей данными, можно рассмотреть, что это значит для практического глубокого обучения. Если у вас есть неограниченный объем данных, памяти и времени, то совет будет прост: обучайте огромную модель на всех имеющихся данных в течение длительного времени. Однако сложность глубокого обучения в том, что обычно количество данных, а также ресурсы памяти и времени ограничены. В таких случаях решением будет обучение небольшой модели. Если же у вас нет возможности проводить обучение достаточно долго для достижения переобучения, то вы не задействуете преимущество емкости вашей модели.

Так что первым шагом будет достичь точки, где начинается переобучение, и после решить вопрос о его уменьшении. На рис. 15.3 показан рекомендованный порядок шагов начиная с этого момента.

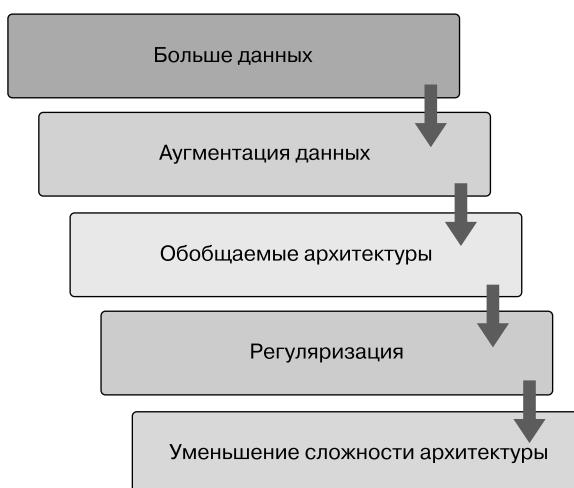


Рис. 15.3. Шаги для уменьшения переобучения

Многие практикующие специалисты, сталкиваясь с переобучением модели, начинают с другого края. Они задействуют уменьшенную модель или дополнительную регуляризацию. Уменьшение модели должно однозначно стоять в конце этой лестницы, если только обучение не требует слишком много памяти и времени, потому что, уменьшив модель, вы также снизите ее способность изучать тонкие связи данных.

Вместо этого первым шагом должен идти поиск способа *создания большего объема данных*. Это может подразумевать добавление дополнительных меток к уже имеющимся данным, нахождение дополнительных задач, которые можно поручить модели (или рассмотрение ее с другой стороны для определения разных видов меток, которые можно смоделировать), либо создание дополнительных

синтетических данных, используя больше техник аугментации или их иные варианты. Благодаря разработке Mixup и аналогичных подходов эффективная аугментация данных стала доступна практически для всех их видов.

Если вы собрали достаточное, на ваш взгляд, количество данных и начали использовать их максимально эффективно, используя все доступные метки и применяя все подходящие типы аугментации, но при этом переобучение сохраняется, тогда следует рассмотреть вариант использования более обобщающихся архитектур. Как вариант, обобщаемость можно повысить добавлением пакетной нормализации.

Если и после этого происходит переобучение, то стоит обратиться к регуляризации. Говоря обобщенно, эту задачу можно успешно реализовать, добавив дропаут в один или два последних слоя. Тем не менее, как показывает опыт разработки AWD-LSTM, добавление дропаута разного типа по всей модели может оказаться даже лучший эффект. Как правило, более крупная модель с большим объемом регуляризации является более гибкой и, следовательно, более точной, чем ее меньшие аналоги с меньшим объемом регуляризации.

Только после рассмотрения всех этих вариантов мы советуем прибегать к уменьшению архитектуры.

Вопросник

1. Что такое вершина нейронной сети?
2. Что такое тело нейронной сети?
3. Что значит «обрезать» нейронную сеть? Зачем это нужно при переносе обучения?
4. Что такое `model_meta`? Попробуйте вывести его и посмотреть, что внутри.
5. Прочтите исходный код `create_head`, убедившись, что поняли роль каждой строки.
6. Проанализируйте выход `create_head`, разобравшись, почему появился каждый слой и как исходный код `create_head` его создал.
7. Выясните, как изменить дропаут, размер слоя и количество слоев, создаваемых `create_cnn`, а затем попробуйте найти значения, которые приведут к повышению точности распознавателя пород домашних животных.
8. Что делает `AdaptiveConcatPool2d`?
9. Что такое интерполяция методом ближайшего соседа? Как ее можно использовать для увеличения количества сверточных активаций?
10. Что такое транспонированная свертка?

11. Создайте сверточный слой с `transpose=True` и примените его к изображению. Проверьте форму на выходе.
12. Нарисуйте архитектуру U-Net.
13. Что такое ВРТТ для классификации текста (ВРТЗС)?
14. Как мы обрабатываем последовательности разной длины в ВРТЗС?
15. Выполните в блокноте все строки `TabularModel.forward` по отдельности в разных ячейках и посмотрите, какой будет форма входа и выхода на каждом шаге.
16. Как в `TabularModel` определяется `self.layers`?
17. Какие пять шагов предпринимаются для уменьшения переобучения?
18. Почему мы не уменьшаем сложность архитектуры до применения других подходов по снижению переобучения?

Дополнительные задания

1. Напишите собственную вершину и обучите с ней распознаватель пород домашних животных. Посмотрите, удастся ли вам добиться лучшего результата, чем в предустановленном варианте `fastai`.
2. Попробуйте переключиться в вершине CNN между `AdaptiveConcatPool2d` и `AdaptiveAvgPool2d`, проанализировав, что это изменит.
3. Напишите собственный разделитель для создания отдельной группы параметров для каждого блока ResNet и отдельной группы для опоры. Попробуйте выполнить с его применением обучение и посмотрите, улучшит ли это распознаватель животных.
4. Прочтите онлайн-главу, посвященную генеративным зрительным моделям, и создайте собственный колоризатор, модель сверхвысокого разрешения или переноса стиля.
5. Создайте собственную вершину, используя интерполяцию методом ближайшего соседа, и задействуйте ее для выполнения сегментации на датасете CamVid.

ГЛАВА 16

Процесс обучения

Теперь вы умеете создавать эталонные архитектуры для компьютерного зрения, обработки изображений, табличного анализа и колаборативной фильтрации, при этом вы также умеете быстро их обучать. Но на этом мы не заканчиваем, так как нам еще предстоит изучить некоторые особенности процесса обучения.

В главе 4 мы объясняли основы стохастического градиентного спуска: передать мини-пакет в модель, сравнить его с целью с помощью функции потерь, затем вычислить градиенты этой функции в отношении каждого веса, после чего обновить веса через следующую формулу:

```
new_weight = weight - lr * weight.grad
```

Мы реализовывали это с нуля в цикле обучения и видели, что в PyTorch есть простой класс `nn.SGD`, который выполняет это вычисление для каждого параметра за нас. В данной главе мы построим более быстрые оптимизаторы, используя гибкую основу. Но это не единственное, что мы хотим изменить в процессе обучения. Для внесения любых доработок в цикл обучения нам нужен способ добавления определенного кода в основу SGD. В fastai для этого есть система обратных вызовов, о которой мы вскоре вам все расскажем.

Начнем со стандартного SGD, чтобы получить базовый вариант, а затем представим наиболее распространенные оптимизаторы.

Создание базовой модели

Сначала создадим базовую модель, используя простой SGD, и сравним ее с предустановленным в fastai оптимизатором. Начнем с извлечения данных из Imagenette, используя метод `get_data`, который применяли в главе 14:

```
dls = get_data(URLs.IMAGENETTE_160, 160, 128)
```

Создадим ResNet-34 без предварительного обучения и передадим все полученные аргументы:

```
def get_learner(**kwargs):
    return cnn_learner(dls, resnet34, pretrained=False,
                       metrics=accuracy, **kwargs).to_fp16()
```

Вот предустановленный оптимизатор fastai с обычной скоростью обучения Зе-3:

```
learn = get_learner() learn.fit_one_cycle(3, 0.003)
```

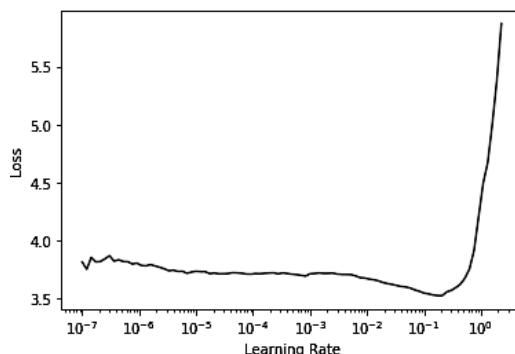
epoch	train_loss	valid_loss	accuracy	time
0	2.571932	2.685040	0.322548	00:11
1	1.904674	1.852589	0.437452	00:11
2	1.586909	1.374908	0.594904	00:11

Теперь применим стандартный SGD. Мы передадим функцию оптимизации `opt_func` в `cnn_learner`, чтобы fastai использовала любой оптимизатор:

```
learn = get_learner(opt_func=SGD)
```

Сначала рассмотрим `lr_find`:

```
learn.lr_find()
(0.017378008365631102, 3.019951861915615e-07)
```



Похоже, нам нужно повысить скорость обучения:

```
learn.fit_one_cycle(3, 0.03, mom=(0,0,0))
```

epoch	train_loss	valid_loss	accuracy	time
0	2.969412	2.214596	0.242038	00:09
1	2.442730	1.845950	0.362548	00:09
2	2.157159	1.741143	0.408917	00:09

Поскольку ускорение SGD-импульсом очень эффективно, fastai делает это в `fit_one_cycle` по умолчанию, так что пока мы отключим это с помощью `moms=(0,0,0)`, а импульс рассмотрим несколько позже.

Очевидно, что простой SGD выполняет обучение с недостаточной для нас скоростью. Так что освоим некоторые приемы, которые позволят этот процесс ускорить.

Универсальный оптимизатор

Для создания ускоренных версий SGD нам потребуется начать с удобной гибкой основы оптимизатора. fastai стала первой библиотекой с поддержкой такой основы, но в процессе разработки fastai мы поняли, что все улучшения оптимизатора, которые встречались нам в академической литературе, можно организовать с помощью *обратных вызовов оптимизатора*. Они представляют собой небольшие части кода, которые можно совмещать, перемешивать и со-поставлять в оптимизаторе, реализуя этап оптимизации. Вызываются же эти обратные вызовы легковесным fastai-классом `Optimizer`. Вот определения двух ключевых методов в `Optimizer`, которые мы использовали на протяжении книги:

```
def zero_grad(self):
    for p,*_ in self.all_params():
        p.grad.detach_()
        p.grad.zero_()

def step(self):
    for p,pg,state,hyper in self.all_params():
        for cb in self.cbs:
            state = _update(state, cb(p, **{**state, **hyper}))
        self.state[p] = state
```

Как вы видели при обучении модели на MNIST, `zero_grad` просто перебирает параметры модели и устанавливает градиенты на ноль. Он также вызывает метод `detach_`, который удаляет историю вычисления градиентов, поскольку после `zero_grad` она уже не понадобится.

Более интересен метод `step`, который перебирает обратные вызовы (`cbs`) и вызывает их для обновления параметров (функция `_update` просто вызывает `state_update`, если `cb` что-либо вернул). Как видите, сам по себе `Optimizer` не выполняет никаких шагов SGD. Теперь посмотрим, как добавить SGD в `Optimizer`.

Вот обратный вызов оптимизатора, выполняющий один шаг SGD, умножая `-lr` на градиенты и прибавляя итог к этому параметру (когда в PyTorch в `Tensor.add` передается два параметра, они перед сложением перемножаются).

```
def sgd_cb(p, lr, **kwargs): p.data.add_(-lr, p.grad.data)
```

Это мы передадим в `Optimizer`, используя параметр `cbs`. Для этого нужно применить `partial`, так как `Learner` вызовет эту функцию для создания оптимизатора позже:

```
opt_func = partial(Optimizer, cbs=[sgd_cb])
```

Посмотрим на результат обучения:

```
learn = get_learner(opt_func=opt_func)
learn.fit(3, 0.03)
```

epoch	train_loss	valid_loss	accuracy	time
0	2.730918	2.009971	0.332739	00:09
1	2.204893	1.747202	0.441529	00:09
2	1.875621	1.684515	0.445350	00:09

Работает! Вот как создается SGD с нуля в fastai. Теперь посмотрим, что такое импульс.

Импульс

Как описывалось в главе 4, SGD можно представить себе в виде пошагового спуска с вершины горы, при котором каждый ваш шаг ориентирован в направлении самого крутого уклона. Но что, если у нас мяч, который просто с этой горы катится? Он не будет в каждой заданной точке в точности следовать направлению градиента, так как будет подвержен *импульсу*. Мяч с большим импульсом (например, более тяжелый) будет проскаакивать через небольшие выпуклости и ямы и с большей вероятностью достигнет нижней точки горы. И наоборот, шарик для пинг-понга будет застревать в каждой мелкой впадинке.

Как же перенести эту идею в SGD? Можно задействовать для совершения шага скользящее среднее дополнительно к текущему градиенту:

```
weight.avg = beta * weight.avg + (1-beta) * weight.grad
new_weight = weight - lr * weight.avg
```

Здесь `beta` — это выбираемое нами число, которое определяет величину импульса. Если `beta` будет равно 0, то первое уравнение станет `weight.avg = weight.grad`, и мы в итоге получим простой SGD. Но если это число сделать близким к 1, то основное направление будет выбираться как среднее на основе предыдущих шагов. (Если вы занимались статистикой, то можете узнать в первом уравнении *экспоненциально взвешенное скользящее среднее*, которое зачастую используется для уменьшения шума данных и получения их внутренней тенденции.)

Обратите внимание, что мы пишем `weight.avg`, подчеркивая тот факт, что нам нужно сохранить скользящее среднее для каждого параметра модели (все они — свои собственные независимые скользящие средние).

На рис. 16.1 показан пример запущенных данных для одного параметра, где красной кривой отражен импульс, а синими точками — градиенты. Градиенты возрастают, затем понижаются, а импульс успешно делает свое дело, следуя общей тенденции и не подвергаясь излишнему вмешательству шума.

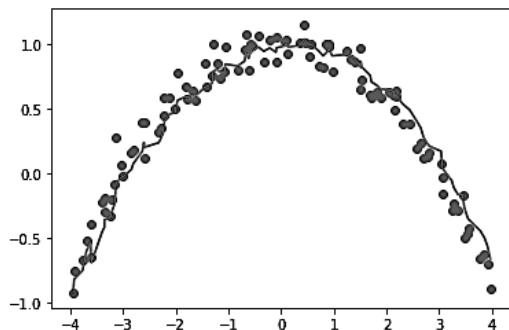


Рис. 16.1. Пример импульса

В особенности хорошо он работает, если в функции потерь есть узкие каньоны, по которым нужно пройти: чистый SGD отправил бы нас блуждать из стороны в сторону, в то время как SGD с импульсом усреднит эти блуждания, и спуск произойдет гладко вдоль склона. Параметр `beta` определяет силу импульса: при его низком значении мы держимся ближе к фактическим значениям градиентов, а при высоком будем больше отходить в направлении среднего этих градиентов, и для смещения этой тенденции любым изменением градиентов потребуется некоторое время.

При высоком `beta` мы можем пропустить факт изменения градиентами направления и перескочить малый локальный минимум. Это желательный побочный эффект: очевидно, что при показе модели новых входных данных они будут похожи на что-то из обучающей выборки, но не *в точности*. Эти данные будут соответствовать точке на функции потерь, которая близка к минимуму, на котором мы закончили обучение, но в точности этим минимумом являться не будет. Поэтому мы предпочтет закончить обучение в широком минимуме, где ближние точки имеют приблизительно одинаковые потери (или, если хотите, в точке, где функция максимально плоская). На рис. 16.2 показано, как график с рис. 16.1 меняется при изменении `beta`.

Здесь мы видим, что чересчур высокое значение `beta` приводит к общему игнорированию изменений в градиентах. Чаще всего в SGD с импульсом для `beta` устанавливается значение 0,9.

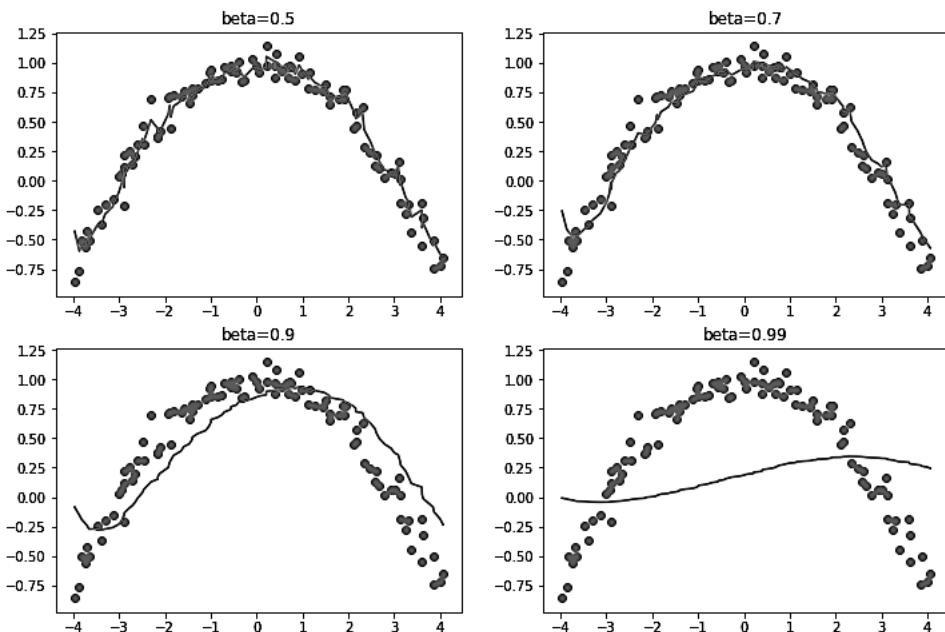


Рис. 16.2. Импульс при разных значениях beta

`fit_one_cycle` по умолчанию начинает с `beta`, равного 0,95, постепенно подстраивает его до 0,85, а к концу обучения возвращает к 0,95. Посмотрим, как пойдет наше обучение после добавления в SGD импульса. Чтобы добавить в оптимизатор импульс, сначала нужно начать отслеживать скользящий средний градиент, что можно сделать с помощью еще одного обратного вызова. Когда обратный вызов оптимизатора возвращает `dict`, он используется для обновления состояния этого оптимизатора и передается ему обратно на следующем шаге. Этот обратный вызов будет отслеживать средние градиентов в параметре `grad_avg`:

```
def average_grad(p, mom, grad_avg=None, **kwargs):
    if grad_avg is None: grad_avg = torch.zeros_like(p.grad.data)
    return {'grad_avg': grad_avg*mom + p.grad.data}
```

Для его использования нужно просто заменить `p.grad.data` на `grad_avg` в ступенчатой функции:

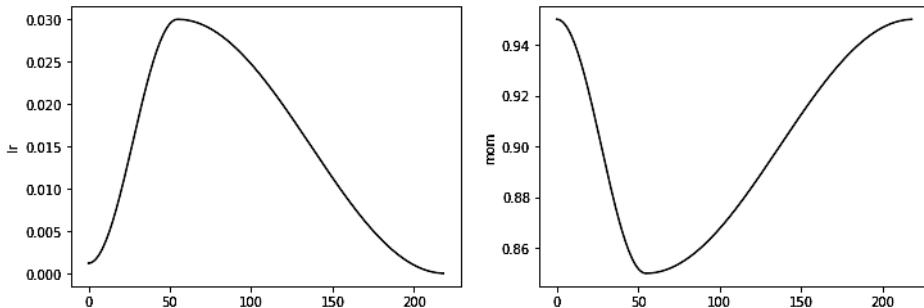
```
def momentum_step(p, lr, grad_avg, **kwargs): p.data.add_(-lr, grad_avg)
opt_func = partial(Optimizer, cbs=[average_grad,momentum_step], mom=0.9)
```

`Learner` автоматически распланирует `mom` и `lr`, так что `fit_one_cycle` будет работать даже в настраиваемом `Optimizer`:

```
learn = get_learner(opt_func=opt_func)
learn.fit_one_cycle(3, 0.03)
```

epoch	train_loss	valid_loss	accuracy	time
0	2.856000	2.493429	0.246115	00:10
1	2.504205	2.463813	0.348280	00:10
2	2.187387	1.755670	0.418853	00:10

```
learn.recorder.plot_sched()
```



Мы по-прежнему не добились отличных результатов, поэтому посмотрим, что еще можно сделать.

RMSProp

RMSProp — это еще один вариант SGD, представленный Джейфри Хинтоном (Geoffrey Hinton) в лекции его курса *Coursera Neural Networks for Machine Learning* (<https://oreil.ly/FVcIE>) («Нейронные сети для машинного обучения»). Основное отличие от SGD в том, что он использует адаптивную скорость обучения: вместо применения одинаковой скорости для всех параметров каждый параметр получает собственное ее значение, управляемое глобальной скоростью обучения. Таким образом, мы ускоряем обучение, задавая повышенную скорость весам, которые требуют существенного изменения, и сохраняя низкую для тех, которые уже достаточно хороши.

Как решить, каким параметрам нужно назначить высокую скорость, а каким — низкую? Можно взглянуть на градиенты. Если градиенты параметра какое-то время оставались близкими к нулю, этот параметр потребует повышенной скорости, потому что функция потеря плоская. И наоборот, если градиенты сильно разбросаны, то скорее всего, нужно быть осторожными и выбрать низкую скорость, чтобы избежать расхождения. Мы не можем просто усреднить градиенты, чтобы оценить, сильно ли они изменяются, так как среднее для очень больших и очень малых значений окажется близко к нулю. Вместо этого мы используем обычный прием, получая либо абсолютное значение, либо значения, возведенные в квадрат (а затем извлекая из них среднего квадратный корень).

И снова, чтобы определить общую тенденцию шума, мы будем использовать скользящее среднее, в частности скользящее среднее градиентов, возвещенных в квадрат. Затем мы обновим соответствующий вес с помощью текущего градиента (для направления), разделенного на квадратный корень этого скользящего среднего (таким образом, если он окажется мал, эффективная скорость обучения будет выше, если же высок, то ниже):

```
w.square_avg = alpha * w.square_avg + (1-alpha) * (w.grad ** 2)
new_w = w - lr * w.grad / math.sqrt(w.square_avg + eps)
```

`eps` (эпсилон) добавляется для численной стабильности (обычно устанавливается как `1e-8`), а предустановленное значение — для `alpha` обычно `0.99`. Можно добавить это к `Optimizer`, проделав в основном то же, что мы делали для `avg_grad`, но с дополнительным `**2`:

```
def average_sqr_grad(p, sqr_mom, sqr_avg=None, **kwargs):
    if sqr_avg is None: sqr_avg = torch.zeros_like(p.grad.data)
    return {'sqr_avg': sqr_avg*sqr_mom + p.grad.data**2}
```

Определить ступенчатую функцию и оптимизатор можно как и ранее:

```
def rms_prop_step(p, lr, sqr_avg, eps, grad_avg=None, **kwargs):
    denom = sqr_avg.sqrt().add_(eps)
    p.data.addcdiv_(-lr, p.grad, denom)

opt_func = partial(Optimizer, cbs=[average_sqr_grad,rms_prop_step],
                   sqr_mom=0.99, eps=1e-7)
```

Проверим все это в деле:

```
learn = get_learner(opt_func=opt_func)
learn.fit_one_cycle(3, 0.003)
```

epoch	train_loss	valid_loss	accuracy	time
0	2.766912	1.845900	0.402548	00:11
1	2.194586	1.510269	0.504459	00:11
2	1.869099	1.447939	0.544968	00:11

Уже лучше! Теперь нужно только объединить эти идеи, и у нас получится Adam, оптимизатор, установленный в fastai по умолчанию.

Adam

Adam смешивает принципы SGD с импульсом и RMSProp: он использует скользящее среднее градиентов в качестве направления и делит на квадратный

корень скользящего среднего градиентов, возведенных в квадрат, вычисляя для каждого параметра аддитивную скорость обучения.

Есть и еще одно отличие в том, как Adam вычисляет скользящие средние. Он получает *несмешенное* скользящее среднее, а именно

```
w.avg = beta * w.avg + (1-beta) * w.grad
unbias_avg = w.avg / (1 - (beta**(i+1)))
```

Если мы выполняем i -ю итерацию (начиная с 0, как делает Python). Этот делитель $1 - (\text{beta}^{**}(i+1))$ гарантирует, что несмешенное среднее больше похоже на градиенты в начале (поскольку $\text{beta} < 1$, знаменатель очень быстро приближается к 1).

Объединив все воедино, мы получим следующий шаг обновления:

```
w.avg = beta1 * w.avg + (1-beta1) * w.grad
unbias_avg = w.avg / (1 - (beta1**(i+1)))
w.sqr_avg = beta2 * w.sqr_avg + (1-beta2) * (w.grad ** 2)
new_w = w - lr * unbias_avg / sqrt(w.sqr_avg + eps)
```

Как и для RMSProp, eps обычно устанавливается как $1e-8$, а для ($\text{beta1}, \text{beta2}$) в пособиях предлагается значение (0.9, 0.999).

В fastai Adam задействуется в качестве оптимизатора по умолчанию, так как он позволяет ускорить обучение, но мы выяснили, что для используемого нами вида планирования лучше подходит `beta2=0.99`. `beta1` — это параметр импульса, который мы указываем с помощью аргумента `moms` в вызове к `fit_one_cycle`. Что касается eps , fastai по умолчанию устанавливает для него $1e-5$. При этом eps полезен не только для поддержания численной устойчивости. Высокий eps ограничивает максимальное значение подстраиваемой скорости обучения. Вот яркий пример: если eps равен 1, то подстраиваемая скорость никогда не превысит базовую.

Вместо того чтобы показывать весь этот код в книге, мы предлагаем вам заглянуть в блокнот оптимизатора по адресу https://oreil.ly/24_O [GitHub] (откройте каталог `_nbs` и найдите блокнот под названием *optimizer*). В нем вы увидите весь показанный нами код с Adam и другими оптимизаторами, а также множество примеров и тестов.

При переходе от SGD к Adam изменяется то, как мы применяем сокращение весов, что может иметь важные последствия.

Раздельное сокращение весов

Сокращение весов, которое мы рассматривали в главе 8, является эквивалентом (в случае чистого SGD) обновления параметров следующим путем:

```
new_weight = weight - lr*weight.grad - lr*wd*weight
```

Последняя часть формулы объясняет название данной техники: каждый вес сокращается на $lr * wd$.

Иначе сокращение весов называют *регуляризацией L2*, которая состоит из прибавления суммы всех весов в квадрате к потерям (умноженным на сокращение веса). Как мы видели в главе 8, это можно непосредственно выразить в градиентах:

```
weight.grad += wd*weight
```

Для SGD эти две формулы равнозначны. Тем не менее эта равнозначность сохраняется только для стандартного SGD, поскольку, как мы видели в случае с импульсом, RMSProp или в Adam, обновление содержит некоторые дополнительные формулы для градиента.

Большинство библиотек используют вторую формулу, но в работе Ильи Лошилова и Фрэнка Хаттера (Ilya Loshchilov, Frank Hutter) *Decoupled Weight Decay Regularization* (<https://oreil.ly/w37Ac>) («Регуляризация раздельного сокращения весов») говорится, что первый вариант является единственным верным подходом при использовании Adam или импульса, в связи с чем он установлен в fastai по умолчанию.

Теперь вам известно все, что кроется за строкой `learn.fit_one_cycle`.

Тем не менее оптимизаторы — это только одна часть процесса обучения. Когда вам нужно изменить его цикл с помощью fastai, вы не можете непосредственно изменить код внутри библиотеки. Вместо этого мы разработали систему обратных вызовов, позволяющую вам писать любые нужные доработки в независимых блоках, которые можно смешивать и сопоставлять.

Обратные вызовы

Иногда вам требуется несколько изменить принцип работы некоторых компонентов. По правде говоря, мы с вами уже видели примеры этого: Mixup, обучение fp16, сброс модели после каждой эпохи для обучения RNN и т. д. Как же мы реализуем подобные доработки в процессе обучения?

Мы видели основной цикл обучения, который при участии класса `Optimizer` для одной эпохи выглядит так:

```
for xb,yb in dl:
    loss = loss_func(model(xb), yb)
    loss.backward()
    opt.step()
    opt.zero_grad()
```

На рис. 16.3 это отражено графически.

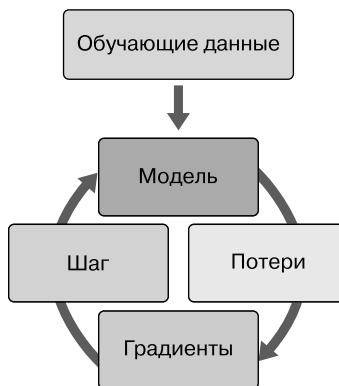


Рис. 16.3. Базовый цикл обучения

Обычно для настройки цикла обучения делается копия существующего цикла, а затем в нее вставляется код, необходимый для конкретных целей. Так выглядит почти весь код, который можно найти в Сети. Но здесь есть серьезные проблемы.

Маловероятно, что какой-то конкретный доработанный цикл обучения будет отвечать вашим конкретным нуждам. В нем могут быть произведены сотни изменений, что подразумевает миллионы возможных комбинаций. Вы не можете просто скопировать одну доработку из этого цикла обучения, другую — из того, ожидая, что все они дружно заработают вместе. Каждый из них будет опираться на другие предположения о среде функционирования, использовать другие условные имена и ожидать данные в другом формате.

Нам нужен способ дать пользователям возможность вставлять свой код в любую часть цикла обучения, но согласованным и грамотно определенным способом. Ученые уже придумали для этого неплохое решение: обратный вызов. *Обратный вызов* — это часть кода, которую вы пишете и вставляете в другую часть кода в предопределенном для этого месте. На самом деле обратные вызовы использовались в циклах глубокого обучения годами. Проблема в том, что в более ранних библиотеках допускалась вставка кода в очень ограниченный спектр областей, где это могло потребоваться, и при этом обратные вызовы не могли делать всё, что от них требовалось.

Чтобы быть столь же гибким, как ручное копирование кода с последующей вставкой в него дополнительных частей, обратный вызов должен иметь возможность считывать каждый доступный элемент информации цикла обучения, при необходимости изменять его и полностью контролировать, когда нужно остановить пакет, эпоху или даже весь цикл обучения. fastai стала первой библиотекой, в которой вся эта функциональность была представлена. Она изменяет цикл обучения, как показано на рис. 16.4.

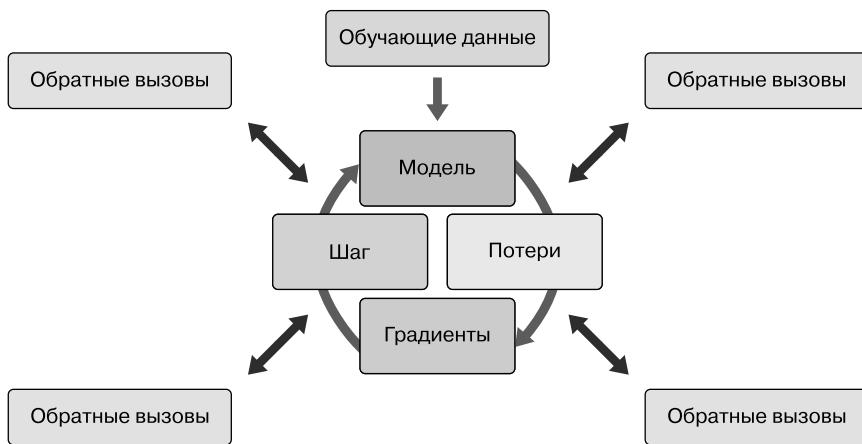


Рис. 16.4. Цикл обучения с обратными вызовами

Эффективность данного подхода была подтверждена в течение последних двух лет: используя систему обратных вызовов fastai, мы смогли реализовать каждую новую работу и удовлетворить каждое пользовательское обращение по изменению цикла обучения. При этом изменять сам цикл не требовалось. На рис. 16.5 показано лишь несколько вариантов из общего числа добавленных обратных вызовов.

Это важно, поскольку таким образом мы можем реализовать любые возникающие у нас идеи. Нам никогда не придется для этого рыться в исходном коде PyTorch или fastai и создавать разовую систему, чтобы опробовать свои замыслы. А когда мы реализуем собственные обратные вызовы для разработки своих идей, мы знаем, что они будут работать согласованно со всей остальной функциональностью fastai, то есть мы получим индикаторы прогресса, обучение со смешанной точностью, нормализацию гиперпараметров и т. д.

Еще одно преимущество в том, что это облегчает постепенное удаление или добавление функционала и выполнение экспериментов по аблации. Нужен лишь список обратных вызовов, который передается в функцию подстройки.

Вот пример исходного кода fastai, который выполняется для каждого пакета в цикле обучения:

```

try:
    self._split(b);
    self.pred = self.model(*self.xb);
    self.loss = self.loss_func(self.pred, *self.yb);
    if not self.training: return
    self.loss.backward();
    self.opt.step();
    self.opt.zero_grad()
except CancelBatchException:
    finally:
        self('begin_batch')
        self('after_pred')
        self('after_loss')
        self('after_backward')
        self('after_step')
        self('after_cancel_batch')
        self('after_batch')

```

SGDR	Запись метрик	Мульти-ГПУ	Смешанная точность
GAN	Отсечение градиентов	Icycle	SWA
Mixup	Сохранение лучшей модели	Tensorboard	Ранняя остановка
Накопление градиентов	Штрафы за потери	AdamW	Ваши!

Рис. 16.5. Некоторые из обратных вызовов fastai

Вызовы с формой `self('...')` находятся там, где выполняются обратные вызовы. Как видите, это происходит после каждого шага. Обратный вызов получает все состояние обучения, а также может его изменить. Например, входные данные и целевые метки находятся в `self.xb` и `self.yb` соответственно. Обратный вызов может скорректировать их, изменив таким образом данные, которые видит цикл обучения. Он также может скорректировать `self.loss` или даже градиенты.

Посмотрим, как это работает на практике, написав такой обратный вызов.

Создание обратного вызова

Когда вам нужно написать собственный обратный вызов, можете обратиться к следующему списку доступных событий.

- `begin_fit` — вызывается до начала выполнения обучения; отлично подходит для начальной настройки.
- `begin_epoch` — вызывается в начале каждой эпохи; полезен при необходимости сброса любого поведения каждой эпохи.
- `begin_train` — вызывается в начале обучающего этапа эпохи.
- `begin_batch` — вызывается в начале каждого пакета, сразу после его извлечения. Может использоваться для выполнения любой настройки пакета (например, планирования гиперпараметров) или для изменения входа/цели до их передачи в модель (например, применением Mixup).
- `after_pred` — вызывается после вычисления моделью выходного результата пакета. Может использоваться для изменения этого результата до его передачи функции потерь.

- `after_loss` — вызывается после вычисления потерь, но до обратного распространения. Может применяться для добавления штрафа к потерям (например, AR или TAR в обучении RNN).
- `after_backward` — вызывается после обратного распространения, но до обновления параметров. Можно использовать для внесения изменений в градиенты перед их обновлением (например, через усечение градиентов).
- `after_step` — вызывается после шага градиентов до их обнуления.
- `after_batch` — вызывается в конце пакета для выполнения необходимой очистки перед обработкой следующего.
- `after_train` — вызывается в конце фазы обучения.
- `begin_validate` — вызывается в начале фазы контроля; позволяет делать настройки непосредственно для контрольного этапа.
- `after_validate` — вызывается в конце контрольной фазы.
- `after_epoch` — вызывается в конце каждой эпохи для выполнения необходимой очистки перед началом следующей.
- `after_fit` — вызывается в конце обучения для итоговой очистки.

Элементы этого списка доступны в качестве атрибутов специальной переменной `event`, так что вы можете просто набрать в блокноте `event` и нажать `Tab`, вызвав таким образом список всех вариантов.

Взглянем на пример. Помните ли вы, как в главе 12 нам потребовалось обеспечить, чтобы особый метод `reset` вызывался в начале обучения и контроля в каждой эпохе? Мы использовали обратный вызов `ModelResetter`, предоставляемый `fastai`, который делал это автоматически. Но как конкретно он работает? Вот весь исходный код для этого класса:

```
class ModelResetter(Callback):
    def begin_train(self):    self.model.reset()
    def begin_validate(self): self.model.reset()
```

Да, все именно так! Он просто делает то, что мы только что описали: вызывает метод `reset` после завершения обучения или контроля в каждой эпохе.

Обратные вызовы зачастую «кратки и понятны», как только что приведенный. Но посмотрим еще на один. Вот исходный код `fastai` для обратного вызова, добавляющего регуляризацию RNN (AR и TAR):

```
class RNNRegularizer(Callback):
    def __init__(self, alpha=0., beta=0.): self.alpha,self.beta = alpha,beta
    def after_pred(self):
        self.raw_out,self.out = self.pred[1],self.pred[2]
        self.learn.pred = self.pred[0]
```

```
def after_loss(self):
    if not self.training: return
    if self.alpha != 0.:
        self.learn.loss += self.alpha * self.out[-1].float().pow(2).mean()
    if self.beta != 0.:
        h = self.raw_out[-1]
        if len(h)>1:
            self.learn.loss += self.beta * (h[:,1:] - h[:,:-1])
                .float().pow(2).mean()
```



РАЗБЕРИТЕСЬ В КОДЕ

Вернитесь назад и перечитайте подраздел «Регуляризация активаций и регуляризация временных активаций» на с. 447, а затем еще раз взгляните на приведенный здесь код. Как следует разберитесь, что именно он делает и почему.

Обратите внимание, как в обоих примерах мы можем обращаться к атрибутам цикла обучения, напрямую проверяя `self.model` или `self.pred`. Это возможно, потому что `Callback` всегда будет стараться получить атрибут, которого нет внутри `Learner`, с ним связанного. Это сокращенные версии `self.learn.model` или `self.learn.pred`. Заметьте, что они работают для считывания атрибутов, но не для их записи, поэтому когда `RNNRegularizer` изменяет потери или прогнозы, вы видите `self.learn.loss =` или `self.learn.pred =`.

При написании обратного вызова доступны следующие атрибуты `Learner`.

- `model` — модель, используемая для обучения/контроля.
- `data` — определяющий `DataLoaders`.
- `loss_func` — используемая функция потерь.
- `opt` — оптимизатор, используемый для обновления параметров модели.
- `opt_func` — функция, применяемая для создания оптимизатора.
- `cbs d1` — список, содержащий все `Callback`.
- `d1` — текущий `DataLoader`, используемый для итерации.
- `x/xb` — последние входные данные, полученные из `self.d1` (потенциально измененные обратными вызовами). `xb` — это всегда кортеж (потенциально с одним элементом), а `x` извлечен из кортежа. Присваивание можно делать только к `xb`.
- `y/yb` — последняя цель, полученная из `self.d1` (потенциально измененная обратными вызовами). `yb` — это всегда кортеж (потенциально с одним элементом), а `y` извлечен из кортежа. Присваивать можно только к `yb`.
- `pred` — последние прогнозы из `self.model` (потенциально измененные обратными вызовами).

- `loss` — последние вычисленные потери (потенциально измененные обратными вызовами).
- `n_epoch` — число эпох в текущем обучении.
- `n_iter` — число итераций в текущем `self.d1`.
- `epoch` — индекс текущей эпохи (от 0 до `n_epoch-1`).
- `iter` — индекс текущей итерации в `self.d1` (от 0 до `n_iter-1`).

Следующие атрибуты добавляются `TrainEvalCallback` и должны быть доступны, если только вы не удалили этот обратный вызов.

- `train_iter` — число итераций, выполненных с начала обучения.
- `pct_train` — процент завершенных итераций обучения (от 0 до 1).
- `training` — флаг, обозначающий нахождение в режиме обучения.

Следующий атрибут добавляется `Recorder` и должен быть доступен, если только вы не удалили этот обратный вызов:

- `smooth_loss` — экспоненциально усредненная версия потерь обучения.

Обратные вызовы также могут прерывать любой этап цикла обучения, используя систему исключений.

Упорядочивание обратных вызовов и исключения

Иногда может потребоваться дать `fastai` команду пропустить пакет, эпоху или даже полностью остановить обучение, для чего опять же задействуются обратные вызовы. Рассмотрите в качестве примера `TerminateOnNaNCallback`. Это удобный обратный вызов, который будет автоматически останавливать обучение каждый раз, когда потери будут становиться бесконечными или `Nan` (*не число*). Вот его исходный код:

```
class TerminateOnNaNCallback(Callback):
    run_before=Recorder
    def after_batch(self):
        if torch.isinf(self.loss) or torch.isnan(self.loss):
            raise CancelFitException
```

Строка `raise CancelFitException` указывает циклу обучения на необходимость прервать обучение в этой точке. Цикл перехватывает это исключение и прекращает выполнение обучения или контроля. Для потока управления программы доступны следующие варианты исключений.

- `CancelBatchException` — пропустить оставшуюся часть пакета и перейти к `after_batch`.

- `CancelEpochException` — пропустить оставшуюся часть обучения эпохи и перейти к `after_train`.
- `CancelTrainException` — пропустить оставшуюся часть этапа контроля и перейти к `after_validate`.
- `CancelValidException` — пропустить остаток эпохи и перейти к `after_epoch`.
- `CancelFitException` — прервать обучение и перейти к `after_fit`.

Вы можете определять момент срабатывания исключений и добавлять код, выполняющийся сразу за ними, с помощью следующих событий:

- `after_cancel_batch` — срабатывает сразу после `CancelBatchException` до перехода к `after_batch`.
- `after_cancel_train` — срабатывает сразу после `CancelTrainException` до перехода к `after_epoch`.
- `after_cancel_valid` — срабатывает сразу после `CancelValidException` до перехода к `after_epoch`.
- `after_cancel_epoch` — срабатывает сразу после `CancelEpochException` до перехода к `after_epoch`.
- `after_cancel_fit` — срабатывает сразу после `CancelFitException` до перехода к `after_fit`.

Иногда требуется выполнение обратных вызовов в определенном порядке. Например, в случае с `TerminateOnNaNCallback` важно, чтобы `Recorder` выполнял `after_batch` после этого обратного вызова во избежание регистрации потерь со значением `NaN`. Можно указать `run_before` (этот обратный вызов должен выполняться до...) или `run_after` (этот обратный вызов должен выполняться после...) в обратном вызове, обеспечивая тем самым требуемый порядок.

Резюме

В этой главе мы подробно рассмотрели цикл обучения, разобрав варианты SGD и то, почему они могут быть более эффективными. На момент написания книги активно разрабатываются новые оптимизаторы, так что к моменту прочтения вами этой главы они уже могут быть добавлены в качестве приложения на сайте (<https://book.fast.ai>). Обязательно изучите возможности быстрой реализации новых оптимизаторов на основе предложенной нами общей схемы.

Помимо этого, мы познакомились с мощной системой обратных вызовов, которая дает возможность инспектировать и изменять любой параметр между каждым шагом обучения, настраивая, таким образом, каждую часть его цикла.

Вопросник

1. Напишите уравнение для шага SGD в математическом виде или в виде кода, как вам больше нравится.
2. Что мы передаем в `cnn_learner` для использования непредустановленного оптимизатора?
3. Что такое обратные вызовы оптимизатора?
4. За что в оптимизаторе отвечает `zero_grad`?
5. За что в оптимизаторе отвечает `step`? Как он реализуется в общем оптимизаторе?
6. Перепишите `sgd_cb` для использования оператора `+=` вместо `add_`.
7. Что такое импульс? Напишите его уравнение.
8. Приведите физическую аналогию импульса. Как она применима в настройках обучения модели?
9. Как увеличение значения импульса влияет на градиенты?
10. Какие для импульса предустановлены значения в обучении 1cycle?
11. Что такое RMSProp? Напишите уравнение.
12. На что указывают возведенные в квадрат значения градиентов?
13. Чем Adam отличается от импульса и RMSProp?
14. Напишите уравнение для Adam.
15. Вычислите значения `unbias_avg` и `w.avg` для нескольких пакетов фиктивных значений.
16. Какое влияние оказывает высокий `eps` в Adam?
17. Прочтите блокнот оптимизатора в репозитории fastai и выполните его.
18. В каких ситуациях методы адаптивных скоростей обучения наподобие Adam изменяют поведение сокращения весов?
19. Из каких четырех этапов состоит цикл обучения?
20. Почему для добавления доработок лучше использовать обратные вызовы, чем писать новый цикл обучения?
21. Какие аспекты системы обратных вызовов fastai делают ее такой же гибкой, как копирование кода с последующей вставкой в него новых сегментов?
22. Как получить список доступных событий при написании обратного вызова?
23. Напишите обратный вызов `ModelResetter` (не подглядывая).
24. Как можно обратиться к нужным атрибутам цикла обучения внутри обратного вызова? Когда можно или нельзя использовать соответствующие им сокращения?

25. Как может обратный вызов повлиять на поток управления цикла обучения?
26. Напишите обратный вызов `TerminateOnNaN` (по возможности не подглядывая).
27. Как обеспечить выполнение одного обратного вызова до или после другого?

Дополнительные задания

1. Изучите работу «Улучшенный Adam», реализуйте его на основе общей схемы оптимизатора и опробуйте. Поиските другие эффективные новейшие оптимизаторы и выберите один для реализации.
2. Изучите обратный вызов смешанной точности в документации (<https://docs.fast.ai/>). Попытайтесь понять, что в нем делает каждое событие и строка кода.
3. Реализуйте с нуля собственный вариант искателя скорости обучения и сравните его с версией fastai.
4. Загляните в исходный код обратных вызовов, предлагаемых fastai. Попробуйте найти среди них подходящий для планируемой вами реализации, чтобы почерпнуть полезную информацию.

Основы глубокого обучения: итог

Поздравляем! Вы добрались до конца части III «Основы глубокого обучения». Теперь вы понимаете, как реализуются все возможности fastai и большинство важнейших архитектур, а также знаете рекомендуемые способы их обучения и имеете всю необходимую информацию для их создания с чистого листа. Несмотря на то что вам вряд ли потребуется разрабатывать, например, собственный цикл обучения или слой пакетной нормализации, понимание сути этих процессов очень пригодится вам при отладке, профайлинге и развертывании готовых решений.

Поскольку вы разобрались с основами возможностей fastai, обязательно уделите время подробному изучению исходных блокнотов и экспериментированию с ними. Так вы максимально эффективно усвоите все принципы внутренних реализаций fastai.

В следующей части мы заглянем еще глубже: подробнее изучим выполнение прямого и обратного прохода в нейронных сетях и посмотрим, какие в нашем распоряжении есть инструменты для повышения производительности. Затем мы перейдем к проекту, который объединит весь пройденный в книге материал, и создадим инструмент для интерпретирования сверточных нейронных сетей. Последним же, но не менее важным, будет создание класса `Learner` с нуля.

Часть IV

Глубокое обучение с чистого листа

ГЛАВА 17

Продвинутые основы нейронной сети

В этой главе начинается путешествие, которое погрузит нас в глубины моделей, использованных до этого момента. Мы разберем многое из того, что уже видели, но на этот раз уделим больше внимания именно деталям реализации, а не практическим аспектам.

Создавать мы все будем с чистого листа, используя только базовую индексацию тензора. Напишем нейронную сеть с самых ее основ, а затем реализуем обратное распространение вручную, разобрав, что конкретно происходит в PyTorch при вызове `loss.backward`. Мы также рассмотрим возможность расширения PyTorch собственными функциями *autograd*, которые позволяют указывать свои варианты прямого и обратного вычисления.

Создание слоя нейронной сети с нуля

Мы начнем с того, что вспомним принцип использования матричного умножения в базовой нейронной сети. Поскольку мы все будем создавать с нуля, то и использовать изначально будем только чистый Python (за исключением индексации тензоров PyTorch). Затем мы заменим чистый Python на функциональность PyTorch, когда увидим, как она создается.

Моделирование нейрона

Нейрон получает заданное число входов и имеет внутренний вес для каждого из них. Он суммирует взвешенные входы, создавая выход, и добавляет внутреннее смещение. Если мы назовем входы ($x_1 \dots x_n$), веса ($w_1 \dots w_n$), а смещение b , то на математическом языке это будет выглядеть так:

$$out = \sum_{i=1}^n x_i w_i + b.$$

В коде эквивалент будет следующим:

```
output = sum([x*w for x,w in zip(inputs,weights)]) + bias
```

Далее, прежде чем отправляться в следующий нейрон, этот выход передается в нелинейную функцию, а именно *функцию активации*. В глубоком обучении самой распространенной из них является блок линейной ректификации (ReLU), который, как мы видели, представляет просто заумную формулировку следующего кода:

```
def relu(x): return x if x >= 0 else 0
```

В итоге модель глубокого обучения создается путем составления множества таких нейронов в последовательные слои. Мы создаем первый слой с определенным числом нейронов (*скрытым размером*) и соединяем входы с каждым из этих нейронов. Такой слой часто называют *полносвязным, плотным или линейным слоем*.

Это требует вычисления скалярного произведения для каждого *input* и каждого нейрона с заданным *weight*:

```
sum([x*w for x,w in zip(input,weight)])
```

Если вы знакомы с линейной алгеброй, то можете вспомнить, что много таких скалярных произведений возникает при выполнении *матричного умножения*. Говоря точнее, если наши входы находятся в матрице *x* размером *batch_size* на *n_inputs* и если мы сгруппировали веса всех нейронов в матрицу *w* размером *n_neurons* на *n_inputs* (каждый нейрон должен иметь столько весов, сколько у него входов), а все смещения — в вектор размером *n_neurons*, то выход данного полносвязного слоя будет таким:

```
y = x @ w.t() + b
```

Здесь *@* представляет матричное произведение, а *w.t()* — это транспонированная матрица *w*. Размер выхода *y* будет равен *batch_size* на *n_neurons*, а в позиции *(i,j)* у нас получится следующее:

$$y_{i,j} = \sum_{k=1}^n x_{i,k} w_{k,j} + b_j.$$

Эквивалент в коде:

```
y[i,j] = sum([a * b for a,b in zip(x[i,:],w[j,:])]) + b[j]
```

Транспонирование необходимо, так как в математическом определении матричного произведения *m @ n* коэффициент *(i,j)* следующий:

```
sum([a * b for a,b in zip(m[i,:],n[:,j])])
```

Поэтому основной нужной нам операцией является матричное умножение, так как именно оно кроется в сердце нейронной сети.

Матричное умножение

Напишем функцию, вычисляющую матричное произведение двух тензоров, а уже потом перейдем к использованию ее версии из PyTorch. Мы задействуем только индексирование тензоров:

```
import torch
from torch import tensor
```

Нам потребуется три вложенных цикла `for`: один для индексов строк, второй для индексов столбцов и третий — для внутренней суммы. `ac` и `ar` означают количество столбцов `a` и количество строк `a` соответственно (то же применимо и к `b`). При этом мы проверяем, чтобы в `a` было столько же столбцов, сколько в `b` строк, подтверждая возможность вычисления матричного произведения:

```
def matmul(a,b):
    ar,ac = a.shape # n_rows * n_cols
    br,bc = b.shape
    assert ac==br
    c = torch.zeros(ar, bc)
    for i in range(ar):
        for j in range(bc):
            for k in range(ac): c[i,j] += a[i,k] * b[k,j]
    return c
```

Чтобы это протестировать, мы представим (используя случайные матрицы), что работаем с небольшим пакетом из пяти сглаженных до векторов 28×28 изображений MNIST, которые линейная модель преобразует в десять активаций:

```
m1 = torch.randn(5,28*28)
m2 = torch.randn(784,10)
```

Давайте замерим время выполнения нашей функции с помощью «магической» команды Jupyter `%time`:

```
%time t1=matmul(m1, m2)
CPU times: user 1.15 s, sys: 4.09 ms, total: 1.15 s
Wall time: 1.15 s
```

И сравним результат со встроенным в PyTorch @:

```
%timeit -n 20 t2=m1@m2
14 µs ± 8.95 µs per loop (mean ± std. dev. of 7 runs, 20 loops each)
```

Очевидно, что в Python три вложенных цикла не являются хорошим решением. Это медленный язык, и такой подход эффективным не будет. В данном случае PyTorch оказывается примерно в 100 000 быстрее Python, и это до начала использования GPU.

Откуда берется такая разница? Матричное умножение PyTorch писалось не на Python, а на C++. Как правило, при каждом выполнении вычислений над тензорами нам потребуется *векторизовать* их, чтобы задействовать преимущества скорости PyTorch. Для этого обычно используются две техники: поэлементная арифметика и уширение.

Поэлементная арифметика

Все базовые операторы ($+, -, *, /, >, <, ==$) можно применять поэлементно. Это означает, что если написать `a + b` для двух тензоров `a` и `b` одной формы, то мы получим тензор, составленный из сумм элементов `a` и `b`:

```
a = tensor([10., 6, -4])
b = tensor([2., 8, 7])
a + b
```

```
tensor([12., 14., 3.])
```

Логические операторы будут возвращать массив логических значений:

```
a < b
```

```
tensor([False, True, True])
```

Если мы хотим узнать, меньше ли каждый элемент в `a` соответствующего ему элемента в `b` или равны ли между собой два тензора, то нужно совместить эти поэлементные операции с помощью `torch.all`:

```
(a < b).all(), (a==b).all()
```

```
(tensor(False), tensor(False))
```

Операции приведения, такие как `all`, `sum` и `mean`, возвращают тензоры только с одним элементом, называемые *тензорами ранга 0*. Если вам нужно преобразовать такой тензор в логическое значение Python или число, то нужно вызвать `.item`:

```
(a + b).mean().item()
```

```
9.666666984558105
```

Поэлементные операции работают для тензоров любого ранга, только если их форма совпадает:

```
m = tensor([[1., 2, 3], [4,5,6], [7,8,9]])
m*m

tensor([[ 1.,  4.,  9.],
       [16., 25., 36.],
       [49., 64., 81.]])
```

Тем не менее нельзя выполнять такие операции для тензоров разной формы (если только они не подвергаются уширению, о чем мы будем говорить в следующем разделе):

```
n = tensor([[1., 2, 3], [4,5,6]])
m*n
```

```
RuntimeError: The size of tensor a (3) must match the size of tensor b (2) at
dimension 0
```

При использовании поэлементной арифметики можно удалить один из наших трех вложенных циклов: умножить тензоры, соответствующие *i*-й строке *a* и *j*-му столбцу *b* перед сложением всех элементов, что ускорит процесс в целом, так как внутренний цикл теперь будет выполняться PyTorch на скорости С.

Чтобы обратиться к одному столбцу или строке, нужно просто написать *a[i, :]* или *b[:, j]*. Здесь *:* означает «взять все в этом измерении». Можно ограничиться и взять только часть этого измерения, передав диапазон, например *1:5*, а не просто *:*. В этом случае мы получим элементы из столбцов от 1 до 4 (второе число (5) не включается).

В качестве упрощения всегда допустимо убрать двоеточие, тогда *a[i, :]* сократится до *a[i]*. С учетом всего этого можно написать новый вариант матричного умножения:

```
def matmul(a,b):
    ar,ac = a.shape
    br,bc = b.shape
    assert ac==br
    c = torch.zeros(ar, bc)
    for i in range(ar):
        for j in range(bc): c[i,j] = (a[i] * b[:,j]).sum()
    return c

%timeit -n 20 t3 = matmul(m1,m2)
```

```
1.7 ms ± 88.1 µs per loop (mean ± std. dev. of 7 runs, 20 loops each)
```

Мы уже ускорились примерно в 700 раз, просто удалив внутренний цикл *for*. И это только начало, так как с помощью уширения мы сможем удалить еще один цикл и получить даже более существенное ускорение.

Уширение (broadcasting)

Как мы уже говорили в главе 4, *уширение* — это термин, впервые использованный в библиотеке NumPy, который описывает взаимодействие тензоров разных рангов в процессе арифметических операций. Например, очевидно, нельзя сложить матрицу 3×3 с матрицей 4×5 . Но что, если мы хотим сложить с матрицей один скаляр (который можно представить как тензор 1×1)? Или вектор с размером 3 и матрицу 3×4 ? В обоих случаях можно найти способ реализации этой операции.

Уширение предоставляет ряд правил, по которым мы определяем случаи совместимости форм для выполнения поэлементных операций, а также расширяем тензоры меньшей формы для соответствия тензорам большей формы. Важно хорошо освоить эти правила, если вы хотите иметь возможность писать быстрый код. В текущем разделе мы лучше разберем суть процесса уширения, чтобы эти правила как следует понять.

Уширение скаляром

Это самый простой вид уширения. Когда у нас есть тензор a и скаляр, мы просто представляем тензор той же формы a , дополненной этим скаляром, и выполняем операцию:

```
a = tensor([10., 6, -4])
a > 0

tensor([ True,  True, False])
```

Как нам удается произвести это сравнение? 0 *уширяется*, получая столько же измерений, сколько есть у a . Обратите внимание: операция выполняется без создания тензора нулей в памяти, что было бы неэффективно.

Это пригодится, когда вы хотите нормализовать датасет, отняв среднее (скаляр) от всего датасета (матрицы) и произведя деление на стандартное отклонение (еще один скаляр):

```
m = tensor([[1., 2, 3], [4,5,6], [7,8,9]])
(m - 5) / 2.73
tensor([[-1.4652, -1.0989, -0.7326],
       [-0.3663,  0.0000,  0.3663],
       [ 0.7326,  1.0989,  1.4652]])
```

А что, если у вас разные средние для каждой строки матрицы? В этом случае потребуется уширить вектор в матрицу.

Уширение вектора в матрицу

Это можно сделать так:

```
c = tensor([10.,20,30])
m = tensor([[1., 2, 3], [4,5,6], [7,8,9]])
m.shape,c.shape

(torch.Size([3, 3]), torch.Size([3]))

m + c

tensor([[11., 22., 33.],
       [14., 25., 36.],
       [17., 28., 39.]])
```

Здесь элементы `c` расширяются, формируя три сопоставимые строки, что делает операцию возможной. И снова PyTorch не создает копии `c` в памяти. Процесс выполняет метод `expand_as` за кадром:

```
c.expand_as(m)
tensor([[10., 20., 30.],
       [10., 20., 30.],
       [10., 20., 30.]])
```

Если мы посмотрим на соответствующий тензор, то можем обратиться к его свойству `storage` (показывающему фактическое содержимое памяти, используемое для тензора), чтобы убедиться в отсутствии хранения бесполезных данных:

```
t = c.expand_as(m)
t.storage()

10.0
20.0
30.0
[torch.FloatTensor of size 3]
```

Несмотря на то что официально тензор содержит девять элементов, в памяти хранится только три скаляра. Это возможно благодаря умному приему, определяющему для этого измерения *шаг 0* (это означает, что PyTorch при поиске следующей строки путем добавления шага не переместится):

```
t.stride(), t.shape

((0, 1), torch.Size([3, 3]))
```

Поскольку `m` имеет размер 3x3, уширение можно выполнить двумя способами. Тот факт, что выше это было сделано для последнего измерения, является одним из условий трансформации и никак не связан с тем, как мы упорядочили тензоры. Если вместо этого написать следующее, то получится тот же результат:

```
c + m
tensor([[11., 22., 33.],
       [14., 25., 36.],
       [17., 28., 39.]])
```

Фактически уширение вектора с размером n возможно только в матрицу размером m на n :

```
c = tensor([10.,20,30])
m = tensor([[1., 2, 3], [4,5,6]])
c+m
```

```
tensor([[11., 22., 33.],
       [14., 25., 36.]])
```

Это не сработает:

```
c = tensor([10.,20])
m = tensor([[1., 2, 3], [4,5,6]])
c+m
```

```
RuntimeError: Размер тензора a (2) должен соответствовать размеру тензора b (3)
в измерении 1
```

Если нам нужно выполнить уширение в другом измерении, то потребуется изменить форму вектора, сделав его матрицей 3×1 . Это выполняется методом `unsqueeze`:

```
c = tensor([10.,20,30])
m = tensor([[1., 2, 3], [4,5,6], [7,8,9]])
c = c.unsqueeze(1).m.shape,c.shape

(torch.Size([3, 3]), torch.Size([3, 1]))
```

На этот раз `c` расширяется по столбцам:

```
c+m
tensor([[11., 12., 13.],
       [24., 25., 26.],
       [37., 38., 39.]])
```

Как и ранее, в памяти сохраняются только три скаляра:

```
t = c.expand_as(m)
t.storage()

10.0
20.0
30.0
[torch.FloatTensor of size 3]
```

А расширенный тензор имеет правильную форму, так как шаг измерения столбцов равен 0:

```
t.stride(), t.shape  
((1, 0), torch.Size([3, 3]))
```

Если при уширении требуется добавить измерения, это делается по умолчанию в самом начале. Когда мы делали уширение ранее, PyTorch выполнял `c.unsqueeze(0)` за кадром:

```
c = tensor([10.,20,30])  
c.shape, c.unsqueeze(0).shape,c.unsqueeze(1).shape  
  
(torch.Size([3]), torch.Size([1, 3]), torch.Size([3, 1]))
```

Команду `unsqueeze` можно заменить на индексирование `None`:

```
c.shape, c[None,:].shape,c[:,None].shape  
  
(torch.Size([3]), torch.Size([1, 3]), torch.Size([3, 1]))  
  
Двоеточия всегда можно опускать, и ... означает все предыдущие измерения:  
  
c[None].shape,c[...,None].shape  
  
(torch.Size([1, 3]), torch.Size([3, 1]))
```

Используя это, можно удалить еще один цикл `for` из нашей функции матричного умножения. Теперь вместо умножения `a[i]` на `b[:,j]` мы умножим `a[i]` на всю матрицу `b`, используя уширение, а затем сложим результаты:

```
def matmul(a,b):  
    ar,ac = a.shape  
    br,bc = b.shape  
    assert ac==br  
    c = torch.zeros(ar, bc)  
    for i in range(ar):  
        c[i,j] = (a[i,:]* b[:,j]).sum() # предыдущий  
        c[i]   = (a[i ].unsqueeze(-1) * b).sum(dim=0)  
    return c  
  
%timeit -n 20 t4 = matmul(m1,m2)
```

```
357 µs ± 7.2 µs per loop (mean ± std. dev. of 7 runs, 20 loops each)
```

Вот теперь мы ускорились в 3700 раз в сравнении с нашей начальной реализацией. Прежде чем продолжать, рассмотрим правила уширения несколько подробнее.

Правила уширения

При работе с двумя тензорами PyTorch сравнивает их формы поэлементно. Начинает он с конечных измерений и перемещается к началу, добавляя 1 в местах их отсутствия. Два измерения считаются *совместимыми*, если выполнено одно из следующих условий:

- Они равны.
- Одно из них равно 1, и в этом случае оно уширяется до формы другого измерения.

Массивам не обязательно иметь одинаковое количество измерений. Например, если у вас есть массив RGB-значений размером $256 \times 256 \times 3$ и требуется увеличить каждый цвет изображения на разное значение, то можно умножить это изображение на одномерный массив с тремя значениями. Выстраивание размеров конечных осей этих массивов согласно правилам уширения показывает их совместимость:

```
Image (3d tensor): 256 x 256 x 3
Scale (1d tensor): (1) (1) 3
Result (3d tensor): 256 x 256 x 3
```

Тем не менее двумерный тензор размером 256×256 окажется несовместим с этим изображением:

```
Image (3d tensor): 256 x 256 x 3
Scale (2d tensor): (1) 256 x 256
Error
```

В более ранних примерах с матрицей 3×3 и вектором с размером 3 уширение выполнялось для строк:

```
Matrix (2d tensor): 3 x 3
Vector (1d tensor): (1) 3
Result (2d tensor): 3 x 3
```

Попробуйте в качестве упражнения определить, какие измерения добавить (и куда), если нужно нормализовать пакет изображений размером $64 \times 3 \times 256 \times 256$ с векторами из трех элементов (один для среднего и один для стандартного отклонения).

Еще один удобный способ упростить умножение тензоров — использовать соглашение Эйнштейна о суммировании.

Соглашение Эйнштейна

Прежде чем использовать PyTorch-операцию @ или `torch.matmul`, попробуйте еще один способ реализации матричного умножения — соглашение Эйнштейна

о суммировании (`einsum`). Это компактное представление для совмещения произведений и сумм в общем. Уравнение записывается так:

`ik,kj -> ij`

Левая часть обозначает измерения операндов, разделенные запятой. Здесь у нас два тензора, каждый из которых имеет два измерения (i, k и k, j). Правая сторона представляет результирующие измерения, и здесь у нас получается тензор с двумя измерениями: i, j .

Правила соглашения о суммировании Эйнштейна.

1. Суммирование ведется по всем индексам, повторяющимся дважды в одном терме.
2. В каждом терме не может встречаться более двух одинаковых индексов.
3. Неповторяющиеся индексы с левой стороны должны встречаться с правой стороны.

Возвращаясь к нашему примеру, поскольку k повторяется, мы суммируем по этому индексу. В итоге формула представляет матрицу, получаемую при помножении в (i, j) суммы всех коэффициентов $(i.k)$ в первом тензоре, умноженных на коэффициенты $(k.j)$ во втором тензоре... что является матричным произведением.

В PyTorch это можно записать так:

```
def matmul(a,b): return torch.einsum('ik,kj->ij', a, b)
```

Соглашение Эйнштейна является очень практичным способом выражения операций, использующим индексирование и сумму произведений. Обратите внимание, что в левой части у вас может быть один член. Например:

```
torch.einsum('ij->ji', a)
```

Возвращает транспонированную матрицу a . При этом у вас также может быть три и более члена:

```
torch.einsum('bi,ij,bj->b', x, y, z)
```

Это действие вернет вектор размером b , где k -я координата является суммой $a[k,i] b[i,j] c[k,j]$. Эта нотация особенно удобна, когда из-за пакетов у вас получается больше измерений. Например, если у вас два пакета матриц и требуется вычислить матричное произведение для каждого, можно сделать так:

```
torch.einsum('bik,bkj->bij', x, y)
```

Давайте вернемся к нашей реализации `matmul`, задействовав `einsum`, и оценим ее скорость:

```
%timeit -n 20 t5 = matmul(m1,m2)
68.7 µs ± 4.06 µs per loop (mean ± std. dev. of 7 runs, 20 loops each)
```

Как видите, это не только практично, но и *очень* быстро. Einsum зачастую оказывается самым быстрым способом выполнения операций в PyTorch без погружения в C++ и CUDA. (Но обычно этот способ медленнее хорошо оптимизированного кода CUDA, что очевидно из раздела «Матричное умножение».)

Теперь, когда мы знаем, как реализовывать матричное умножение с чистого листа, можно переходить к построению нейронной сети, а именно ее прямого и обратного проходов с помощью одного только матричного умножения.

Прямой и обратный проход

В главе 4 мы видели, что для обучения модели необходимо вычислять градиенты полученных потерь в отношении ее параметров, что называется *обратным проходом*. В *прямом проходе* мы вычисляем выход модели для заданного входа на основе матричных произведений. При определении первой нейронной сети мы также углубимся в проблему правильной инициализации весов, что чрезвычайно важно для корректного начала обучения.

Определение и инициализация слоя

Для начала возьмем в качестве примера двухслойную нейронную сеть. Нам известно, что один слой можно выразить как $y = x @ w + b$, где x — это входные данные, y — выходные, w представляет веса слоя (размер которого будет равен количеству входных данных, умноженному на количество нейронов, если только мы не выполним транспонирование, как ранее), а b означает вектор смещений:

```
def lin(x, w, b): return x @ w + b
```

Второй слой можно надстроить над первым, но так как математически композиция двух линейных операций является еще одной линейной операцией, это имеет смысл, только если мы поместим между ними нелинейную операцию, называемую функцией активации. Как мы говорили в начале главы, в глубоком обучении в качестве такой функции чаще всего используется ReLU, которая возвращает максимум из x и 0 .

В этой главе мы не будем обучать модель, поэтому используем для входных данных и целей случайные тензоры. Предположим, на входе у нас 200 векторов размером 100, которые мы группируем в один пакет. В качестве целей мы задействуем 200 случайных чисел с плавающей запятой:

```
x = torch.randn(200, 100)
y = torch.randn(200)
```

Для нашей двухслойной модели потребуется две матрицы весов и два вектора смещений. Предположим, что скрытый размер будет равен 50, а размер выхода — 1 (в этом игрушечном примере для одного нашего входа соответствующий выход будет одним числом с плавающей запятой). Веса мы инициализируем случайные, а смещение с нулем.

```
w1 = torch.randn(100,50)
b1 = torch.zeros(50)
w2 = torch.randn(50,1)
b2 = torch.zeros(1)
```

Тогда результатом первого слоя будет:

```
l1 = lin(x, w1, b1)
l1.shape

torch.Size([200, 50])
```

Обратите внимание, что эта формула работает с пакетом входных данных и возвращает пакет скрытого состояния: `l1` — это матрица размером 200 (размер пакета) на 50 (скрытый размер).

Тем не менее в инициализации модели есть проблема. Чтобы ее понять, нужно взглянуть на среднее и стандартное отклонение (`std`) `l1`:

```
l1.mean(), l1.std()

(tensor(0.0019), tensor(10.1058))
```

Среднее близко к нулю, что логично, так как средние матрицы входных данных и матрицы весов близки к нулю. Но стандартное отклонение, которое показывает, насколько далеки наши активации от среднего, переместилось с 1 к 10. Это серьезная проблема, потому что произошла уже в одном слое. В современных же нейронных сетях могут присутствовать сотни слоев, и если каждый из них будет увеличивать масштаб активаций в десять раз, то к концу последнего слоя мы получим число, которое компьютер не сможет выразить.

К примеру, если мы выполним 50 умножений между `x` и случайными матрицами размером 100×100 , то получим:

```
x = torch.randn(200, 100)
for i in range(50): x = x @ torch.randn(100,100)
x[0:5,0:5]

tensor([[nan, nan, nan, nan, nan],
       [nan, nan, nan, nan, nan],
```

```
[nan, nan, nan, nan, nan],
[nan, nan, nan, nan, nan],
[nan, nan, nan, nan, nan]])
```

Результатом везде стало значение `nan`. Может, масштаб матриц был слишком велик и нужно использовать меньшие веса? Но если мы задействуем слишком низкие веса, то получим обратную проблему — масштаб активаций уменьшится с 1 до 0,1, и спустя 50 слоев у нас везде останутся нули:

```
x = torch.randn(200, 100)
for i in range(50): x = x @ (torch.randn(100,100) * 0.01)
x[0:5,0:5]

tensor([[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.]])
```

Поэтому нам нужно масштабировать матрицы весов ровно так, чтобы стандартное отклонение активаций оставалось равным 1. Мы можем вычислить точное подходящее значение математически, как это демонстрируют Ксавье Глорот и Йошуа Бенджио (Xavier Glorot, Yoshua Bengio) в своей работе *Understanding the Difficulty of Training Deep Feedforward Neural Networks* (<https://oreil.ly/9tiTC>) («Понимание сложности обучения глубоких нейронных сетей прямого распространения»). Правильным масштабом для заданного слоя будет $1/\sqrt{n_{in}}$, где n_{in} представляет количество входных данных. В нашем случае, если используется 100 входных данных, нужно масштабировать матрицы весов в 0,1 раза:

```
x = torch.randn(200, 100)
for i in range(50): x = x @ (torch.randn(100,100) * 0.1)
x[0:5,0:5]

tensor([[ 0.7554,  0.6167, -0.1757, -1.5662,  0.5644],
       [-0.1987,  0.6292,  0.3283, -1.1538,  0.5416],
       [ 0.6106,  0.2556, -0.0618, -0.9463,  0.4445],
       [ 0.4484,  0.7144,  0.1164, -0.8626,  0.4413],
       [ 0.3463,  0.5930,  0.3375, -0.9486,  0.5643]])
```

Вот теперь мы получили адекватные числа. Обратите внимание, насколько стабилен масштаб активаций даже спустя 50 фиктивных слоев:

```
x.std()
tensor(0.7042)
```

Если немного поэкспериментировать со значением для масштабирования, то можно заметить, что даже совсем небольшое отклонение от 0,1 приведет либо к слишком маленьким числам, либо к слишком большим. Поэтому очень важно правильно инициализировать веса.

Вернемся к нашей нейронной сети. Поскольку мы внесли некоторый беспорядок во входные данные, нужно их переопределить:

```
x = torch.randn(200, 100)
y = torch.randn(200)
```

И для весов мы используем правильный масштаб, известный как *инициализация Ксавье*:

```
from math import sqrt
w1 = torch.randn(100,50) / sqrt(100)
b1 = torch.zeros(50)
w2 = torch.randn(50,1) / sqrt(50)
b2 = torch.zeros(1)
```

Если теперь мы вычислим результат первого слоя, то можем убедиться, что среднее и стандартное отклонение находятся под контролем:

```
l1 = lin(x, w1, b1)
l1.mean(), l1.std()
(tensor(-0.0050), tensor(1.0000))
```

Очень хорошо! Далее нужно пройти через ReLU, так что определим ее. ReLU удаляет отрицательные значения, заменяя их нулями:

```
def relu(x): return x.clamp_min(0.)
```

Активации мы передаем через:

```
l2 = relu(l1)
l2.mean(), l2.std()
(tensor(0.3961), tensor(0.5783))
```

Тут мы вернулись к исходной точке: среднее активаций стало 0,4 (что можно понять, так как мы удалили отрицательные значения), а std опустилось до 0,58. Так что, как и в прошлом случае, через несколько слоев мы, скорее всего, получим нули:

```
x = torch.randn(200, 100)
for i in range(50): x = relu(x @ (torch.randn(100,100) * 0.1))
x[0:5,0:5]

tensor([[0.0000e+00, 1.9689e-08, 4.2820e-08, 0.0000e+00, 0.0000e+00],
       [0.0000e+00, 1.6701e-08, 4.3501e-08, 0.0000e+00, 0.0000e+00],
       [0.0000e+00, 1.0976e-08, 3.0411e-08, 0.0000e+00, 0.0000e+00],
       [0.0000e+00, 1.8457e-08, 4.9469e-08, 0.0000e+00, 0.0000e+00],
       [0.0000e+00, 1.9949e-08, 4.1643e-08, 0.0000e+00, 0.0000e+00]])
```

Значит, инициализация была ошибочной. Но почему? Когда Глорот и Бенжио писали свою статью, самой популярной функцией активации в нейронной сети

был гиперболический тангенс (\tanh , который они и использовали), и такая инициализация не подходит для нашей ReLU. Тем не менее нам повезло, так как другие ученые произвели необходимые математические вычисления правильного масштаба. В работе *Delving Deep into Rectifiers: Surpassing Human-Level Performance* (<https://oreil.ly/-qua>) («Углубление в ректификаторы: алгоритм превосходит человека») (которую мы уже видели — это работа, в которой была впервые представлена ResNet). Кайминг Хе и др. показывают, что нужно использовать следующий масштаб: $\sqrt{2/n_{in}}$, где n_{in} — это количество входных данных модели. Посмотрим, что мы получим:

```
x = torch.randn(200, 100)
for i in range(50): x = relu(x @ (torch.randn(100,100) * sqrt(2/100)))
x[0:5,0:5]

tensor([[0.2871, 0.0000, 0.0000, 0.0000, 0.0026],
       [0.4546, 0.0000, 0.0000, 0.0000, 0.0015],
       [0.6178, 0.0000, 0.0000, 0.0180, 0.0079],
       [0.3333, 0.0000, 0.0000, 0.0545, 0.0000],
       [0.1940, 0.0000, 0.0000, 0.0000, 0.0096]])
```

Уже лучше: числа на этот раз не обнулились. Так что вернемся к определению нейронной сети и используем именно эту инициализацию (которая в честь своих авторов названа *инициализацией Кайминга, или инициализацией Xe*):

```
x = torch.randn(200, 100)
y = torch.randn(200)

w1 = torch.randn(100,50) * sqrt(2 / 100)
b1 = torch.zeros(50)
w2 = torch.randn(50,1) * sqrt(2 / 50)
b2 = torch.zeros(1)
```

Взглянем на масштаб активаций после прохождения первого линейного слоя и ReLU:

```
l1 = lin(x, w1, b1)
l2 = relu(l1) l2.mean(),
l2.std()

(tensor(0.5661), tensor(0.8339))
```

Намного лучше! Теперь, когда веса у нас должным образом инициализированы, можно определить всю модель:

```
def model(x):
    l1 = lin(x, w1, b1) l2 = relu(l1)
    l3 = lin(l2, w2, b2)
    return l3
```

Это прямой проход. Далее осталось лишь сравнить выход модели с метками (в данном примере это случайные числа) с помощью функции потерь. В этом случае ею

будет среднеквадратичная ошибка. (Это не игрушечная задача, и данная функция является самой простой для последующего этапа — вычисления градиентов.)

Единственная тонкость в том, что выходные данные и цели имеют разную форму — после прохода через модель мы получаем следующий выход:

```
out = model(x)
out.shape

torch.Size([200, 1])
```

От конечного измерения 1 мы избавимся с помощью функции `squeeze`:

```
def mse(output, targ): return (output.squeeze(-1) - targ).pow(2).mean()
```

Вот теперь можно вычислять потери:

```
loss = mse(out, y)
```

На этом прямой проход заканчивается, так что взглянем на наши градиенты.

Градиенты и обратный проход

Мы видели, что PyTorch вычисляет все нужные градиенты с помощью вызова `loss.backward`, но давайте заглянем за кулисы и разберемся, что же именно там происходит.

Настал момент, когда нам нужно вычислить градиенты потерь в отношении всех весов модели, то есть все числа с плавающей запятой в `w1`, `b1`, `w2` и `b2`. Для этого нам понадобится обратиться к математике, а именно к *цепному правилу (дифференцированию сложной функции)*. Это правило относится к математическому анализу и определяет вычисление производной сложной функции:

$$(g \circ f)'(x) = g'(f(x))f'(x).$$

СЛОВО ДЖЕРЕМИ



Мне лично сложно в полной мере понять эту нотацию, поэтому я представляю ее так: если $y = g(u)$ и $u = f(x)$, тогда $dy/dx = dy/du * du/dx$. Обе эти записи равнозначны, так что можете использовать более, на ваш взгляд, удобную.

Наши потери представляют собой сочетание разных функций: среднеквадратичной ошибки (которая, в свою очередь, состоит из среднего и степени числа 2), второго линейного слоя, ReLU и первого линейного слоя. Например, если нам нужно получить градиенты для `b2` и потери определяются следующим выражением:

```
loss = mse(out,y) = mse(lin(12, w2, b2), y)
```

согласно правилу цепи получаем:

$$\frac{dloss}{db_2} = \frac{dloss}{dout} \times \frac{dout}{db_2} = \frac{d}{dout} mse(out, y) \times \frac{d}{db_2} lin(l_2, w_2, b_2).$$

Чтобы вычислить градиенты в отношении b_2 , сначала нужно получить градиенты в отношении выхода out . То же самое касается и вычисления градиентов в отношении w_2 . Тогда для получения градиентов в отношении b_1 или w_1 понадобятся градиенты в отношении l_1 , что, в свою очередь, потребует получить градиенты для l_2 , которой понадобятся градиенты для out .

Поэтому чтобы вычислить все необходимые для обновления градиенты, нужно начать с выхода модели и слой за слоем перемещаться *назад*, — вот почему этот этап называется *обратным распространением*. Его можно автоматизировать, если обеспечить предоставление каждой реализованной нами функцией (`relu`, `mse`, `lin`) своего обратного шага, то есть способа получения градиентов потерь в отношении входа(-ов) из градиентов потерь в отношении выхода.

Здесь мы заполняем эти градиенты в атрибуте каждого тензора. Нечто похожее делает PyTorch с помощью `.grad`.

Первыми идут градиенты потерь в отношении выхода модели (являющиеся входом для функции потерь). Мы отменяем `squeeze`, которую выполнили в `mse`, и затем используем формулу, которая получает производную от $x^2: 2x$. Эта производная среднего представляет собой просто l / n , где n — это число элементов на входе:

```
def mse_grad(inp, targ):
    # градиент потерь в отношении выхода предыдущего слоя
    inp.g = 2. * (inp.squeeze() - targ).unsqueeze(-1) / inp.shape[0]
```

Для получения градиентов ReLU и линейного слоя мы используем градиенты в отношении выхода (в `out.g`) и применяем цепное правило для вычисления градиентов в отношении входа (в `inp.g`). Цепное правило гласит, что `inp.g = relu'(inp.g) * out.g`. Производная `relu` равна либо 0 (когда входы отрицательные), либо 1 (когда входы положительные). В итоге мы получаем:

```
def relu_grad(inp, out):
    # градиенты relu в отношении входных активаций
    inp.g = (inp>0).float() * out.g
```

Та же схема применяется для вычисления градиентов в отношении входных данных, весов и смещения в линейном слое:

```
def lin_grad(inp, out, w, b):
    # градиенты matmul в отношении входных данных
    inp.g = out.g @ w.t()
    w.g = inp.t() @ out.g
    b.g = out.g.sum(0)
```

Мы не будем останавливаться на математических формулах, которые их определяют, так как для наших задач они не столь важны. Если же эта тема вас интересует, обязательно ознакомьтесь с уроками по математическому анализу в Академии Хана.

SYMPY

SymPy — это библиотека для символьных вычислений, которая очень пригодится при работе с математическим анализом. Вот выдержка из документации (<https://oreil.ly/i1IK9>): «Символьные вычисления подразумевают преобразования и работу с математическими объектами как с последовательностью символов. Это означает, что математические объекты представляются в точной форме, а не приблизительной, а математические выражения с невычисленными переменными остаются в символьной форме».

Для выполнения символьных вычислений сначала мы определяем *символ*, после чего делаем расчеты:

```
from sympy import symbols, diff
sx, sy = symbols('sx sy')
diff(sx**2, sx)

2*sx
```

Здесь SymPy получила за нас производную от sx^{**2} . Эта библиотека может получать производную сложных выражений, упрощать и масштабировать уравнения, а также многое другое.

Сегодня нет особой необходимости заниматься подобными математическими вычислениями самостоятельно — для нахождения градиентов за нас это делает PyTorch, а для показа уравнений — SymPy.

Определив все эти функции, можно использовать их для реализации обратного прохода. Поскольку каждый градиент автоматически заполняется в нужном тензоре, нам не требуется сохранять результаты для функций `_grad` — нужно лишь выполнить их в порядке, обратном прямому проходу, чтобы убедиться, что в каждой функции `out.g` присутствует:

```
def forward_and_backward(inp, targ):
    # прямой проход:
    l1 = inp @ w1 + b1
    l2 = relu(l1)
    out = l2 @ w2 + b2
    # при обратном проходе нам не требуются потери!
    loss = mse(out, targ)
```

```
# обратный проход:
mse_grad(out, targ)
lin_grad(l2, out, w2, b2)
relu_grad(l1, l2)
lin_grad(inp, l1, w1, b1)
```

Теперь можно обратиться к градиентам параметров модели в `w1.g`, `b1.g`, `w2.g` и `b2.g`. На этом мы успешно закончили определение модели — теперь пора сделать ее больше похожей на модуль PyTorch.

Рефакторинг модели

С тремя использованными нами выше функциями связаны еще две функции: прямой проход и обратный проход. Вместо того чтобы писать их раздельно, мы создадим для них обертывающий класс. Этот класс также будет хранить входные и выходные данные для обратного прохода, так что нам просто потребуется вызывать `backward`:

```
class Relu():
    def __call__(self, inp):
        self.inp = inp
        self.out = inp.clamp_min(0.)
        return self.out

def backward(self): self.inp.g = (self.inp>0).float() * self.out.g
```

`__call__` — это «магическое» имя в Python, которое делает наш класс вызываемым. Именно эта инструкция будет выполнена, когда мы наберем `y = Relu()(x)`. То же самое можно сделать для линейного слоя и потерь MSE.

```
class Lin():
    def __init__(self, w, b): self.w, self.b = w,b

    def __call__(self, inp):
        self.inp = inp
        self.out = inp@self.w + self.b
        return self.out

    def backward(self):
        self.inp.g = self.out.g @ self.w.t()
        self.w.g = self.inp.t() @ self.out.g
        self.b.g = self.out.g.sum(0)

class Mse():
    def __call__(self, inp, targ):
        self.inp = inp
        self.targ = targ
```

```

        self.out = (inp.squeeze() - targ).pow(2).mean()
        return self.out

    def backward(self):
        x = (self.inp.squeeze()-self.targ).unsqueeze(-1)
        self.inp.g = 2.*x/self.targ.shape[0]

```

Затем поместим все в модель, которую инициализируем с тензорами $w1$, $b1$, $w2$ и $b2$:

```

class Model():
    def __init__(self, w1, b1, w2, b2):
        self.layers = [Lin(w1,b1), Relu(), Lin(w2,b2)]
        self.loss = Mse()

    def __call__(self, x, targ):
        for l in self.layers: x = l(x)
        return self.loss(x, targ)

    def backward(self):
        self.loss.backward()
        for l in reversed(self.layers): l.backward()

```

Этот рефакторинг и регистрация всего в виде слоев модели облегчают написание прямого и обратного проходов. Теперь для инстанцирования модели достаточно этого:

```
model = Model(w1, b1, w2, b2)
```

Далее прямой проход выполняется так:

```
loss = model(x, y)
```

А обратный так:

```
model.backward()
```

Переходим в PyTorch

Написанные нами классы `Lin`, `Mse` и `Relu` имеют много общего, поэтому мы можем создать их наследованием от одного базового класса:

```

class LayerFunction():
    def __call__(self, *args):
        self.args = args
        self.out = self.forward(*args)
        return self.out

    def forward(self):    raise Exception('not implemented')
    def bwd(self):       raise Exception('not implemented')
    def backward(self):  self.bwd(self.out, *self.args)

```

Затем нам просто нужно реализовать в каждом подклассе `forward` и `bwd`:

```
class Relu(LayerFunction):
    def forward(self, inp): return inp.clamp_min(0.)
    def bwd(self, out, inp): inp.g = (inp>0).float() * out.g

class Lin(LayerFunction):
    def __init__(self, w, b): self.w, self.b = w, b

    def forward(self, inp): return inp@self.w + self.b

    def bwd(self, out, inp):
        inp.g = out.g @ self.w.t()
        self.w.g = self.inp.t() @ self.out.g
        self.b.g = out.g.sum(0)

class Mse(LayerFunction):
    def forward (self, inp, targ): return (inp.squeeze() - targ).pow(2).mean()
    def bwd(self, out, inp, targ):
        inp.g = 2*(inp.squeeze()-targ).unsqueeze(-1) / targ.shape[0]
```

Остальную часть модели можно оставить прежней. Так мы постепенно приближаемся к тому, что делает PyTorch. Каждая базовая функция, которую нам нужно дифференцировать, написана как объект `torch.autograd.Function`, содержащий методы `forward` и `backward`. После этого PyTorch будет отслеживать все наши вычисления для правильного выполнения обратного прохода до тех пор, пока мы не установим в тензорах атрибут `requires_grad` как `False`.

Написать такой класс (почти) так же легко, как исходный. Разница в том, что мы выбираем, что сохранять, а что помещать в переменные контекста (избегая сохранения ненужных данных), и возвращаем градиенты в проходе `backward`.

Писать собственную `Function` приходится редко, но если вам однажды понадобится нечто экзотическое или вы захотите поработать с градиентами стандартной функции, то написать ее можно так:

```
from torch.autograd import Function

class MyRelu(Function):
    @staticmethod
    def forward(ctx, i):
        result = i.clamp_min(0.)
        ctx.save_for_backward(i)
        return result

    @staticmethod
    def backward(ctx, grad_output):
        i, = ctx.saved_tensors
        return grad_output * (i>0).float()
```

В качестве структуры для построения более сложной модели, использующей эти `Function`, применяется `torch.nn.Module`. Это базовая структура для всех моделей, и все нейронные сети, которые вы видели до сих пор, принадлежали к этому классу. Он помогает регистрировать все обучаемые параметры, которые, как мы видели, могут использоваться в цикле обучения.

Для реализации `nn.Module` нужно сделать следующее.

1. Обеспечить, чтобы при инициализации суперкласса сначала вызывался метод `__init__`.
2. Определить все параметры модели в качестве атрибутов с помощью `nn.Parameter`.
3. Определить функцию `forward`, возвращающую выход модели.

Вот пример линейного слоя, написанного с нуля:

```
import torch.nn as nn

class LinearLayer(nn.Module):
    def __init__(self, n_in, n_out):
        super().__init__()
        self.weight = nn.Parameter(torch.randn(n_out, n_in) * sqrt(2/n_in))
        self.bias = nn.Parameter(torch.zeros(n_out))

    def forward(self, x): return x @ self.weight.t() + self.bias
```

Этот класс автоматически отслеживает, какие параметры определены:

```
lin = LinearLayer(10,2)
p1,p2 = lin.parameters()
p1.shape,p2.shape

(torch.Size([2, 10]), torch.Size([2]))
```

Именно благодаря этой особенности `nn.Module` мы можем просто дать команду `opt.step`, и оптимизатор переберет все параметры, обновив каждый из них.

Обратите внимание, что в PyTorch веса сохраняются как матрица `n_out x n_in`, поэтому при прямом проходе мы выполняем транспонирование.

При использовании линейного слоя из PyTorch (который тоже применяет инициализацию Кайминга) создаваемую нами в этой главе модель можно записать так:

```
class Model(nn.Module):
    def __init__(self, n_in, nh, n_out):
        super().__init__()
        self.layers = nn.Sequential(
```

```

        nn.Linear(n_in,nh), nn.ReLU(), nn.Linear(nh,n_out))
self.loss = mse

def forward(self, x, targ): return self.loss(self.layers(x).squeeze(), targ)

```

fastai предоставляет свой собственный вариант `Module`, который идентичен `nn.Module`, но не требует от вас вызова `super().__init__()`, так как делает это автоматически:

```

class Model(Module):
    def __init__(self, n_in, nh, n_out):
        self.layers = nn.Sequential(
            nn.Linear(n_in,nh), nn.ReLU(), nn.Linear(nh,n_out))
        self.loss = mse

    def forward(self, x, targ): return self.loss(self.layers(x).squeeze(), targ)

```

В главе 19 мы начнем с подобной модели и рассмотрим построение цикла обучения с нуля с последующим его рефакторингом согласно пройденному материалу.

Резюме

В этой главе мы изучили основы глубокого обучения, начав с матричного умножения и рассмотрев реализацию прямого и обратного проходов нейронной сети с нуля. После этого мы реорганизовали код, чтобы показать, как работает PyTorch изнутри.

Вот основное, что нужно запомнить.

- Нейронная сеть в своей основе является набором матричных умножений, перемежающихся с нелинейными функциями.
- Python медленный, поэтому для написания быстрого кода нужно перевести его в векторы, после чего задействовать техники поэлементной арифметики и уширения.
- Уширение одного тензора в отношении другого возможно, если их измерения от конца к началу совпадают (если они одинаковы или одно из них представлено как 1). Чтобы сделать уширение тензора возможным, нужно добавить измерения размером 1 с помощью функции `unsqueeze` или индекса `None`.
- Для начала обучения очень важно правильно инициализировать нейронную сеть. При использовании нелинейных функций ReLU нужно применять инициализацию Кайминга.
- Обратный проход представляет собой итеративный метод обратного распространения ошибки, при котором градиенты вычисляются на основе выхода модели, после чего происходит возвращение к входу через все ее слои.

- При создании подкласса `nn.Module` (если не используется `Module` `fastai`) нужно вызывать из нашего метода `__init__` метод `__init__` суперкласса, а также определить функцию `forward`, получающую входные данные и возвращающую нужный выход.

Вопросник

1. Напишите код Python для реализации одного нейрона.
2. Напишите код Python для реализации ReLU.
3. Напишите код Python для линейного слоя в форме матричного умножения.
4. Напишите код Python для линейного слоя в простом Python (то есть с помощью генераторов списков и его встроенной функциональности).
5. Что такое «скрытый размер» слоя?
6. За что в Python отвечает метод `t`?
7. Почему написанное на Python матричное умножение выполняется очень медленно?
8. Почему в `matmul ac==br`?
9. Как в блокноте Jupyter измерить время выполнения одной ячейки?
10. Что такое поэлементная арифметика?
11. Напишите код Python для проверки, больше ли каждый элемент из `a` соответствующего ему элемента из `b`.
12. Что такое тензор ранга 0? Как преобразовать его в тип данных простого Python?
13. Что вернет этот код и почему?
`tensor([1,2]) + tensor([1])`

14. Что вернет этот код и почему?
`tensor([1,2]) + tensor([1,2,3])`

15. Как поэлементная арифметика помогает ускорить `matmul`?
16. Назовите условия для выполнения уширения.
17. Что такое `expand_as`? Приведите пример его использования для сопоставления результатов уширения.
18. Как `unsqueeze` помогает решить определенные сложности уширения?
19. Как можно использовать индексирование для выполнения операции альтернативной `unsqueeze`?

20. Как просмотреть фактическое содержимое памяти, используемое для тензора?
21. При сложении вектора размером 3 и матрицы размером 3×3 элементы вектора прибавляются к каждой строке или к каждому столбцу матрицы? (Проверьте ответ, выполнив этот код в блокноте.)
22. Увеличивают ли потребление памяти операции уширения и `expand_as`? Почему?
23. Реализуйте `matmul`, используя суммирование Эйнштейна.
24. Что обозначает повторяющаяся буква индекса в левой стороне `einsum`?
25. Назовите три правила выражения суммирования Эйнштейна и объясните их смысл.
26. Что такое прямой и обратный проход нейронной сети?
27. Почему при прямом проходе нужно сохранять активации, вычисленные для промежуточных слоев?
28. В чем выражается проблема, когда стандартное отклонение активаций уходит от 1?
29. Как в этом случае может помочь инициализация весов?
30. Какой формулой мы инициализируем веса, чтобы получить стандартное отклонение 1 для простого линейного слоя и линейного слоя, сопровождаемого ReLU?
31. Почему нам иногда приходится использовать в функции потерь метод `squeeze`?
32. Что делает аргумент метода `squeeze`? Почему добавление этого аргумента может оказаться важным, несмотря на то что PyTorch его не требует?
33. Что такое цепное правило? Приведите уравнение в одной из двух представленных в главе форм.
34. Вычислите градиенты `mse(lin(12, w2, b2), y)`, используя цепное правило.
35. Что такое градиент ReLU? Приведите его в математической форме или в коде. (Это не обязательно запоминать — попробуйте разобраться, используя знания о формах функции.)
36. В каком порядке нужно вызывать функции `*_grad` при обратном проходе? Почему?
37. Что такое `__call__`?
38. Какие методы нужно реализовать при написании `torch.autograd.Function`?
39. Напишите `nn.Linear` с нуля и проверьте, как он работает.
40. Чем отличается `nn.Module` от `Module` в fastai?

Дополнительные задания

1. Реализуйте ReLU как `torch.autograd.Function` и обучите с ее помощью модель.
2. Если вы хороши в математике, определите градиенты линейного слоя в математическом выражении и сопоставьте с реализацией из данной главы.
3. Изучите PyTorch-метод `unfold` и используйте его вместе с матричным умножением для реализации собственной функции двумерной свертки, после чего обучите использующую ее CNN.
4. Реализуйте все материалы главы, используя NumPy вместо PyTorch.

ГЛАВА 18

Интерпретация CNN с помощью CAM

Теперь мы уже многое умеем создавать с нуля, так что используем эти знания для создания совершенно новой и очень полезной функциональности: *карты активаций класса*. С ее помощью мы сможем лучше понять, почему CNN дает те или иные прогнозы.

В течение этого процесса мы познакомимся с одной ранее нам не встречавшейся возможностью PyTorch — *хуком*, а также задействуем ряд концепций, которые будут введены в оставшейся части книги. Если вам в конце этой главы захочется проверить знания всех пройденных материалов, попробуйте убрать книгу в сторонку и воссоздать рассмотренные здесь идеи самостоятельно, не подглядывая.

CAM и хуки

Карта активаций класса (class activation map CAM) была представлена Болей Чжоу (Bolei Zhou) и др. в работе *Learning Deep Features for Discriminative Localization* (<https://oreil.ly/5hik3>) («Обучение глубоких признаков для дискриминативной локализации»). В ней используется выход последнего сверточного слоя (сразу перед слоем среднего пулинга) совместно с прогнозами, на основе чего формируется изображение тепловой карты, отражающей причины принятия моделью ее решений. Это очень помогает при интерпретации.

Говоря точнее, в каждой позиции последнего сверточного слоя у нас столько же фильтров, сколько в последнем линейном слое. Следовательно, мы можем вычислить скалярное произведение их активаций с конечными весами, получив для каждого участка карты интенсивность признаков, использованных для принятия решения.

Для этого нам понадобится способ получить доступ к активациям внутри модели в процессе ее обучения. В PyTorch это можно сделать с помощью *хука*. Хуки

представляют собой эквивалент обратных вызовов в fastai. Тем не менее вместо внедрения кода в цикл обучения, как это делает обратный вызов `Learner`, хуки дают возможность внедрить код в вычисления непосредственно прямого и обратного проходов. Мы можем прикрепить хук к любому слою модели, и он будет выполнен при вычислении выходов (прямой хук) или в процессе обратного распространения (обратный хук). Прямой хук — это функция, получающая три компонента: модуль, его вход и выход. При этом она может реализовывать любое необходимое поведение. (fastai тоже предоставляет удобный `HookCallback`, несколько упрощающий работу с хуками; здесь мы его рассматривать не будем, но рекомендуем к самостоятельному изучению в документации.)

Все описанное мы продемонстрируем на той же модели кошек и собак, которую обучали в главе 1:

```
path = untar_data(URLs.PETS) / 'images'
def is_cat(x): return x[0].isupper()
dls = ImageDataLoaders.from_name_func(
    path, get_image_files(path), valid_pct=0.2, seed=21,
    label_func=is_cat, item_tfms=Resize(224))
learn = cnn_learner(dls, resnet34, metrics=error_rate)
learn.fine_tune(1)
```

epoch	train_loss	valid_loss	error_rate	time
0	0.141987	0.018823	0.007442	00:16

epoch	train_loss	valid_loss	error_rate	time
0	0.050934	0.015366	0.006766	00:21

Для начала выберем фото кошки и пакет данных:

```
img = PILImage.create('images/chapter1_cat_example.jpg')
x, = first(dls.test_dl([img]))
```

Для CAM нам нужно сохранять активации последнего сверточного слоя. Мы помещаем функцию хука в класс, определяя для него состояние, к которому мы сможем обратиться позже, а затем просто сохраним копию выхода:

```
class Hook():
    def hook_func(self, m, i, o): self.stored = o.detach().clone()
```

После этого можно инстанцировать `Hook` и прикрепить его к последнему слою тела CNN:

```
hook_output = Hook()
hook = learn.model[0].register_forward_hook(hook_output.hook_func)
```

Теперь извлекаем пакет и передаем его через модель:

```
with torch.no_grad(): output = learn.model.eval()(x)
```

Обращаемся к сохраненным активациям:

```
act = hook_output.stored[0]
```

И перепроверяем прогнозы:

```
F.softmax(output, dim=-1)
```

```
tensor([[7.3566e-07, 1.0000e+00]], device='cuda:0')
```

Нам известно, что 0 (для `False`) означает *dog*, потому что в `fastai` классы сортируются автоматически, но мы можем перепроверить это, заглянув в `dls.vocab`:

```
dls.vocab
```

```
(#2) [False, True]
```

Наша модель очень уверена, что на фото была кошка.

Чтобы получить скалярное произведение матрицы весов ($2 * \text{число активаций}$) и активаций (размер пакета * строки * столбцы), мы задействуем настраиваемое `einsum`:

```
act.shape
```

```
torch.Size([512, 7, 7])
```

```
cam_map = torch.einsum('ck,kij->ciij', learn.model[1][-1].weight, act)
cam_map.shape
```

```
torch.Size([2, 7, 7])
```

Для каждого изображения пакета и для каждого класса мы получаем карту признаков 7×7 , по которой видно, где активации были выше, а где ниже. Это позволит нам понять, какие области изображений повлияли на решение модели.

Например, мы можем выяснить, на основе каких областей модель решила, что перед ней кошка (обратите внимание, что нам потребуется `decode` вход `x`, так как он был нормализован `DataLoader`, а также выполнить приведение к `TensorName`, поскольку на момент написания книги PyTorch не поддерживает типы при индексировании — к моменту прочтения вами этого материала данная проблема может быть исправлена):

```
x_dec = TensorImage(dls.train.decode((x,))[0][0])
_,ax = plt.subplots()
```

```
x_dec.show(ctx=ax)
ax.imshow(cam_map[1].detach().cpu(), alpha=0.6, extent=(0,224,224,0),
           interpolation='bilinear', cmap='magma');
```



В данном случае ярко-желтые области соответствуют высоким активациям, а фиолетовые — низким. По этой карте мы видим, что голова и передняя лапа оказались двумя основными областями, на основании которых модель решила, что это кошка.

Как только вы получили от хука все, что хотели, его стоит удалить, иначе может возникнуть утечка памяти:

```
hook.remove()
```

Поэтому, как правило, рекомендуется создать класс `Hook` в роли *менеджера контекста*, который будет регистрировать хук, когда вы в него входите, и удалять, когда выходите. Менеджер контекста — это конструкция Python, которая при создании объекта в инструкции `with` вызывает `__enter__`, а при завершении инструкции `with` происходит вызов `__exit__`. Например, вот как Python обрабатывает часто встречающуюся конструкцию `with open(...) as f`: для открывания файлов, не требуя явного указания `close(f)` в конце.

Если мы определим `Hook` так:

```
class Hook():
    def __init__(self, m):
        self.hook = m.register_forward_hook(self.hook_func)
    def hook_func(self, m, i, o): self.stored = o.detach().clone()
    def __enter__(self, *args): return self
    def __exit__(self, *args): self.hook.remove()
```

то можем спокойно использовать его так:

```
with Hook(learn.model[0]) as hook:
    with torch.no_grad(): output = learn.model.eval()(x.cuda())
act = hook.stored
```

fastai облегчает работу с хуками, предоставляя вам как показанный класс `Hook`, так и ряд других удобных классов.

Этот метод полезен, но работает только для последнего слоя. Тем не менее эту проблему решает САМ градиентов.

САМ градиентов

Только что рассмотренный метод позволяет нам вычислять только тепловую карту последних активаций, так как при наличии признаков нам нужно умножать их на последнюю матрицу весов. Такой способ не будет работать для внутренних слоев сети. Его вариация, представленная в 2016 году Рампраасатом Р. Сельвараджу (Ramprasaath R. Selvaraju) и др. в их работе *Grad-CAM: Why Did You Say That?* (<https://oreil.ly/4krXH>) («Grad-Cam: зачем вы это сказали?»), использует градиенты последней активации для нужного класса. Если вы помните тему обратного прохода, то градиенты выхода последнего слоя в отношении входа этого слоя равны весам слоя, поскольку он линейный.

В случае более глубоких слоев нам также нужны градиенты, но они уже не будут просто равны весам, и нам понадобится их вычислять. Градиенты каждого слоя за нас вычисляет PyTorch в процессе обратного прохода, но они не сохраняются (кроме как для тензоров, где `requires_grad = True`). Тем не менее мы можем зарегистрировать при обратном проходе хук, которому PyTorch передаст градиенты в качестве параметра, и мы их там сохраним. Для этого мы используем класс `HookBwd`, который работает аналогично `Hook`, но вместо активаций перехватывает и сохраняет градиенты:

```
class HookBwd():
    def __init__(self, m):
        self.hook = m.register_backward_hook(self.hook_func)
    def hook_func(self, m, gi, go): self.stored = go[0].detach().clone()
    def __enter__(self, *args): return self
    def __exit__(self, *args): self.hook.remove()
```

Затем для индекса класса 1 (для `True`, означающего *cat*) мы, как и ранее, перехватываем признаки последнего сверточного слоя и вычисляем градиенты выходных активаций класса. Мы не можем просто вызвать `output.backward`, потому что вычисление градиентов имеет смысл только в отношении скаляра (которым обычно являются потери), а `output` — это тензор ранга 2. Но если мы возьмем одно изображение (используем 0) и один класс (используем 1), то сможем с помощью `output[0,cls].backward` вычислить градиенты любого веса или активации в отношении этого одного значения. Наш хук перехватывает градиенты, которые мы задействуем в качестве весов:

```
cls = 1
with HookBwd(learn.model[0]) as hookg:
    with Hook(learn.model[0]) as hook:
        output = learn.model.eval()(x.cuda())
        act = hook.stored
output[0,cls].backward()
grad = hookg.stored
```

Веса для Grad-CAM задаются средним наших градиентов по карте признаков. Так мы получаем в точности то, что и прежде:

```
w = grad[0].mean(dim=[1,2], keepdim=True)
cam_map = (w * act[0]).sum(0)

_,ax = plt.subplots()
x_dec.show(ctx=ax)
ax.imshow(cam_map.detach().cpu(), alpha=0.6, extent=(0,224,224,0),
           interpolation='bilinear', cmap='magma');
```



Особенность Grad-CAM в том, что его можно использовать для любого слоя. Например, здесь мы применяем его для выхода предпоследней группы ResNet:

```
with HookBwd(learn.model[0][-2]) as hookg:
    with Hook(learn.model[0][-2]) as hook:
        output = learn.model.eval()(x.cuda())
        act = hook.stored
output[0,cls].backward()
grad = hookg.stored

w = grad[0].mean(dim=[1,2], keepdim=True)
cam_map = (w * act[0]).sum(0)
```

Теперь можно просмотреть карту активаций для этого слоя:

```
_,ax = plt.subplots()
x_dec.show(ctx=ax)
```

```
ax.imshow(cam_map.detach().cpu(), alpha=0.6, extent=(0,224,224,0),  
         interpolation='bilinear', cmap='magma');
```



Резюме

Интерпретация моделей является активно исследуемой областью, и мы затронули лишь малую часть материала. Карты активаций класса дают представление о том, почему модель спрогнозировала те или иные результаты, показывая области изображений, которые оказали на этот прогноз максимальное влияние. Эта информация помогает проанализировать ложноположительные результаты и выяснить, какие данные нужно дополнить, чтобы таких результатов избежать.

Вопросник

1. Что такое хук в PyTorch?
2. Выходы какого слоя использует CAM?
3. Зачем CAM нужен хук?
4. Просмотрите исходный код класса `ActivationStats` и разберитесь, как он использует хуки.
5. Напишите хук, который сохраняет активации заданного слоя модели (по возможности не подглядывайте).
6. Почему мы вызываем `eval` до получения активаций? Зачем мы используем `no_grad`?
7. Примените `torch.einsum` для вычисления количественного показателя `dog` или `cat` для каждого участка последней активации тела модели.
8. Как проверить, в каком порядке расположены категории (то есть соответствие индекса → категории)?
9. Почему при отображении входного изображения мы используем `decode`?

10. Что такое менеджер контекстов? Какие особые методы необходимо определить для его создания?
11. Почему нельзя использовать простой CAM для внутренних слоев сети?
12. Почему для выполнения Grad-CAM нужно регистрировать хук в обратном проходе?
13. Почему нельзя вызвать `output.backward`, когда `output` является тензором ранга 2 выходных активаций каждого изображения каждого класса?

Дополнительные задания

1. Удалите `keepdim`, и посмотрите, к чему это приведет. Найдите этот параметр в документации PyTorch. Зачем он нужен в этом блокноте?
2. Создайте блокнот, подобный этому, но для NLP и используйте его для определения слов отзыва на фильм, максимально повлиявших на определение его тональности.

ГЛАВА 19

Класс Learner с нуля

Эта последняя глава (без учета заключения и онлайн-глав) будет несколько отличаться по своей структуре от предыдущих, так как содержит намного больше кода, чем текста. Здесь мы представим новые ключевые слова и библиотеки Python, не углубляясь в их описание. При этом подразумевается, что эта глава положит начало вашему серьезному исследовательскому проекту. Мы собираемся реализовать многие ключевые элементы API fastai и PyTorch с чистого листа, используя для построения только компоненты, которые изучили в главе 17. Основная цель — получить в итоге собственный класс `Learner` и ряд обратных вызовов, чего будет достаточно для обучения модели на Imagenette, включая примеры ключевых техник, которые мы изучили. В процессе построения `Learner` мы создадим собственные версии `Module`, `Parameter` и параллельного `DataLoader`, чтобы у вас сформировалось устойчивое понимание выполняемых этими классами ролей.

Вопросник в этой главе будет особенно важен. В нем мы укажем многие интересные направления, в которых можно продолжать двигаться, начиная с этой главы. Советуем вам следовать изложенному здесь материалу, используя компьютер для экспериментирования, поиска по интернету и любых других действий, которые помогут вам лучше понять описываемые процессы. Вы уже достаточно хорошо освоили материал предыдущей части книги и однозначно готовы справиться с оставшимся!

Начнем с ручного сбора данных.

Данные

Загляните в исходный код `untar_data`, чтобы понять, как он работает. Здесь мы с его помощью получаем выборку Imagenette с размером изображений 160×160 пикселей:

```
path = untar_data(URLs.IMAGENETTE_160)
```

Для обращения к файлам изображений можно задействовать `get_image_files`:

```
t = get_image_files(path)
t[0]

Path('/home/jhoward/.fastai/data/imagenette2-160/val/n03417042/
n03417042_3752.JPEG')
```

Или получить то же самое, задействовав стандартный модуль Python `glob`:

```
from glob import glob
files = L(glob(f'{path}/**/*.{JPEG}', recursive=True)).map(Path)
files[0]

Path('/home/jhoward/.fastai/data/imagenette2-160/val/n03417042/
n03417042_3752.JPEG')
```

Если же посмотреть на исходный код `get_image_files`, то можно увидеть, что в нем используется `os.walk`. Эта функция более быстрая и более гибкая, чем `glob`, так что обязательно ее опробуйте.

Изображение можно открыть с помощью класса `Image` из Python Imaging Library (PIL):

```
im = Image.open(files[0])
im
```



```
im_t = tensor(im)
im_t.shape

torch.Size([160, 213, 3])
```

Это будет основа нашей независимой переменной. Для зависимой же переменной мы используем `Path.parent` из `pathlib`. Сначала нам понадобится словарь:

```
lbls = files.map(Self.parent.name()).unique(); lbls

(#10) ['n03417042', 'n03445777', 'n03888257', 'n03394916', 'n02979186', 'n03000684',
'n03425413', 'n01440764', 'n03028079', 'n02102040']
```

И обратное сопоставление, выполняемое `L.val2idx`:

```
v2i = lbls.val2idx(); v2i

{'n03417042': 0,
 'n03445777': 1,
 'n03888257': 2,
 'n03394916': 3,
 'n02979186': 4,
 'n03000684': 5,
 'n03425413': 6,
 'n01440764': 7,
 'n03028079': 8,
 'n02102040': 9}
```

Это все, что нам нужно для создания `Dataset`.

Dataset

`Dataset` в PyTorch может быть чем угодно, поддерживающим индексирование (`__getitem__`) и `len`:

```
class Dataset:
    def __init__(self, fns): self.fns=fns
    def __len__(self): return len(self.fns)
    def __getitem__(self, i):
        im = Image.open(self.fns[i]).resize((64,64)).convert('RGB')
        y = v2i[self.fns[i].parent.name]
        return tensor(im).float()/255, tensor(y)
```

Нам нужен список имен файлов для обучения и контроля, который мы передадим в `Dataset.__init__`:

```
train_filt = L(o.parent.parent.name=='train' for o in files)
train,valid = files[train_filt],files[~train_filt]
len(train),len(valid)

(9469, 3925)
```

Проверим его:

```
train_ds,valid_ds = Dataset(train),Dataset(valid)
x,y = train_ds[0]
x.shape,y

(torch.Size([64, 64, 3]), tensor(0))

show_image(x, title=lbls[y]);
```

n03417042



Как видите, датасет возвращает независимые и зависимые переменные в виде кортежа, что нам и нужно. Нам также нужна возможность объединить их в мини-пакет. Обычно это делается с помощью `torch.stack`, который мы здесь и используем:

```
def collate(idxs, ds):
    xb,yb = zip(*[ds[i] for i in idxs])
    return torch.stack(xb),torch.stack(yb)
```

Вот мини-пакет с двумя элементами для проверки `collate`:

```
x,y = collate([1,2], train_ds)
x.shape,y
(torch.Size([2, 64, 64, 3]), tensor([0, 0]))
```

Теперь, когда у нас есть датасет и функция объединения, можно переходить к созданию `DataLoader`. Здесь мы добавим еще два компонента: optionalный `shuffle` для обучающей выборки и `ProcessPoolExecutor` для параллельного препроцессинга.

Параллельный загрузчик данных очень важен, так как открывание и декодирование изображений JPEG – это медленный процесс. Один CPU не сможет достаточно быстро декодировать изображения, чтобы полноценно загрузить современный GPU. Вот наш класс `DataLoader`:

```
class DataLoader:
    def __init__(self, ds, bs=128, shuffle=False, n_workers=1):
        self.ds,self.bs,self.shuffle,self.n_workers = ds,bs,shuffle,n_workers

    def __len__(self): return (len(self.ds)-1)//self.bs+1

    def __iter__(self):
        idxs = L.range(self.ds)
        if self.shuffle: idxs = idxs.shuffle()
        chunks = [idxs[n:n+self.bs] for n in range(0, len(self.ds), self.bs)]
        with ProcessPoolExecutor(self.n_workers) as ex:
            yield from ex.map(collate, chunks, ds=self.ds)
```

Опробуем его на обучающей и контрольной выборках:

```
n_workers = min(16, defaults.cpus)
train_dl = DataLoader(train_ds, bs=128, shuffle=True, n_workers=n_workers)
valid_dl = DataLoader(valid_ds, bs=256, shuffle=False, n_workers=n_workers)
xb,yb = first(train_dl)
xb.shape,yb.shape,len(train_dl)

(torch.Size([128, 64, 64, 3]), torch.Size([128]), 74)
```

Этот загрузчик данных не намного медленнее PyTorch, но при этом намного проще. Так что если вам потребуется отладка сложного процесса загрузки данных, то не бойтесь делать ее вручную — это поможет лучше понять происходящее.

Для нормализации нам потребуются статистики изображений. Обычно их вполне можно рассчитать на одном обучающем мини-пакете, так как точность здесь не важна:

```
stats = [xb.mean((0,1,2)),xb.std((0,1,2))]
stats

[tensor([0.4544, 0.4453, 0.4141]), tensor([0.2812, 0.2766, 0.2981])]
```

Класс `Normalize` должен просто хранить эти статистики и применять их (чтобы понять назначение `to_device`, закомментируйте его и посмотрите, что произойдет позже в этом блокноте):

```
class Normalize:
    def __init__(self, stats): self.stats=stats
    def __call__(self, x):
        if x.device != self.stats[0].device:
            self.stats = to_device(self.stats, x.device)
        return (x-self.stats[0])/self.stats[1]
```

Нам нравится проверять все создаваемое сразу в блокноте:

```
norm = Normalize(stats)
def tfm_x(x): return norm(x).permute((0,3,1,2))
t = tfm_x(x)
t.mean((0,2,3)),t.std((0,2,3))

(tensor([0.3732, 0.4907, 0.5633]), tensor([1.0212, 1.0311, 1.0131]))
```

Здесь `tfm_x` не просто применяет `Normalize`, она также пермутирует порядок осей из NHWC в NCHW (эти акронимы описывались в главе 13). PIL использует порядок осей HWC, который мы не можем применить в PyTorch, поэтому и требуется `permute`.

Это все, что касается подготовки данных для модели. Теперь нам нужна сама модель!

Module и Parameter

Для создания модели нам понадобится `Module`, а чтобы создать `Module` — `Parameter`, так что с него мы и начнем. В главе 8 мы говорили, что класс `Parameter` «не добавляет функциональности (кроме автоматического вызова `requires_grad`) и используется только как “маркер”, показывающий, что требуется включить в `parameters`». Вот определение, которое именно это и делает:

```
class Parameter(Tensor):
    def __new__(self, x): return Tensor._make_subclass(Parameter, x, True)
    def __init__(self, *args, **kwargs): self.requires_grad_()
```

Данная реализация немного странновата: нам приходится определить особый метод Python `__new__` и использовать внутренний метод PyTorch `_make_subclass`, потому что на момент написания книги PyTorch иначе не работает корректно при таком способе создания подкласса и не предоставляет для этого официально поддерживаемый API. К моменту прочтения вами этих слов данная проблема может быть исправлена, так что загляните на сайт книги, где об этом будет обязательно объявлено.

Теперь наш `Parameter` ведет себя как тензор, чего мы и добивались:

```
Parameter(tensor(3.))
tensor(3., requires_grad=True)
```

Далее можно перейти к определению `Module`:

```
class Module:
    def __init__(self):
        self.hook, self.params, self.children, self._training = None, [], [], False

    def register_parameters(self, *ps): self.params += ps
    def register_modules(self, *ms): self.children += ms

    @property
    def training(self): return self._training
    @training.setter
    def training(self,v):
        self._training = v
        for m in self.children: m.training=v

    def parameters(self):
        return self.params + sum([m.parameters() for m in self.children], [])

    def __setattr__(self,k,v):
        super().__setattr__(k,v)
        if isinstance(v,Parameter): self.register_parameters(v)
        if isinstance(v,Module): self.register_modules(v)
```

```

def __call__(self, *args, **kwargs):
    res = self.forward(*args, **kwargs)
    if self.hook is not None: self.hook(res, args)
    return res

def cuda(self):
    for p in self.parameters(): p.data = p.data.cuda()

```

Основная функциональность содержится в определении `parameters`:

```
self.params + sum([m.parameters() for m in self.children], [])
```

Это означает, что мы можем запросить у любого `Module` его параметры, и он их вернет, включая параметры всех дочерних модулей (рекурсивно). Но откуда он знает, что это за параметры? В этом помогает реализация особого метода Python `__setattr__`, который вызывается каждый раз, когда Python устанавливает атрибут в классе. В нашей реализации присутствует следующая строка:

```
if isinstance(v, Parameter): self.register_parameters(v)
```

Как видите, здесь мы используем новый класс `Parameter` в качестве «маркера» — в итоге все из этого класса добавляется в `params`.

Метод `__call__` позволяет нам определить, что должно происходить, когда наш объект трактуется как функция. Мы просто вызываем `forward` (которого здесь нет, поэтому его должны будут добавить подклассы). После этого мы вызываем хук, если он определен. Теперь вы видите, что хуки PyTorch не делают ничего особенного — они просто вызывают зарегистрированные хуки.

Помимо этой функциональности, `Module` также предоставляет атрибуты `cuda` и `training`, которые мы вскоре задействуем.

Теперь мы готовы создать наш первый `Module`, являющийся `ConvLayer`:

```

class ConvLayer(Module):
    def __init__(self, ni, nf, stride=1, bias=True, act=True):
        super().__init__()
        self.w = Parameter(torch.zeros(nf, ni, 3, 3))
        self.b = Parameter(torch.zeros(nf)) if bias else None
        self.act, self.stride = act, stride
        init = nn.init.kaiming_normal_ if act else nn.init.xavier_normal_
        init(self.w)

    def forward(self, x):
        x = F.conv2d(x, self.w, self.b, stride=self.stride, padding=1)
        if self.act: x = F.relu(x)
        return x

```

Мы не реализуем `F.conv2d` с нуля, так как вы уже должны были это сделать (используя `unfold`) в вопроснике к главе 17. Вместо этого мы просто создаем

небольшой класс, оберывающий его вместе со смещением и инициализацией весов. Убедимся в том, что он успешно работает с `Module.parameters`:

```
l = ConvLayer(3, 4)
len(l.parameters())
```

2

И доступен для вызова (что приводит к вызову `forward`):

```
x = torch.randn(128, 3, 64, 64)
x = x.unsqueeze(1)
xb = x.clone().detach_()

l = ConvLayer(3, 4)
r = l(xb)
r.shape
```



```
torch.Size([128, 4, 64, 64])
```

Аналогичным образом реализуем `Linear`:

```
class Linear(Module):
    def __init__(self, ni, nf):
        super().__init__()
        self.w = Parameter(torch.zeros(nf, ni))
        self.b = Parameter(torch.zeros(nf))
        nn.init.xavier_normal_(self.w)

    def forward(self, x): return x@self.w.t() + self.b
```

И проверим его работоспособность:

```
l = Linear(4,2)
r = l(torch.ones(3,4))
r.shape
```



```
torch.Size([3, 2])
```

Давайте также создадим модуль тестирования для проверки правильного регистрация добавляемых в качестве атрибутов параметров:

```
class T(Module):
    def __init__(self):
        super().__init__()
        self.c, self.l = ConvLayer(3,4), Linear(4,2)
```

Так как у нас есть сверточный и линейный слои, каждый из которых содержит веса и смещения, то всего мы получаем четыре параметра:

```
t = T()
len(t.parameters())
```

4

Нужно также знать, что вызов `cuda` для этого класса помещает все параметры в GPU:

```
t.cuda()
t.l.w.device
device(type='cuda', index=5)
```

Теперь можно из всего этого создать CNN.

Простая CNN

Как мы видели, класс `Sequential` облегчает реализацию любой архитектуры, так что начнем с него:

```
class Sequential(Module):
    def __init__(self, *layers):
        super().__init__()
        self.layers = layers
        self.register_modules(*layers)

    def forward(self, x):
        for l in self.layers: x = l(x)
        return x
```

Здесь метод `forward` просто вызывает каждый слой по очереди. Обратите внимание, что нам нужно использовать метод `register_modules`, который мы определили в `Module`, поскольку иначе в `parameters` содержимое `layers` не появится.



ВЕСЬ КОД ЗДЕСЬ

Помните, что здесь для модулей мы не используем функциональность PyTorch и определяем все сами. Поэтому если вы не до конца понимаете роль `register_modules`, посмотрите еще раз, что мы написали в коде для `Module`.

Далее мы создадим упрощенный `AdaptivePool`, выполняющий адаптивный пулинг в выход 1×1 , который он сглаживает посредством `mean`:

```
class AdaptivePool(Module):
    def forward(self, x): return x.mean((2,3))
```

Этого достаточно для создания CNN.

```
def simple_cnn():
    return Sequential(
        ConvLayer(3 ,16 ,stride=2), #32
        ConvLayer(16,32 ,stride=2), #16
        ConvLayer(32,64 ,stride=2), # 8
        ConvLayer(64,128,stride=2), # 4
```

```
        AdaptivePool(),
        Linear(128, 10)
    )
```

Посмотрим, правильно ли регистрируются все параметры:

```
m = simple_cnn()
len(m.parameters())
10
```

Теперь можно перейти к добавлению хука. Обратите внимание, что мы оставили в `Module` место только для одного хука. Вы можете сделать его списком или использовать нечто вроде `Pipeline` для выполнения нескольких хуков как одной функции:

```
def print_stats(outp, inp): print (outp.mean().item(), outp.std().item())
for i in range(4): m.layers[i].hook = print_stats

r = m(xbt)
r.shape
0.5239089727401733 0.8776043057441711
0.43470510840415955 0.8347987532615662
0.4357188045978546 0.7621666193008423
0.46562111377716064 0.7416611313819885
torch.Size([128, 10])
```

Данные и модель подготовлены, теперь нужна функция потерь.

Функция потерь

Мы уже видели, как определять «отрицательное логарифмическое правдоподобие»:

```
def nll(input, target): return -input[range(target.shape[0])], target].mean()
```

Вообще-то логарифма здесь нет, так как мы используем то же определение, что и PyTorch. Это значит, что логарифм нужно поместить вместе с softmax:

```
def log_softmax(x): return (x.exp()/(x.exp().sum(-1,keepdim=True))).log()

sm = log_softmax(r); sm[0][0]
tensor(-1.2790, grad_fn=<SelectBackward>)
```

Их совмещение даст нам перекрестную энтропию:

```
loss = nll(sm, yb)
loss
tensor(2.5666, grad_fn=<NegBackward>)
```

Заметьте, что формула:

$$\log(a / b) = \log(a) - \log(b)$$

дает упрощение, когда мы вычисляем `log_softmax`, которая раньше была определена как `(x.exp()/(x.exp().sum(-1))).log()`:

```
def log_softmax(x): return x - x.exp().sum(-1,keepdim=True).log()
sm = log_softmax(r); sm[0][0]

tensor(-1.2790, grad_fn=<SelectBackward>)
```

Далее есть более стабильный способ вычисления логарифма суммы экспонент, называемый LogSumExp (<https://oreil.ly/9UB0b>). Он подразумевает использование следующей формулы:

$$\log\left(\sum_{j=1}^n e^{x_j}\right) = \log\left(e^a \sum_{j=1}^n e^{x_j-a}\right) = a + \log\left(\sum_{j=1}^n e^{x_j-a}\right),$$

где a — это максимум от x_j . Вот эквивалент в коде:

```
x = torch.rand(5)
a = x.max()
x.exp().sum().log() == a + (x-a).exp().sum().log()

tensor(True)
```

Поместим это в функцию:

```
def logsumexp(x):
    m = x.max(-1)[0]
    return m + (x-m[:,None]).exp().sum(-1).log()

logsumexp(r)[0]

tensor(3.9784, grad_fn=<SelectBackward>)
```

И получим возможность использовать это для `log_softmax`:

```
def log_softmax(x): return x - x.logsumexp(-1,keepdim=True)
```

Так мы получим тот же результат, что и ранее:

```
sm = log_softmax(r); sm[0][0]

tensor(-1.2790, grad_fn=<SelectBackward>)
```

С помощью этого можно создать `cross_entropy`:

```
def cross_entropy(preds, yb): return nll(log_softmax(preds), yb).mean()
```

Теперь совместим все эти компоненты в `Learner`.

Learner

У нас есть данные, модели и функция потерь. Теперь для начала подстройки модели недостает только одного — оптимизатора. Вот SGD:

```
class SGD:
    def __init__(self, params, lr, wd=0.):
        store_attr(self, 'params,lr,wd')
    def step(self):
        for p in self.params:
            p.data -= (p.grad.data + p.data*self.wd) * self.lr
            p.grad.data.zero_()
```

Очевидно, что Learner упрощает нам жизнь. Ему нужно знать наши обучающий и контрольный наборы, а значит, нам потребуется `DataLoaders` для их хранения. Другая функциональность нас не интересует, нужно лишь место, где они смогут храниться и куда мы сможем за ними обращаться.

```
class DataLoaders:
    def __init__(self, *dls): self.train, self.valid = dls

dls = DataLoaders(train_dl, valid_dl)
```

Теперь можно создавать сам класс Learner:

```
class Learner:
    def __init__(self, model, dls, loss_func, lr, cbs, opt_func=SGD):
        store_attr(self, 'model,dls,loss_func,lr,cbs,opt_func')
        for cb in cbs: cb.learner = self

    def one_batch(self):
        self('before_batch')
        xb,yb = self.batch
        self.preds = self.model(xb)
        self.loss = self.loss_func(self.preds, yb)
        if self.model.training:
            self.loss.backward()
            self.opt.step()
        self('after_batch')

    def one_epoch(self, train):
        self.model.training = train
        self('before_epoch')
        dl = self.dls.train if train else self.dls.valid
        for self.num, self.batch in enumerate(progress_bar(dl, leave=False)):
            self.one_batch()
        self('after_epoch')

    def fit(self, n_epochs):
        self('before_fit')
        self.opt = self.opt_func(self.model.parameters(),
```

```

    self.lr) self.n_epochs = n_epochs
try:
    for self.epoch in range(n_epochs):
        self.one_epoch(True)
        self.one_epoch(False)
except CancelFitException: pass
self('after_fit')

def __call__(self, name):
    for cb in self.cbs: getattr(cb, name, noop)()

```

Это самый большой класс из всех, какие мы создавали на протяжении книги, но при этом каждый его метод достаточно мал, так что, просматривая все их по очереди, вы должны разобраться, что здесь происходит.

Главным методом, который мы будем вызывать, является `fit`. Он циклически вызывается

```
for self.epoch in range(n_epochs)
```

и каждую эпоху вызывает `self.one_epoch` для каждого `train=True`, а затем для `train=False`. Далее `self.one_epoch` вызывает `self.one_batch` для каждого пакета в `dls.train` или `dls.valid`, в зависимости от ситуации (после обертывания `DataLoader` в `fastprogress.progress_bar`). В завершение `self.one_batch` следует стандартному набору шагов для подстройки одного мини-пакета, что мы неоднократно видели в книге.

До и после каждого шага `Learner` вызывает `self`, который вызывает `__call__` (являющийся стандартной функциональностью Python). `__call__` использует `getattr(cb, name)` в каждом обратном вызове в `self.cbs`, которая представляет собой встроенную функцию Python, возвращающую атрибут (в данном случае метод) с запрошенным именем. Например, `self('before_fit')` будет вызывать `cb.before_fit()` для каждого обратного вызова, где этот метод определен.

Как видите, `Learner` просто использует стандартный цикл обучения, а дополнительно только вызывает в нужное время обратные вызовы. Так что теперь пора определить эти обратные вызовы.

Обратные вызовы

В `Learner.__init__` у нас есть:

```
for cb in cbs: cb.learner = self
```

Иначе говоря, каждому обратному вызову известно, в какой сущности `learner` он используется. Это необходимо, так как в противном случае обратный вызов не сможет получить от него информацию или изменить его компоненты. Поскольку

получение информации от `learner` происходит регулярно, мы упрощаем этот процесс, определяя `Callback` в качестве подкласса `GetAttr`, указывая в нем атрибут по умолчанию `learner`:

```
class Callback(GetAttr): _default='learner'
```

`GetAttr` — это класс `fastai`, автоматически реализующий стандартные методы Python `__getattr__` и `__dir__`, в связи с чем при каждой вашей попытке обратиться к несуществующему атрибуту он передает запрос в объект, который был определен как `_default`.

Например, нам нужно передавать все параметры модели в GPU автоматически при запуске `fit`. Это можно сделать, определив `before_fit` как `self.learner.model.cuda`. Тем не менее, так как атрибутом по умолчанию является `learner` и `SetupLearnerCB` наследует от `Callback` (который наследует от `GetAttr`), то можно удалить `.learner` и просто вызвать `self.model.cuda`:

```
class SetupLearnerCB(Callback):
    def before_batch(self):
        xb,yb = to_device(self.batch)
        self.learner.batch = tfm_x(xb),yb

    def before_fit(self): self.model.cuda()
```

В `SetupLearnerCB` мы тоже передаем каждый мини-пакет в GPU, вызывая `to_device(self.batch)` (можно было использовать более длинный `to_device(self.learner.batch)`). При этом нужно иметь в виду, что `.learner` в строке `self.learner.batch = tfm_x(xb),yb` удалить нельзя, потому что здесь мы *устанавливаем* этот атрибут, а не получаем.

Прежде чем задействовать `Learner`, создадим обратный вызов для отслеживания и вывода его прогресса. Если этого не сделать, мы не будем знать, правильно ли он работает:

```
class TrackResults(Callback):
    def before_epoch(self): self.accs,self.losses,self.ns = [],[],[]

    def after_epoch(self):
        n = sum(self.ns)
        print(self.epoch, self.model.training, sum(self.losses).item()/n,
              sum(self.accs).item()/n)

    def after_batch(self):
        xb,yb = self.batch
        acc = (self.preds.argmax(dim=1)==yb).float().sum()
        self.accs.append(acc)
        n = len(xb)
        self.losses.append(self.loss*n)
        self.ns.append(n)
```

Вот теперь мы готовы к первому использованию Learner!

```
cbs = [SetupLearnerCB(), TrackResults()]
learn = Learner(simple_cnn(), dls, cross_entropy, lr=0.1, cbs=cbs)
learn.fit(1)

0 True 2.1275552130636814 0.2314922378287042

0 False 1.9942575636942674 0.2991082802547771
```

Очень приятно осознавать возможность реализации всех ключевых идей Learner в столь небольшом объеме кода. Теперь же добавим в процесс график скорости обучения.

Планирование скорости обучения

Для получения хороших результатов нам понадобится искатель скорости обучения и метод 1cycle. Они оба представляют обратные вызовы *отжига*, то есть постепенно изменяют гиперпараметры по ходу обучения. Вот LRFinder:

```
class LRFinder(Callback):
    def before_fit(self):
        self.losses, self.lrs = [], []
        self.learner.lr = 1e-6

    def before_batch(self):
        if not self.model.training: return
        self.opt.lr *= 1.2

    def after_batch(self):
        if not self.model.training: return
        if self.opt.lr>10 or torch.isnan(self.loss): raise CancelFitException
        self.losses.append(self.loss.item())
        self.lrs.append(self.opt.lr)
```

Здесь показано, как мы используем `CancelFitException`, который является пустым классом, применяемым для обозначения типа исключения. В Learner видно, что это исключение перехватывается. (Вам стоит самостоятельно добавить и протестировать `CancelBatchException`, `CancelEpochException` и т. д.) Добавим его к списку обратных вызовов и проверим, как он работает:

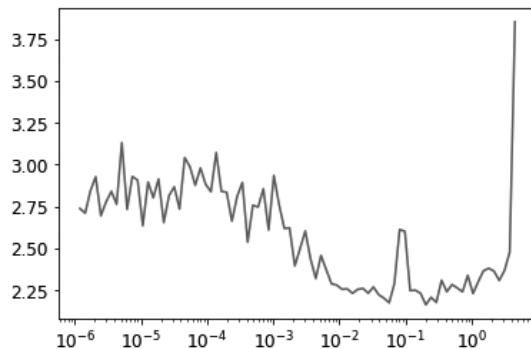
```
lrfind = LRFinder()
learn = Learner(simple_cnn(), dls, cross_entropy, lr=0.1, cbs=cbs+[lrfind])
learn.fit(2)

0 True 2.6336045582954903 0.11014890695955222

0 False 2.230653363853503 0.18318471337579617
```

Взглянем на результаты:

```
plt.plot(lrfind.lrs[:-2], lrfind.losses[:-2])
plt.xscale('log')
```



Теперь можно определить обратный вызов обучения OneCycle:

```
class OneCycle(Callback):
    def __init__(self, base_lr): self.base_lr = base_lr
    def before_fit(self): self.lrs = []

    def before_batch(self):
        if not self.model.training: return
        n = len(self.dls.train)
        bn = self.epoch*n + self.num
        mn = self.n_epochs*n
        pct = bn/mn
        pct_start, div_start = 0.25, 10
        if pct<pct_start:
            pct /= pct_start
            lr = (1-pct)*self.base_lr/div_start + pct*self.base_lr
        else:
            pct = (pct-pct_start)/(1-pct_start)
            lr = (1-pct)*self.base_lr
        self.opt.lr = lr
        self.lrs.append(lr)
```

Попробуем скорость 0,1:

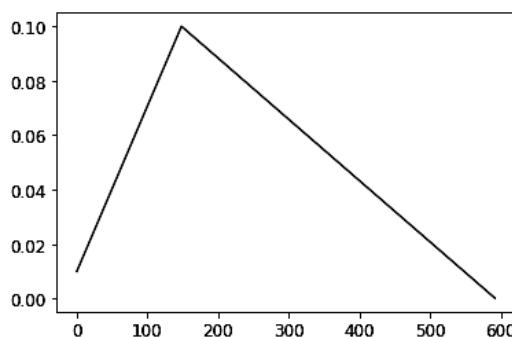
```
onecyc = OneCycle(0.1)
learn = Learner(simple_cnn(), dls, cross_entropy, lr=0.1, cbs=cbs+[onecyc])
```

Обучим еще немного и посмотрим, что получится (мы не будем показывать весь вывод — попробуйте выполнить это в блокноте, чтобы оценить результаты):

```
learn.fit(8)
```

В завершение убедимся, что скорость обучения следовала нашему графику (как видите, косинусный отжиг мы здесь не используем):

```
plt.plot(onecyc.lrs);
```



Резюме

В этой главе мы погрузились в использование ключевых концепций fastai, реализовав их повторно. Поскольку в основном мы здесь работали с кодом, вам стоит самостоятельно поэкспериментировать с ним, сверяясь с соответствующим блокнотом на сайте книги. Теперь, когда вы знаете, из чего и как построена эта библиотека, рекомендуем ознакомиться с руководствами среднего и продвинутого уровней из ее документации, чтобы освоить настройку каждого компонента.

Вопросник



ЭКСПЕРИМЕНТЫ

Для ответа на вопрос, в котором требуется объяснить, что такое функция и класс, вам также нужно выполнить свои практические примеры кода.

1. Что такое `glob`?
2. Как открыть изображение с помощью Python imaging library?
3. Что делает `L.map`?
4. Что делает `Self`?
5. Что такое `L.val2idx`?
6. Какие методы необходимо реализовать для создания `Dataset`?

7. Зачем мы вызываем `convert`, когда открываем изображение из `ImageNet`?
8. Что делает `~`? В чем его польза при разделении обучающей и контрольной выборки?
9. Работает ли `~` с классами `L` или `Tensor`? А с массивами NumPy, списками Python или Pandas DataFrames?
10. Что такое `ProcessPoolExecutor`?
11. Как работает `L.range(self.ds)`?
12. Что такое `__iter__`?
13. Что такое `first`?
14. Что такое `permute` и зачем он нужен?
15. Что такое рекурсивная функция? Как она помогает определять метод `parameters`?
16. Напишите рекурсивную функцию, возвращающую первые 20 элементов последовательности Фибоначчи.
17. Что такое `super`?
18. Почему подклассам `Module` вместо определения `__call__` нужно переопределять `forward`?
19. Почему `init` в `ConvLayer` зависит от `act`?
20. Зачем `Sequential` нужно вызывать `register_modules`?
21. Напишите хук, выводящий форму активаций каждого слоя.
22. Что такое `LogSumExp`?
23. Чем полезна `log_softmax`?
24. Что такое `GetAttr`? В чем его польза для обратных вызовов?
25. Реализуйте повторно один из обратных вызовов этой главы без наследования от `Callback` или `GetAttr`.
26. Что делает `Learner.__call__`?
27. Что такое `getattr`? (Обратите внимание на отличие в регистре `GetAttr`.)
28. Зачем в `fit` блок `try`?
29. Зачем мы проверяем `model.training` в `one_batch`?
30. Что такое `store_attr`?
31. В чем задача `TrackResults.before_epoch`?
32. Что делает `model.cuda`? Как он работает?
33. Зачем нам нужно проверять `model.training` в `LRFinder` и `OneCycle`?
34. Используйте в `OneCycle` косинусный отжиг.

Дополнительные задания

1. Напишите `resnet18` с нуля (при необходимости загляните в главу 14) и обучите ее с помощью `Learner` из этой главы.
2. Реализуйте слой пакетной нормализации с нуля и используйте его в вашей `resnet18`.
3. Напишите обратный вызов `Mixup` для использования в этой главе.
4. Добавьте в `SGD` импульс.
5. Выберите несколько интересных для вас возможностей `fastai` (или другой библиотеки) и реализуйте их с помощью созданных в этой главе объектов.
6. Выберите исследовательскую работу, которая еще не реализовывалась в `fastai` или `PyTorch`, и реализуйте ее с применением объектов, созданных в этой главе. Затем:
 - портируйте работу в `fastai`;
 - отправьте в `fastai` пул-реквест или создайте собственный модуль расширения и сделайте его релиз.

Подсказка: для создания и развертывания пакета может пригодиться `nbdev` (<https://nbdev.fast.ai/>).

ГЛАВА 20

Подведем итог

Поздравляем! Вы справились! Если вы проработали все блокноты вплоть до этого момента, то присоединились к небольшой, но быстро растущей группе людей, способных обуздить мощь глубокого обучения для решения реальных задач, хотя пока что для вас это может быть не столь очевидно. Мы постоянно видим, как студенты, окончившие курсы fast.ai, серьезно недооценивают собственный потенциал как практиков глубокого обучения. Мы также часто видим, что этих людей недооценивают другие специалисты, имеющие за своими плечами классическое академическое образование. Так что если вы готовы превзойти собственные ожидания и ожидания других, то дальнейшие ваши действия по окончании чтения книги будут еще более важными, чем все проделанное до сих пор.

Самое важное — это сохранять импульс. К тому же, как вы знаете из собственного опыта изучения оптимизаторов, импульс — это характеристика, способная к самостоятельному развитию. Так что подумайте о том, как вы можете продолжать и ускорять свое путешествие по миру глубокого обучения. А тем временем рис. 20.1 подбросит вам несколько идей.

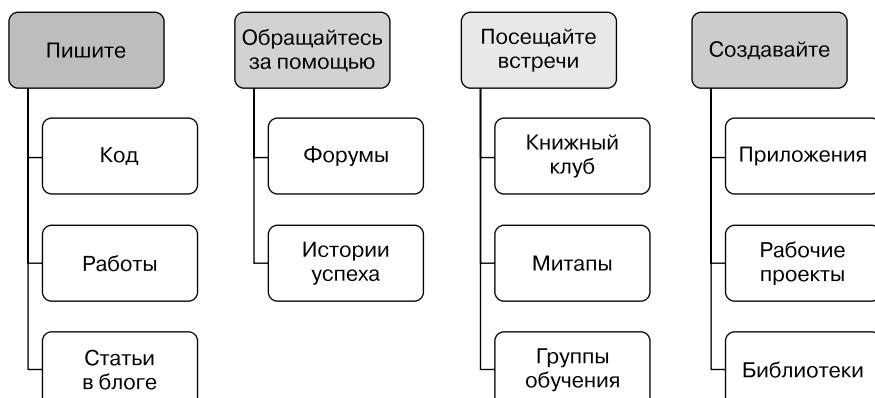


Рис. 20.1. Что делать дальше

В книге мы много говорили о том, насколько важно регулярно практиковать письмо, будь то код или текст. Сейчас вы наверняка написали не столько, сколько хотелось бы, и это нормально. Теперь же как раз есть отличная возможность восполнить этот пробел. Вы уже многое можете сказать. Быть может, вы пробовали экспериментировать с датасетом под таким углом, под которым другие его еще не рассматривали. Расскажите миру об этом! Возможно, вы обдумываете какие-то идеи, пришедшие в голову в процессе чтения книги, — сейчас отличное время перевести все эти замыслы в код.

Если вы хотите поделиться своими мыслями, то одной из скромных площадок для этого может послужить форум fast.ai (<https://forum.fast.ai/>). Там вы встретите очень дружное и готовое помочь сообщество, так что присоединяйтесь и не стесняйтесь рассказывать нам о своих достижениях. Вы также можете обратиться с интересующими вас вопросами к ребятам, немного дальше продвинувшимся на этом пути обучения.

Не стесняйтесь и смело сообщайте на форуме о любом как большом, так и малом своем успехе. Это очень помогает, потому что успех одних студентов может послужить чрезвычайной мотивацией для других.

При этом для многих наиболее эффективным способом поддержания интереса является формирование вокруг него сообщества. Например, вы можете попробовать создать группу учащихся в своем городе или организовать небольшую встречу по теме глубокого обучения и предложить к обсуждению наиболее интересующие вас аспекты. Вполне正常но, что пока вы не стали экспертом мирового уровня, важно помнить, что теперь вы знаете большой объем материала, который не знают другие. Так что ваша точка зрения наверняка вызовет интерес у других.

Еще один очень полезный вид мероприятия — это регулярные книжные встречи или встречи для чтения научных работ. Возможно, в вашем городе такие уже проходят. Если же нет, то вы вполне можете их организовать. Даже если у вас найдется всего один единомышленник, это уже окажется поддержкой и добавит решительности в продвижении.

Если вы живете вдалеке от мест, где можно было бы найти товарищей в этой области интересов, то обращайтесь к форумам, так как люди постоянно организуют виртуальные обучающие встречи. Обычно в них участвует множество ребят, которые еженедельно в формате, например, видеочата обсуждают разные тематики глубокого обучения.

Надеемся, что к этому моменту у вас уже образовалось несколько собранных воедино небольших проектов, а также ряд проделанных экспериментов. В качестве следующего шага мы советуем вам выбрать один из этих проектов и отточить его по максимуму своих возможностей, чтобы вы могли им реально гордиться. Это послужит стимулом к дальнейшему углублению в тему,

попутно проверив уровень ваших знаний и дав возможность оценить свои возможности.

Вас также может заинтересовать бесплатный онлайн-курс [fast.ai](https://course.fast.ai/) (<https://course.fast.ai/>), в котором рассматривается тот же материал, что и в книге. Иногда взгляд на одни и те же темы с разных сторон помогает прояснить их смысл. Исследователи теории человеческого обучения выяснили, что одним из лучших способов усвоения материала является его восприятие с разных сторон и с разным описанием.

Вашей последней миссией, если вы согласитесь ее пройти, будет передать эту книгу другу, чтобы и он встал на путь освоения глубокого обучения!

Приложения

ПРИЛОЖЕНИЕ А

Создание блога

В главе 2 мы предложили вам завести блог, чтобы лучше усвоить материал, который вы изучаете и практикуете. Но что, если у вас еще блога нет? Какую платформу лучше использовать?

К сожалению, когда дело доходит до ведения блога, приходится принимать нелегкое решение: либо задействовать платформу, которая делает это доступным, но нагружает вас и ваших читателей рекламой, платными подписками или ограничениями бесплатных возможностей, либо затратить долгие часы на настройку собственного хостинга и недели на изучение всех связанных с этим тонкостей. Самое, на наш взгляд, важное преимущество самостоятельного сценария в том, что вы становитесь полноправным владельцем собственных публикаций и не зависите от решений провайдера о том, как монетизировать ваш контент в будущем.

Но при всем при этом мы спешим вас обрадовать, что можно найти удачный баланс между двумя этими крайностями.

Блогинг на GitHub Pages

Отличным решением будет разместить блог на бесплатной платформе GitHub Pages (<https://pages.github.com/>), которая не напрягает рекламой, не выдвигает платежных требований и при этом делает вашу информацию доступной, позволяя в любое время без проблем переместить свой блог на другой сервис. Однако все известные нам подходы использования GitHub Pages требовали знания командной строки и сложных инструментов, с которыми знакомы в основном разработчики ПО. Например, в документации (<https://oreil.ly/xemwJ>) GitHub по настройке блога указан длинный список инструкций, включая установку языка Ruby, использование инструмента командной строки `git`, копирование номеров версий и многое другое — в общей сложности 17 шагов!

Чтобы максимально сократить сложный путь, мы создали простой способ, который позволит вам использовать для блогинга *полностью браузерный интерфейс*. Для настройки нового блога вам потребуется буквально около пяти минут. Все это будет для вас бесплатно, и вы также сможете легко добавить в него собственный домен, если захотите. В этом разделе мы объясним, как все это сделать с помощью созданного нами шаблона `fast_template`. (Не забывайте посещать сайт книги для ознакомления с рекомендациями по ведению блога, так как постоянно появляются новые инструменты.)

Создание репозитория

Вам понадобится аккаунт на GitHub. Поэтому отправляйтесь на этот ресурс и создайте учетную запись, если еще этого не сделали. Обычно эта платформа используется разработчиками ПО для написания кода, и они работают с ним через запутанную командную строку. Мы же покажем вам подход, в котором командная строка вам не понадобится вовсе.

Для начала перейдите в браузере по адресу https://github.com/fastai/fast_template/generate (заранее авторизовавшись). Это позволит вам создать *репозиторий* для размещения блога. Вы увидите экран, аналогичный приведенному на рис. А.1.

The screenshot shows the GitHub repository creation interface. At the top, there's a 'Owner' dropdown set to 'jph00' and a 'Repository name' field containing 'jph00.github.io'. Below this, a note says 'Great repository names are short and memorable. Need inspiration?'. A 'Description (optional)' field contains 'Jeremy's blog'. Under 'Visibility', 'Public' is selected (indicated by a checked radio button) with the note 'Anyone can see this repository. You choose who can commit.'. Below that, 'Private' is another option with the note 'You choose who can see and commit to this repository.'. At the bottom is a large blue 'Create repository from template' button.

Рис. А.1. Создание репозитория

Обратите внимание, что название репозитория необходимо ввести в точности в том формате, который указан здесь, то есть ваше имя пользователя GitHub, сопровождаемое `.github.io`.

После ввода названия и любого сопутствующего описания щелкните на *Create repository from template*. Вы можете сделать репозиторий *private* (закрытым), но так как вы создаете блог, который будет доступен для чтения другим людям, то раскрытие внутренних файлов не должно быть для вас проблемой.

А теперь настроим домашнюю страницу!

Настройка домашней страницы

Первым делом посетители вашего блога видят содержимое файла `index.md`. Это файл в формате Markdown (<https://oreil.ly/aVOhs>). Markdown – это мощный и в то же время простой способ создания форматированного текста, содержащего опорные пункты, шрифт с наклонным начертанием, гиперссылки и т. д. Он широко используется, в том числе для форматирования блокнотов Jupyter, практически для всех частей сайта GitHub и на многих других интернет-ресурсах. Для создания текста в формате Markdown можно просто печатать на русском или другом языке и затем добавлять специальные символы, отвечающие за соответствующее форматирование. Например, если вы поставите знак * до и после слова или фразы, то переведете его (ее) в *курсив*. Давайте попробуем.

Чтобы открыть файл, выберите его в GitHub. Для редактирования щелкните на значке карандаша в правой части экрана, как показано на рис. А.2.



Рис. А.2. Переход к редактированию файла

Вы можете добавлять, редактировать и удалять отображаемый текст. Щелкните на *Preview changes* (рис. А.3), чтобы увидеть, как будет выглядеть текст в блоге. Строки, которые вы только добавили или отредактировали, будут обозначены зелеными полосками с левой стороны.

Для сохранения изменений прокрутите страницу вниз и нажмите кнопку *Commit changes* (рис. А.4). *Коммит* здесь означает сохранение на сервере GitHub.

Далее вам нужно настроить блог. Выберите файл `_config.yml` и нажмите кнопку редактирования, как делали для файла `index.md`. Измените значения заголовка, описания и имени пользователя GitHub (рис. А.5). При этом изменять начало строк до двоеточия не нужно, изменяются лишь последующие значения, отделяемые от двоеточий пробелом. При желании можно добавить электронный адрес и имя пользователя в Twitter, но имейте в виду, что эти данные будут отображены в вашем блоге.



Рис. А.3. Предпросмотр изменений для исправления возможных ошибок



Рис. А.4. Коммит изменений для их сохранения

```

1 # Welcome to Jekyll!
2 #
3 # This config file is meant for settings that affect your whole blog.
4 #
5 # If you need help with YAML syntax, here are some quick references for you:
6 # https://learn-the-web.algonquindesign.ca/topics/markdown-yaml-cheat-sheet/#yaml
7 # https://learnxinyuminutes.com/docs/yaml/
8
9 title: Edit _config.yml to set your title!
10 description: This is where the description of your site will go. You should change it by editing the _config.yml file.
11 github_username: jph00

```

Рис. А.5. Заполнение файла конфигурации

Когда закончите, сделайте коммит изменений так же, как делали с предыдущим файлом. После чего нужно подождать около минуты, пока GitHub обработает ваш новый блог. Перейдите в браузере по адресу <имя_пользователя>.github.

io (заменив <имя_пользователя> своим именем на GitHub). Теперь вы должны увидеть созданный блог, который похож на представленный на рис. А.6.

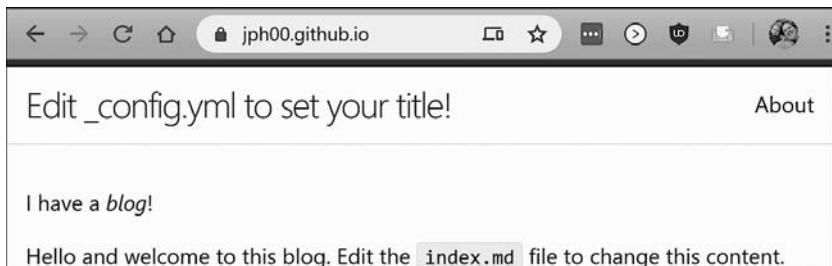


Рис. А.6. Ваш блог запущен!

Создание публикаций

Теперь можете создавать свою первую публикацию. Все публикации будут помещаться в каталог `_posts`. Перейдите в него сейчас, а затем нажмите кнопку `Create file`. Будьте внимательны при именовании файла, используя формат <год>-<месяц>-<день>-<имя>.md, как показано на рис. А.7, где <год> — это четырехзначное число, <месяц> и <день> — двузначные числа, а <имя> может быть чем угодно, что поможет вам запомнить тему публикации. Расширение .md означает документ Markdown.



Рис. А.7. Именование публикаций

Далее можете написать само содержимое публикации. Единственное правило при этом — это то, что первая строка должна быть заголовком Markdown. Для ее создания в начале строки помещается #, как показано на рис. А.8 (так создается заголовок первого уровня, который нужно использовать только один раз в начале документа; заголовки второго уровня создаются вводом ##, третьего — с помощью ### и т. д.).

Как и ранее, вы можете нажать кнопку `Preview` для просмотра итогового форматирования (рис. А.9).

Для дальнейшего сохранения этого нужно будет нажать `Commit new file`, как показано на рис. А.10.

Взгляните еще раз на домашнюю страницу блога, и вы увидите, что на ней появилась созданная публикация. На рис. А.11 показан рассматриваемый нами

пример. Помните, что нужно будет подождать около минуты, пока GitHub обработает запрос для отображения файла.

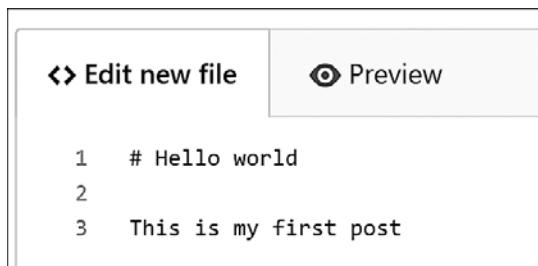


Рис. А.8. Синтаксис Markdown для заголовка



Рис. А.9. Так будет выглядеть пример приведенного выше синтаксиса в вашем блоге

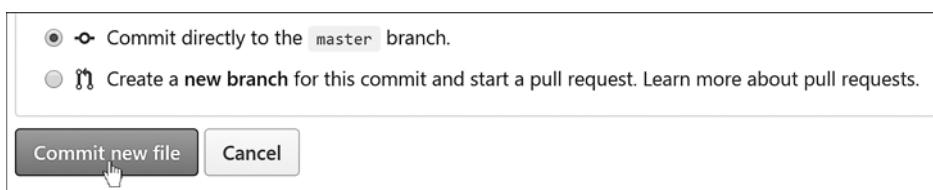


Рис. А.10. Коммит изменений для их сохранения

Возможно, вы заметили, что мы привели пример публикации, который вы можете прямо сейчас удалить. Перейдите в каталог `_posts` и выберите файл `2020-01-14-welcome.md`. Далее щелкните на значке корзины справа (рис. А.12).

В GitHub ничего не изменится, пока вы не выполните коммит, включая удаление файла. Так что после щелчка на значке корзины прокрутите страницу вниз и сделайте коммит изменений.

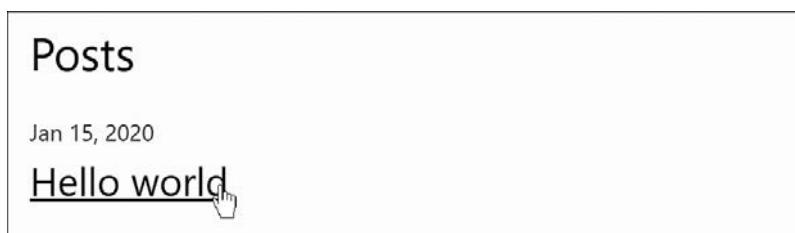


Рис. А.11. Ваша первая публикация готова!

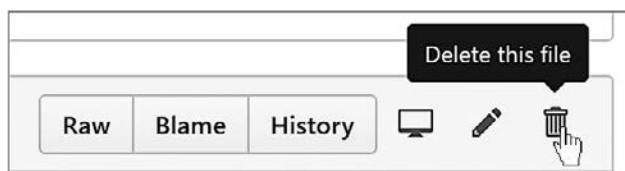


Рис. А.12. Удаление примера публикации

В публикации также можно встраивать изображения путем добавления строки markdown-разметки, аналогичной следующей:

```
![Image description](images/filename.jpg)
```

Чтобы она заработала, нужно будет поместить соответствующее изображение в каталог `images`. Для этого выберите каталог `images` и нажмите кнопку `Upload files` (рис. А.13).

Теперь посмотрим, как все это делать напрямую с компьютера.

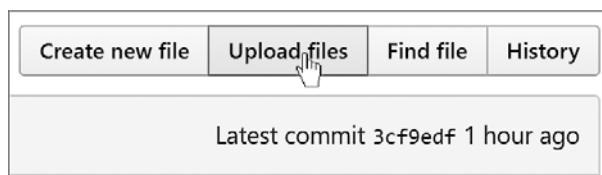


Рис. А.13. Загрузка файла с компьютера

Синхронизация GitHub и компьютера

Есть множество причин, по которым вы можете захотеть скопировать содержимое блога на ПК, — возможно, вам понадобится возможность читать или редактировать публикации офлайн или вы решите сделать резервную копию на случай возникновения проблем с репозиторием.

GitHub не просто позволяет вам скопировать репозиторий на ПК, он дает возможность *синхронизировать* его с ним. Это означает, что вы можете вносить изменения на GitHub и они будут отображаться у вас на компьютере. При этом верно и обратное — при редактировании контента на компьютере изменения будут вноситься в репозиторий на GitHub. Вы даже можете дать другим людям доступ и возможность редактировать ваш блог, и при следующей синхронизации их и ваши изменения будут совмещены.

Чтобы это все заработало, потребуется установить на ПК приложение GitHub Desktop. Оно работает на всех современных платформах: Mac, Windows и Linux. Для его установки следуйте инструкциям, а после запуска вам будет предложено авторизоваться на GitHub и выбрать репозиторий для синхронизации. Выберите **Clone a repository from the Internet**, как показано на рис. А.14.

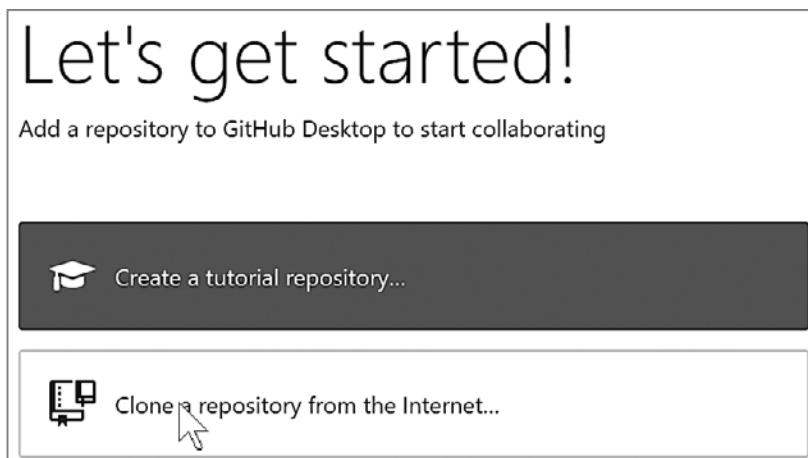


Рис. А.14. Клонирование репозитория в GitHub Desktop

Как только GitHub завершит синхронизацию, вы сможете нажать **View the files of your repository in Explorer** (или **Finder**), как показано на рис. А.15, и вы увидите копию своего блога! Попробуйте отредактировать один из файлов на ПК, а затем вернитесь в GitHub Desktop, где уже будет кнопка **Sync**. Нажав ее, вы скопируете все изменения на GitHub.

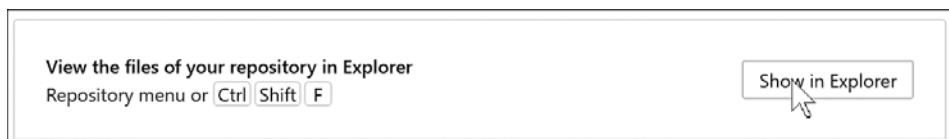


Рис. А.15. Локальный просмотр файлов

Если вы ранее не использовали `git`, то GitHub Desktop — отличный способ начать. Как вы позже узнаете, это основной инструмент, который используется большинством аналитиков данных. Еще один инструмент, который, как мы надеемся, вам полюбился, — Jupyter Notebook, и у вас есть возможность делать публикации непосредственно из него!

Блогинг из Jupyter

Блокноты Jupyter позволяют размещать их содержимое в блоге. Ячейки Markdown, ячейки кода и все выводы будут отображаться в экспортируемой публикации. Наиболее удобный способ реализации этого мог измениться со временем прочтения вами книги, так что рекомендуем обратиться к сайту (<https://book.fast.ai/>) за последней информацией. В момент написания самый простой способ создания блога из блокнота — это использование fastpages (<http://fastpages.fast.ai/>), продвинутой версии `fast_template`.

Для отправки материала в блог из блокнота просто поместите его в каталог `_notebooks` в репозитории, и он появится в списке публикаций. При написании блокнота вы можете без проблем добавлять все, что хотите продемонстрировать читателям. Поскольку на большинстве платформ для блогинга в публикации сложно вставлять код и его результаты, многие из нас выработали привычку включать меньше реальных примеров, чем стоило бы. Так что использование блокнота окажется отличным способом выработать привычку включать в текст больше примеров.

Зачастую вам понадобится скрывать шаблонный код, такой как инструкции импорта. Для этого нужно добавить `#hide` в верхнюю часть ячейки, чтобы скрыть ее показ при выводе. Jupyter отображает результат последней строки ячейки, так что нет необходимости добавлять инструкцию `print`. (Добавление необязательного кода окажет излишнюю мыслительную нагрузку на читателя, так что постарайтесь этого избежать.)

ПРИЛОЖЕНИЕ Б

Схема подготовки проекта по аналитике данных

Для создания полезных проектов по аналитике данных требуется не только точная модель. Когда Джереми занимался консультированием по этим вопросам, он всегда стремился сначала понять конкретные нужды и контекст организации, в чем ему помогала следующая схема вопросов (рис. Б.1).

Стратегия

Чего организация хочет добиться (ее *цель*) и какие действия она может предпринять для повышения эффективности (*рычаги*)?

Данные

Собирает ли организация необходимые данные и делает ли их доступными?

Аналитика

Выявление каких факторов окажется для организации полезным?

Реализация

Какими возможностями организация располагает?

Обслуживание

Какие системы для отслеживания изменений в рабочей среде используются?

Ограничения

Какие ограничения необходимо учитывать в каждом из предыдущих пунктов?



Рис. Б.1. Цепочка ценностей в процессе аналитики

Он разработал анкету с вопросами, которую клиенты заполняли перед началом проекта. Затем по ходу реализации этого проекта Джереми помогал им уточнять ответы. В основе этого вопросника лежит многолетний опыт работы во многих индустриях, включая агрокультуру, добычу полезных ископаемых, банковское дело, пивоварение, телекоммуникации, рынок недвижимости и др.

Прежде чем мы рассмотрим цепочку ценностей аналитики, нужно обратить внимание на первую часть вопросника, которая связана с наймом наиболее важных работников вашего проекта — ученых по данным.

Специалисты по данным

Аналитики данных должны иметь прямые возможности к занятию в компании руководящих ролей. Кроме того, в организации должна быть предусмотрена программа найма таких специалистов непосредственно на должности руководителей. В компании, основывающейся на обработке данных, соответствующие сотрудники должны иметь самый высокий уровень оплаты труда. При этом необходимо организовать систему, которая позволит этим специалистам беспрепятственно сотрудничать и обмениваться опытом.

- Специалисты какого уровня навыков обработки данных есть сейчас в компании?
- Как происходит наем аналитиков данных?
- Как в компании выявляются люди с навыками обработки данных?
- Какие навыки организация ищет и как их оценивает? На основании чего определяется важность именно этих навыков?
- Какая система консалтинга по аналитике данных используется? В каких ситуациях обработка данных поручается сторонним подрядчикам и как передается?

- Сколько получают специалисты по данным? Перед кем они отчитываются? Как поддерживается актуальность их навыков?
- Какие карьерные перспективы есть у таких специалистов?
- Как много людей из руководящего звена имеют уверенные навыки в области аналитики данных?
- Как отбирается и распределяется работа для аналитиков?
- Какое в их распоряжении есть программное и аппаратное обеспечение?

Стратегия

Все проекты по обработке данных должны служить решению стратегически важных задач, а значит, первым в списке должно стоять понимание бизнес-стратегии.

- Каковы пять текущих важнейших стратегических задач организации?
- Какие данные доступны для их решения?
- Используется ли для решения аналитика данных? Привлечены ли специалисты по данным?
- На какие определяющие прибыль факторы организация может серьезно повлиять (рис. Б.2)?
- Какие весомые решения и действия в отношении каждого из таких выявленных факторов может предпринять организация, включая практические действия (например, звонок клиенту) и стратегические (например, выпуск нового продукта)?
- Какие данные доступны для каждого из наиболее важных решений и действий (в самой организации, от вендора или путем сбора в будущем), способные повысить отдачу?
- Какие наиболее результативные возможности обработки данных на основе предыдущего анализа можно выделить в организации?
- Для каждой из этих возможностей:
 - На какой фактор прибыли она повлияет?
 - Какие конкретные решения или действия она задействует?
 - Как эти решения и действия будут связаны с результатами проекта?
 - Какова оценочная рентабельность проекта?
 - Какие временные ограничения и дедлайны могут на нее повлиять?

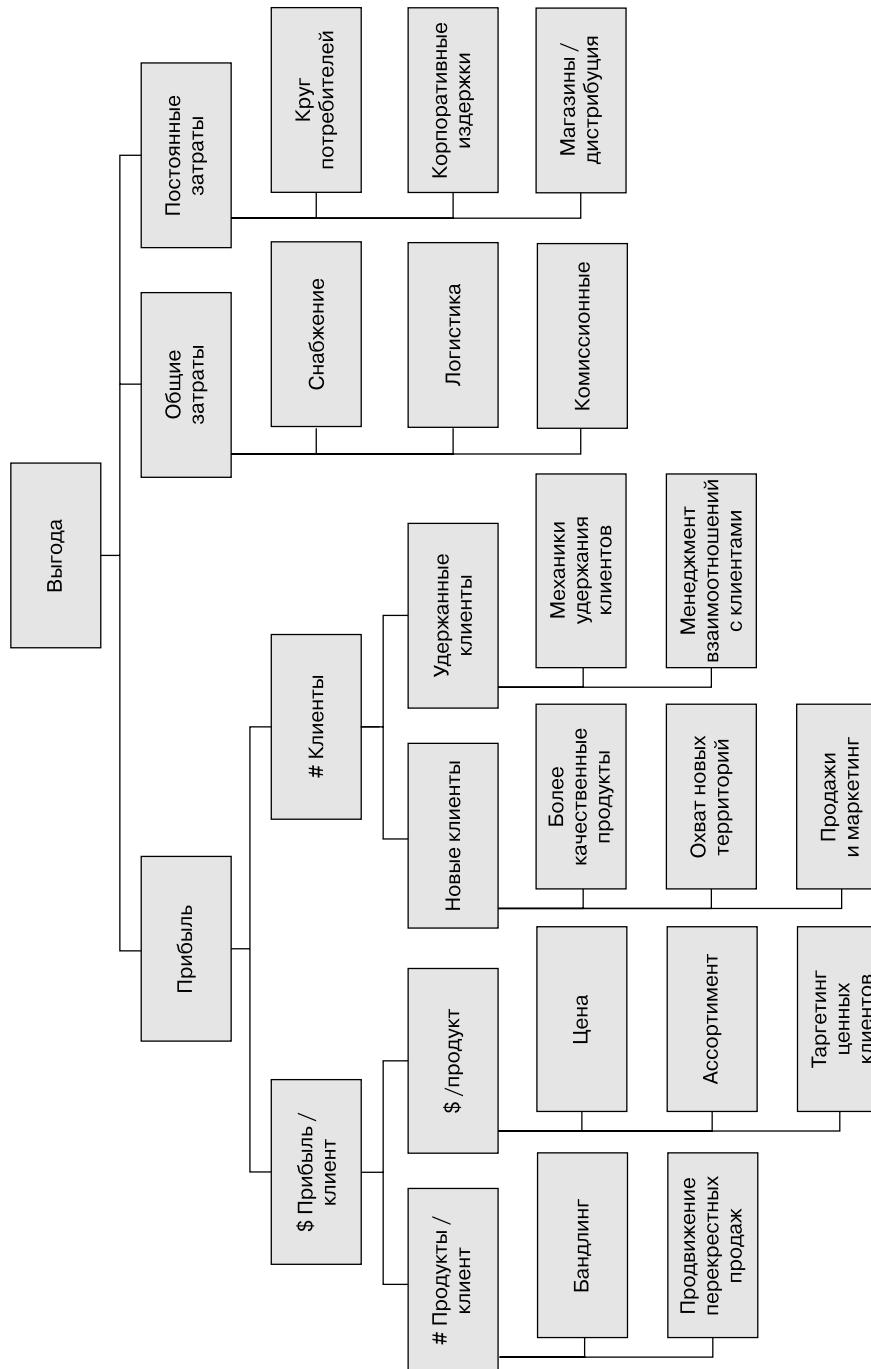


Рис. Б.2. Наиболее актуальные факторы, которые могут определять прибыль организации

Данные

Без данных модель обучить не получится. При этом данные должны быть доступными, интегрированными и проверяемыми.

- Какими платформами данных располагает организация? К ним могут относиться витрины данных, OLAP-кубы, хранилища данных, кластеры Hadoop, системы OLTP, сводные ведомости департаментов и т. д.
- Предоставьте всю собранную информацию по оценке доступности данных в организации, текущей работе и дальнейшим планам по построению платформ данных.
- Какие инструменты и процессы доступны для перемещения данных между системами и форматами?
- Каким образом организован доступ к источникам данных для разных групп пользователей и администраторов?
- Какие инструменты доступа к данным (например, базы данных, клиенты OLAP, внутрифирменное ПО, SAS) имеются в организации? Сколько людей используют каждый инструмент и какие они занимают должности?
- Как пользователей информируют о новых системах, коррективах, новых и измененных элементах данных и т. д.? Приведите примеры.
- На основании чего принимаются решения об ограничении доступа к данным? Как и кем обрабатываются запросы на доступ к защищенным данным? Какие в этом задействованы критерии? Сколько в среднем времени занимает ответ? Какой процент запросов одобряется? Как это все отслеживается?
- Как организация решает, когда собирать дополнительные данные или приобретать внешние? Приведите примеры.
- Какие данные использовались при анализе последних проектов? Что оказалось наиболее полезным? Что было менее значимым? На основе чего происходила оценка?
- Исходя из каких дополнительных внутренних данных можно принять полезные решения для предлагаемых проектов? Как насчет внешних данных?
- Какие возможны ограничения или сложности при доступе к этим данным или их получении?
- Какие произошедшие за последние два года изменения в сборе данных, написании кода, интеграции и т. д. могли повлиять на интерпретацию или доступность собранных данных?

Аналитика

Специалисты по данным должны иметь доступ к актуальным, соответствующим их нуждам инструментам. При этом необходимо постоянно отслеживать появление новых, более производительных и эффективных средств решения задач.

- Какие инструменты аналитики используются в организации и кем? Как происходит их выбор, настройка и обслуживание?
- Как организован процесс настройки дополнительных инструментов аналитики на клиентской машине? Сколько в среднем это занимает времени? Какой процент таких запросов удовлетворяется?
- Как выстраивают системы аналитики сторонние консультанты, приглашенные в организацию? Просят ли их ограничивать используемые системы, чтобы результаты соответствовали внутренней инфраструктуре?
- В каких ситуациях использовалась облачная обработка? Каковы планы по использованию облака?
- В каких ситуациях для специализированной аналитики привлекались сторонние эксперты? Как этот процесс организовывался? По какому принципу происходил выбор этих специалистов?
- Какие аналитические инструменты были опробованы в последних проектах?
- Какие сработали, а какие — нет и почему?
- Предоставьте любые доступные результаты произведенного на сегодняшний день анализа по этим проектам.
- На основании чего оценивались результаты этого анализа? Какие использовались метрики? С какими бенчмарками производились сравнения? Как вы определяете, что модель «достаточно хороша»?
- В каких ситуациях организация использует визуализацию, а не табличную отчетность, предиктивное моделирование (и аналогичные инструменты ML)? Как настраиваются и тестируются модели для более продвинутых подходов моделирования? Приведите примеры.

Реализация

Слабым местом проектов по аналитике данных зачастую оказываются ограничения ИТ-средств, поэтому их следует учитывать наперед.

- Приведите примеры прошлых успешных и безуспешных проектов данных, а также поясните сложности, возникшие в связи с интеграцией ИТ и человеческими ресурсами. Как эти сложности были решены?

- Как подтверждается эффективность аналитических моделей до их реализации?
- Как определяются требования к эффективности реализации проекта по аналитике (в плане скорости и точности)?
- Касательно предложенных проектов предоставьте следующую информацию.
 - Какие ИТ-системы будут использоваться для поддержки основанных на данных решений и действий?
 - Как будет происходить эта интеграция ИТ?
 - Какие альтернативы, требующие меньше интеграции, доступны ИТ?
 - На какие должности повлияют подходы, опирающиеся на данные?
 - Как будет организовано обучение этих сотрудников, а также их контроль и поддержка?
 - Какие предполагаются сложности в реализации?
 - Участие каких заинтересованных лиц потребуется для обеспечения успеха реализации? Как они могут отнестись к этим проектам и потенциальному влиянию проектов на них?

Обслуживание

Если не отслеживать модели должным образом, они могут привести вас к катастрофе.

- Как обслуживаются системы, созданные сторонними подрядчиками? Когда они передаются внутренним командам?
- Как отслеживается эффективность моделей? Когда принимается решение об их пересборке?
- Как происходит внутреннее согласование изменения данных и управление ими?
- Как аналитики данных взаимодействуют с инженерами ПО для обеспечения корректной реализации алгоритмов?
- Когда разрабатываются тестовые случаи и как они администрируются?
- Когда код подвергается рефакторингу? Как в процессе рефакторинга проверяется и поддерживается верность и эффективность моделей?
- Как логируются требования к обслуживанию и поддержке? Как эти логи используются?

Ограничения

Для каждого рассматриваемого проекта необходимо перечислить потенциальные ограничения, способные повлиять на его успех.

- Потребуется ли для использования результатов проекта изменять или разрабатывать ИТ-системы? Есть ли более простые реализации, не подразумевающие существенных изменений в ИТ? Если да, то как использование упрощенной реализации привело бы к существенному снижению их влияния?
- Какие нормативные ограничения присутствуют для сбора данных, анализа или реализации? Изучались ли соответствующие законы и прецеденты? Какие есть обходные пути?
- Какие существуют организационные ограничения, включая культуру, навыки и структуру?
- Есть ли управленческие ограничения и какие?
- Выполнялись ли ранее проекты по аналитике, способные повлиять на отношение организации к подходам, основанным на данных?

Об авторах

Джереми Ховард — предприниматель, бизнес-стратег, разработчик и преподаватель. Основал исследовательский институт fast.ai, цель которого — сделать глубокое обучение максимально доступным каждому. Джереми — выдающийся ученый-исследователь Университета Сан-Франциско, член факультета в Университете сингулярности, обладатель звания молодого мирового лидера, присужденного на международном экономическом форуме.

Последняя организованная им компания Enlitic первой применила глубокое обучение в медицине и, согласно MIT Tech Review, была зачислена в топ-50 интеллектуальных компаний мира в 2015 и 2016 годах. Ранее Джереми занимал пост президента и главного научного консультанта в Kaggle, где на ведущих ролях принимал участие в международных соревнованиях по машинному обучению два года подряд. Стал основоположником двух успешных австралийских стартапов (FastMail и Optimal Decisions Group, которые были приобретены LexisNexis). До этого восемь лет занимался управленческим консультированием в компаниях McKinsey & Co и AT Kearney. Джереми инвестировал во многие стартапы, оказывал им содействие в качестве наставника и просто давал советы. Помимо этого, принимал участие во многих опенсорсных проектах.

Он частый гость известной австралийской утренней программы новостей, выступал на TED.com и выпустил учебники по науке о данных и веб-разработке.

Сильвейн Гутгер — инженер-исследователь в HuggingFace. Ранее в роли ученого-исследователя fast.ai работал над расширением доступности глубокого обучения путем разработки и совершенствования техник, позволяющих ускоренно обучать модели в условиях ограниченных ресурсов.

До этого на протяжении семи лет Сильвейн преподавал computer science и математику по программе CPGE во Франции. CPGE — это особый вид занятий, которые посещают отдельные студенты по завершении высшей школы для подготовки к вступительным экзаменам в ведущие инженерные и бизнес-институты. Помимо этого, Сильвейн написал несколько книг, посвященных всему преподаваемому им курсу обучения, которые были изданы Editions Dunod.

Окончил Высшую школу Нормаль (Париж, Франция), где изучал математику. Более того, Сильвейн получил в этой области степень магистра от IX Парижского университета (Орсе, Франция).

Благодарности

Хотим выделить Алексиса Галлахера и Рейчел Томас, проделавших потрясающую работу. Алексис выступил в роли гораздо большей, нежели просто научный редактор. Его влияние ощущалось в каждой главе. Именно он написал многие наиболее проницательные и убедительные пояснения. Он также предоставил отличное описание структуры библиотеки fastai, особенно API блока данных. Рейчел, в свою очередь, предоставила большую часть материалов для главы 3, а также рассказала о сложностях, присущих области этики данных.

Благодарим сообщество fast.ai, включая сторонних членов forums.fast.ai, пятьсотен участников, приложивших свои усилия к разработке библиотеки fastai, и сотни тысяч студентов course.fast.ai. Особая благодарность наиболее самоотверженным соавторам fastai, включая Захари Мюллера, Радека Осмультски, Эндрю Шоу, Стаса Бекмана, Лукаса Вассеса и Бориса Дайма (Zachary Muller, Radek Osmulski, Andrew Shaw, Stas Bekman, Lucas Vasquez, Boris Dayma). Также спасибо исследователям, использовавшим fastai для своих передовых научных работ: Себастьяну Рудеру, Петру Чапла, Марчину Кардасу, Джилиану Эйзеншлюсу, Нильсу Стродтоффу, Патрику Вагнеру, Маркусу Венцелю, Войцеху Семеку, Полу Марагакису, Хантеру Нисоноффу, Брайану Коулу и Дэвиду Э. Шоу (Sebastian Ruder, Piotr Czapla, Marcin Kardas, Julian Eisenschlos, Nils Strodthoff, Patrick Wagner, Markus Wenzel, Wojciech Samek, Paul Maragakis, Hunter Nisonoff, Brian Cole, and David E. Shaw). Благодарим Хамеля Хуссейна (Hamel Hussain), который создал с помощью fastai одни из наиболее вдохновляющих проектов и выступил основным идеальным организатором платформы для блогинга fastpages. Огромная благодарность Крису Латтнеру (Chris Lattner) за участие в наших беседах, в которые он привнес замечательные идеи из Swift и свои огромные знания в языках программирования. Все это мы смогли успешно использовать для совершенствования структуры fastai.

Спасибо всем сотрудникам O'Reilly за их старания. Особенная благодарность Ребекке Новак (Rebecca Novak), которая обеспечила бесплатную доступность всех блокнотов и полноцветную публикацию. Спасибо также Рейчел Хэд (Rachel Head), чьи комментарии позволили улучшить все части книги. Благодарим Мелиссу Поттер (Melissa Potter), которая помогла успешно завершить весь процесс.

Спасибо всем рецензентам — невероятным людям, которые дали проницательную и осмысленную обратную связь: Орельену Жерону (Aurélien Géron), автору одной из лучших, на наш взгляд, книги по машинному обучению, который помог

доработать и наш труд; Джо Списаку (Joe Spisak), продакт-менеджеру PyTorch; Мигелю Де Иказа, легендарному создателю Gnome, Xamarian и др.; Россу Вайтману (Ross Wightman), разработчику нашего любимого зоопарка моделей в PyTorch; Радеку Осмульски (Radek Osmulski), одному из самых выдающихся выпускников fast.ai, с каким мы имели честь познакомиться; Дмитро Мишкину (Dmytro Mishkin), сооснователю проекта Kornia и автору некоторых из наших любимых работ по глубокому обучению; Фреду Монро (Fred Monroe), который помог нам со множеством проектов; а также Эндрю Шоу (Andrew Shaw), директору WAMRI и создателю прекрасного ресурса musicautobot.com.

Особая благодарность Сумиту Чинтала (Soumith Chintala) и Адаму Пашке (Adam Paszke) за создание PyTorch, а также всей команде PyTorch за то, что сделали его настолько удобным в использовании. И конечно же, огромное спасибо нашим семьям за их поддержку и терпение в процессе реализации столь большого проекта.

06 обложке

На обложке книги — рыба-кабан (*Carpos aper*), единственный известный представитель своего рода. Встречается в основном в тропических и умеренных водах южного полушария. Рыба-кабан обитает на глубине от 40 до 700 м в зоне пелагиали: области открытого моря, которая не находится в непосредственной близости от дна и является наиболее обитаемой частью водной среды на Земле.

Рыба невелика, имеет красно-оранжевый окрас, большие глаза и вытянутый рот. Ее тело продолговатое, сильно сжато с боков и имеет форму правильного ромба. Средний размер особи около 13 см, хотя самки бывают больше. Рекордная зафиксированная длина составляет 27 см. Несмотря на уязвимость для хищников ввиду небольшого размера, эти рыбы перемещаются косяками, что позволяет им эффективно обороняться, а также облегчает процесс поиска еды и спаривания. Их близкий родственник — антигония, или короткоспинная рыба-кабан (*Antigonia combatia*), родом из тропических и субтропических вод, а также глубоководная рыба-кабан (*Antigonia capros*), обитающая в соседних водах Атлантики.

Несмотря на то что нынешний статус рыбы-кабана определен как вызывающий наименьшее беспокойство, многие из животных, публикуемых на обложках книг O'Reilly, находятся под угрозой уничтожения. Все они являются чрезвычайно важными для нашей планеты.

Иллюстрация обложки выполнена Карен Монтгомери на основе черно-белой гравюры из книги «Естественная история».

Джереми Ховард, Сильвейн Гуггер

**Глубокое обучение с fastai и PyTorch:
минимум формул, минимум кода,
максимум эффективности**

Перевел с английского *Д. Брайт*

Руководитель дивизиона
Ведущий редактор
Литературные редакторы
Художественный редактор
Корректоры
Верстка

*Ю. Сергиенко
Н. Гринчик
Ю. Зорина, К. Тульцева
В. Мостипан
С. Беляева, Н. Викторова
Л. Егорова*

Изготовлено в России. Изготовитель: ООО «Прогресс книга». Место нахождения и фактический адрес: 194044, Россия, г. Санкт-Петербург, Б. Сампсониевский пр., д. 29А, пом. 52. Тел.: +78127037373.

Дата изготовления: 05.2022. Наименование: книжная продукция. Срок годности: не ограничен.
Налоговая льгота — общероссийский классификатор продукции ОК 034-2014, 58.11.12 — Книги печатные
профессиональные, технические и научные.

Импортер в Беларусь: ООО «ПИТЕР М», 220020, РБ, г. Минск, ул. Тимирязева, д. 121/3, к. 214,
тел./факс: 208 80 01.
Подписано в печать 30.03.22. Формат 70×100/16. Бумага офсетная. Усл. п. л. 50,310. Тираж 1000. Заказ 0000.
Отпечатано в полном соответствии с качеством предоставленных материалов в ООО «Фотоэксперт».
109316, г. Москва, Волгоградский проспект, д. 42, корп. 5, эт. 1, пом. I, ком. 6.3-23Н

*Сергей Николенко, Артур Кадурин,
Екатерина Архангельская*

ГЛУБОКОЕ ОБУЧЕНИЕ



Перед вами — первая книга о глубоком обучении, написанная на русском языке. Глубокие модели оказались ключом, который подходит ко всем замкам сразу: новые архитектуры и алгоритмы обучения, а также увеличившиеся вычислительные мощности и появившиеся огромные наборы данных привели к революционным прорывам в компьютерном зрении, распознавании речи, обработке естественного языка и многих других типично «человеческих» задачах машинного обучения. Эти захватывающие идеи, вся история и основные компоненты революции глубокого обучения, а также самые современные достижения этой области доступно и интересно изложены в книге. Максимум объяснений, минимум кода, серьезный материал о машинном обучении и увлекательное изложение — в этой уникальной работе замечательных российских ученых и интеллектуалов.

КУПИТЬ