

*Теоретически между теорией и практикой
нет никакой разницы. Но на практике она есть.*
– Бенджамин Брюстер

*Идеальный план проекта –
сначала составить список всех неизвестных.*
– Билл Лэнгли

*Привлечение средств – это искусственный интеллект.
Найм на работу – машинное обучение.
Реализация – это линейная регрессия.
А отладка – это printf().*
– Бэррон Шварц

Содержание

От издательства	14
Вступительное слово	15
Предисловие	17
Глава 1. Введение	19
1.1. Обозначения и определения	19
1.1.1. Структуры данных	19
1.1.2. Заглавная сигма	21
1.2. Что такое машинное обучение.....	21
1.2.1. Обучение с учителем.....	22
1.2.2. Обучение без учителя	23
1.2.3. Обучение с частичным привлечением учителя.....	24
1.2.4. Обучение с подкреплением	24
1.3. Терминология машинного обучения	25
1.3.1. Данные, используемые прямо и косвенно.....	25
1.3.2. Первичные и аккуратные данные	26
1.3.3. Обучающие и зарезервированные наборы.....	27
1.3.4. Ориентир.....	28
1.3.5. Конвейер машинного обучения	28
1.3.6. Параметры и гиперпараметры	29
1.3.7. Классификация и регрессия	29
1.3.8. Обучение на основе модели и обучение на основе экземпляров	30
1.3.9. Поверхностное и глубокое обучение	30
1.3.10. Обучение и оценивание.....	30
1.4. Когда следует использовать машинное обучение	31
1.4.1. Когда задача слишком сложна для кодирования.....	31
1.4.2. Когда задача постоянно меняется.....	32
1.4.3. Когда речь идет о задаче восприятия	32
1.4.4. Когда это неизученное явление.....	32
1.4.5. Когда задача имеет простую целевую функцию	33
1.4.6. Когда это экономически выгодно	33
1.5. Когда не следует использовать машинное обучение.....	34
1.6. Что такое инженерия машинного обучения	34
1.7. Жизненный цикл проекта машинного обучения.....	36
1.8. Резюме	37
Глава 2. Прежде чем приступать к проекту	39
2.1. Определение приоритетов проекта машинного обучения	39
2.1.1. Последствия машинного обучения	39
2.1.2. Стоимость машинного обучения	40

2.2. Оценивание сложности проекта машинного обучения.....	41
2.2.1. Неизвестные	41
2.2.2. Упрощение задачи.....	42
2.2.3. Нелинейный прогресс.....	43
2.3. Определение цели проекта машинного обучения.....	43
2.3.1. Что модель может делать.....	43
2.3.2. Свойства успешной модели	44
2.4. Организация группы машинного обучения	45
2.4.1. Две традиции.....	45
2.4.2. Члены группы машинного обучения.....	46
2.5. Причины провалов проектов машинного обучения	47
2.5.1. Нехватка квалифицированных кадров	47
2.5.2. Отсутствие поддержки со стороны руководства.....	48
2.5.3. Отсутствующая инфраструктура данных.....	48
2.5.4. Трудности с разметкой данных	49
2.5.5. Разобщенные организации и отсутствие сотрудничества.....	49
2.5.6. Технически невыполнимые проекты	50
2.5.7. Нестыковка между техническими и коммерческими группами.....	50
2.6. Резюме	51

Глава 3. Сбор и подготовка данных.....

3.1. Вопросы к данным	53
3.1.1. Доступны ли данные?.....	54
3.1.2. Насколько велик объем данных?.....	54
3.1.3. Пригодны ли данные для использования?	56
3.1.4. Понятны ли данные?	58
3.1.5. Надежны ли данные?.....	58
3.2. Типичные проблемы с данными	60
3.2.1. Высокая стоимость	60
3.2.2. Плохое качество	62
3.2.3. Зашумленность	62
3.2.4. Смещение.....	63
Типы смещения	63
Как избежать смещения	67
3.2.5. Низкая предсказательная способность	69
3.2.6. Устаревшие примеры	70
3.2.7. Выбросы.....	70
3.2.8. Просачивание данных.....	71
3.3. Что считать хорошими данными.....	72
3.3.1. Хорошие данные информативны.....	72
3.3.2. Хорошие данные обладают хорошим покрытием.....	72
3.3.3. Хорошие данные отражают реальные входы	73
3.3.4. Хорошие данные несмещенные	73
3.3.5. Хорошие данные не являются результатом петли обратной связи.....	73
3.3.6. У хороших данных согласованные метки	74
3.3.7. Хорошие данные достаточно велики	74
3.3.8. Сводный перечень свойств хороших данных.....	74

3.4. Обработка данных о взаимодействии	75
3.5. Причины просачивания данных.....	75
3.5.1. Цель является функцией от признака	76
3.5.2. Признак скрывает цель.....	76
3.5.3. Признак из будущего.....	77
3.6. Разбиение данных.....	78
3.6.1. Просачивание во время разбиения.....	79
3.7. Обработка отсутствия атрибутов	80
3.7.1. Методы подстановки данных.....	80
3.7.2. Просачивание во время подстановки	82
3.8. Приращение данных.....	82
3.8.1. Приращение данных для изображений	82
3.8.2. Приращение данных для текста	84
3.9. Обработка несбалансированных данных.....	85
3.9.1. Выборка с избытком.....	86
3.9.2. Выборка с недостатком.....	87
3.9.3. Гибридные стратегии	87
3.10. Стратегии выборки данных.....	88
3.10.1. Простая случайная выборка	89
3.10.2. Систематическая выборка.....	90
3.10.3. Стратифицированная выборка.....	90
3.11. Хранение данных	90
3.11.1. Форматы данных	91
3.11.2. Уровни хранения данных	92
3.11.3. Версионирование данных	94
3.11.4. Документация и метаданные	96
3.11.5. Жизненный цикл данных.....	96
3.12. Дополнительные рекомендации по работе с данными	97
3.12.1. Воспроизводимость.....	97
3.12.2. Сначала данные, потом алгоритм.....	97
3.13. Резюме	98

Глава 4. Конструирование признаков	100
4.1. Зачем конструировать признаки.....	100
4.2. Как конструируются признаки.....	101
4.2.1. Конструирование признаков для текста	102
4.2.2. Почему мешок слов работает.....	105
4.2.3. Преобразование категориальных признаков в числа.....	105
4.2.4. Хеширование признаков	108
4.2.5. Тематическое моделирование	109
4.2.6. Признаки для временных рядов	112
4.2.7. Проявите свои творческие способности.....	114
4.3. Штабелирование признаков.....	115
4.3.1. Штабелирование векторов признаков	115
4.3.2. Штабелирование индивидуальных признаков	116
4.4. Свойства хороших признаков	117

4.4.1. Высокая предсказательная способность	117
4.4.2. Быстрое вычисление	117
4.4.3. Надежность	118
4.4.4. Некоррелированность	118
4.4.5. Другие свойства	118
4.5. Отбор признаков	119
4.5.1. Отрезание длинного хвоста	119
4.5.2. Voruta	120
4.5.3. L1-регуляризация	123
4.5.4. Зависящий от задачи отбор признаков	123
4.6. Синтезирование признаков	123
4.6.1. Дискретизация признаков	124
4.6.2. Синтез признаков из реляционных данных	125
4.6.3. Синтезирование признаков по данным	126
4.6.4. Синтезирование признаков по другим признакам	127
4.7. Обучение признаков на данных	128
4.7.1. Погружения слов	128
4.7.2. Погружения документов	130
4.7.3. Погружения всего, чего угодно	131
4.7.4. Выбор размерности погружения	132
4.8. Понижение размерности	132
4.8.1. Быстрое понижение размерности методом PCA	133
4.8.2. Понижение размерности с целью визуализации	133
4.9. Масштабирование признаков	133
4.9.1. Нормировка	134
4.9.2. Стандартизация	135
4.10. Просачивание данных при конструировании признаков	136
4.10.1. Возможные проблемы	136
4.10.2. Решение	136
4.11. Хранение и документирование признаков	136
4.11.1. Файл схемы	137
4.11.2. Хранилище признаков	138
4.12. Рекомендации по конструированию признаков	141
4.12.1. Генерируйте много простых признаков	141
4.12.2. Повторно используйте унаследованные системы	141
4.12.3. Используйте идентификаторы как признаки, когда это необходимо	142
4.12.4. ...но по возможности уменьшайте количество значений	142
4.12.5. Осторожнее со счетчиками	143
4.12.6. Отбирайте признаки, когда необходимо	143
4.12.7. Тщательно тестируйте код	144
4.12.8. Синхронизируйте код, модель и данные	144
4.12.9. Изолируйте код выделения признаков	144
4.12.10. Сериализуйте модель и экстрактор признаков совместно	145
4.12.11. Протоколируйте значения признаков	145
4.13. Резюме	145

Глава 5. Обучение модели с учителем (часть 1)	147
5.1. Прежде чем приступить к работе над моделью.....	148
5.1.1. Проверка согласованности со схемой.....	148
5.1.2. Определение достижимого уровня качества.....	148
5.1.3. Выбор метрики качества.....	149
5.1.4. Выбирайте правильный ориентир.....	149
5.1.5. Разбиение данных на три набора.....	151
5.1.6. Предварительные условия для обучения с учителем.....	152
5.2. Представление меток для машинного обучения.....	153
5.2.1. Многоклассовая классификация.....	153
5.2.2. Многозначная классификация.....	154
5.3. Выбор алгоритма обучения.....	154
5.3.1. Основные свойства алгоритма обучения.....	155
5.3.2. Выборочная проверка алгоритмов.....	156
5.4. Построение конвейера.....	158
5.5. Оценивание качества модели.....	159
5.5.1. Метрики качества для регрессии.....	159
5.5.2. Метрики качества для классификации.....	160
5.5.3. Метрики качества для ранжирования.....	165
5.6. Настройка гиперпараметров.....	168
5.6.1. Поиск на сетке.....	169
5.6.2. Случайный поиск.....	170
5.6.3. Поиск с измельчением.....	170
5.6.4. Другие методы.....	172
5.6.5. Перекрестная проверка.....	172
5.7. Обучение поверхностной модели.....	173
5.7.1. Стратегия обучения поверхностной модели.....	173
5.7.2. Сохранение и восстановление модели.....	174
5.8. Компромисс между смещением и дисперсией.....	175
5.8.1. Недообучение.....	175
5.8.2. Переобучение.....	176
5.8.3. Компромисс.....	177
5.9. Регуляризация.....	179
5.9.1. L1- и L2-регуляризации.....	179
5.9.2. Другие формы регуляризации.....	180
5.10. Резюме.....	180
Глава 6. Обучение модели с учителем (часть 2)	183
6.1. Стратегия обучения глубоких моделей.....	183
6.1.1. Стратегия обучения нейронной сети.....	184
6.1.2. Метрика качества и функция стоимости.....	184
6.1.3. Стратегии инициализации параметров.....	187
6.1.4. Алгоритмы оптимизации.....	187
6.1.5. Планы уменьшения скорости обучения.....	191
6.1.6. Регуляризация.....	192
6.1.7. Определение размера сети и настройка гиперпараметров.....	193

6.1.8. Работа с несколькими входами	195
6.1.9. Работа с несколькими выходами.....	196
6.1.10. Перенос обучения.....	197
6.2. Штабелирование моделей.....	199
6.2.1. Типы ансамблевого обучения.....	199
6.2.2. Алгоритм штабелирования моделей	200
6.2.3. Просачивание данных при штабелировании моделей.....	201
6.3. Борьба со сдвигом распределения.....	202
6.3.1. Обработка сдвига распределения	202
6.3.2. Состязательная проверка	202
6.4. Обработка несбалансированных наборов данных	203
6.4.1. Взвешивание классов	203
6.4.2. Ансамбль перераспределенных наборов данных.....	204
6.4.3. Другие методы	205
6.5. Калибровка модели.....	205
6.5.1. Хорошо откалиброванные модели.....	205
6.5.2. Методы калибровки	207
6.6. Поиск неполадок и анализ ошибок	208
6.6.1. Причины плохого поведения модели.....	208
6.6.2. Итеративное уточнение модели.....	209
6.6.3. Анализ ошибок.....	209
6.6.4. Анализ ошибок в комплексных системах.....	211
6.6.5. Использование расслоенных метрик.....	212
6.6.6. Исправление неправильных меток.....	212
6.6.7. Нахождение дополнительных примеров для пометки	213
6.6.8. Поиск неполадок при глубоком обучении	213
6.7. Рекомендации	215
6.7.1. Поставляйте хорошую модель.....	215
6.7.2. Доверяйте популярным реализациям с открытым исходным кодом	215
6.7.3. Оптимизируйте важную для бизнеса меру качества.....	216
6.7.4. При обновлении начинайте с нуля.....	216
6.7.5. Избегайте каскадов коррекций.....	217
6.7.6. Используйте каскадирование моделей с осторожностью	217
6.7.7. Пишите эффективный код, компилируйте и распараллеливайте	218
6.7.8. Тестируйте на старых и новых данных.....	219
6.7.9. Больше данных лучше, чем более умный алгоритм	220
6.7.10. Новые данные лучше более изоощренных признаков.....	220
6.7.11. Радуйтесь крохотным достижениям	220
6.7.12. Обеспечьте воспроизводимость	220
6.8. Резюме	221
Глава 7. Оценивание модели	224
7.1. Офлайнное и онлайнное оценивания	225
7.2. A/B-тестирование	227
7.2.1. G-критерий	228
7.2.2. Z-критерий.....	231

7.2.3. Заключительные замечания и предупреждения.....	233
7.3. Многорукий бандит	233
7.4. Статистические границы качества модели	236
7.4.1. Статистический интервал для ошибки классификации	236
7.4.2. Бутстреп статистического интервала	237
7.4.3. Бутстреп интервала предсказания для регрессии	238
7.5. Оценивание адекватности тестового набора.....	239
7.5.1. Нейронное покрытие.....	239
7.5.2. Мутационное тестирование	240
7.6. Оценивание свойств модели	240
7.6.1. Робастность	241
7.6.2. Справедливость	241
7.7. Резюме	243

Глава 8. Развертывание модели

8.1. Статическое развертывание	245
8.2. Динамическое развертывание на устройстве пользователя.....	245
8.2.1. Развертывание параметров модели	246
8.2.2. Развертывание сериализованного объекта	246
8.2.3. Развертывание в браузере.....	246
8.2.4. Плюсы и минусы	246
8.3. Динамическое развертывание на сервере	247
8.3.1. Развертывание на виртуальной машине	247
8.3.2. Развертывание в контейнере	248
8.3.3. Бессерверное развертывание.....	250
8.3.4. Потокное развертывание модели.....	251
8.4. Стратегии развертывания	253
8.4.1. Разовое развертывание	253
8.4.2. Немое развертывание	254
8.4.3. Канареечное развертывание.....	254
8.4.4. Многорукие бандиты	255
8.5. Автоматизированное развертывание, версионирование и метаданные	255
8.5.1. Объекты, сопровождающие модель	255
8.5.2. Синхронизация версий.....	256
8.5.3. Метаданные версии модели.....	256
8.6. Рекомендации по развертыванию модели	257
8.6.1. Эффективность алгоритма	257
8.6.2. Развертывание глубоких моделей.....	260
8.6.3. Кеширование	260
8.6.4. Формат доставки модели и кода.....	261
8.6.5. Начинайте с простой модели.....	263
8.6.6. Тестируйте на посторонних	263
8.7. Резюме.....	264

Глава 9. Выполнение, мониторинг и сопровождение модели.....

9.1. Свойства среды выполнения модели.....	266
9.1.1. Безопасность и корректность	267

9.1.2. Простота развертывания	267
9.1.3. Гарантии правильности модели	268
9.1.4. Простота восстановления	268
9.1.5. Предотвращение расхождений между обучением и выполнением	268
9.1.6. Предотвращение скрытых петель обратной связи	269
9.2. Режимы выполнения модели	269
9.2.1. Выполнение в пакетном режиме.....	270
9.2.2. Обслуживание запроса со стороны человека	270
9.2.3. Обслуживание запроса со стороны машины.....	272
9.3. Выполнение модели на практике	273
9.3.1. Готовность к ошибкам.....	273
9.3.2. Отношение к ошибкам.....	274
9.3.3. Готовность к изменениям и отношение к ним	275
9.3.4. Готовность к особенностям человеческой природы и отношение к ним	277
Избегайте путаницы	277
Умерьте ожидания.....	277
Завоевывайте доверие	277
Не переутомляйте пользователя.....	278
Остерегайтесь фактора отторжения	278
9.4. Мониторинг модели	278
9.4.1. Что может пойти не так?.....	279
9.4.2. Что и как мониторить	280
9.4.3. Что протоколировать	282
9.4.4. Мониторинг неправомерного использования	283
9.5. Сопровождение модели	283
9.5.1. Когда обновлять	284
9.5.2. Как обновлять.....	285
9.6. Резюме	288
Глава 10. Заключение	290
10.1. Сухой остаток.....	290
10.2. Что еще почитать	294
10.3. Благодарности	295
Предметный указатель	297

От издательства

Отзывы и пожелания

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

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

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

Список опечаток

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

Нарушение авторских прав

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

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

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

Вступительное слово

Хочу открыть вам один секрет: когда говорят «машинное обучение», складывается впечатление, будто речь идет только об одной дисциплине. А вот и нет! На самом деле есть два вида машинного обучения, и они различаются так же, как придумывание новых блюд и изобретение кухонных принадлежностей. То и другое занятия достойны уважения, но путать их не надо: не станете же вы поручать шеф-повару разработку кухонной плиты, а инженеру-электрику – печь хлеб!

Увы, почти все путают эти два вида машинного обучения. Поэтому неудивительно, что так много компаний терпят крах на этапе машинного обучения. Начинаящим, похоже, никто не объяснил, что все то, что преподают на курсах машинного обучения и пишут в учебниках, – это теоретическое машинное обучение: как конструировать кухонные плиты (а также микроволновки, блендеры, тостеры, чайники ... и мойку, в которой все смешивается!) с нуля, а не как придумывать новые блюда и готовить на огромную компанию. Иными словами, если ваша цель – создавать инновационные решения конкретных задач на базе МО, то вам нужно прикладное, а не теоретическое машинное обучение. Поэтому большинство книг – не для вас.

А теперь хорошая новость! Перед вами первая книга, посвященная именно прикладному машинному обучению. Да, вы ее нашли! Настоящая прикладная иголка в стоге теоретического сена. Поздравляю, любезный читатель... если, конечно, вы не искали книгу, которая помогла бы отточить навыки проектирования алгоритмов общего назначения. Но тогда я надеюсь, что автор не будет на меня в обиде за то, что я посоветую вам купить чуть ли не любую другую книгу по МО. А эта от них отличается.

В 2016 году, разрабатывая курс Google по прикладному МО «Как подружиться с машинным обучением», полюбившийся нашим инженерам и руководителям групп – а их больше десяти тысяч, – я придала ему примерно такую же структуру, как в этой книге. А все потому, что действовать в правильном порядке очень важно в прикладных дисциплинах. Попытка выполнить определенные шаги, прежде чем будут завершены другие, может иметь печальные последствия: от зря потраченного времени до феерического краха проекта. На самом деле сходство оглавления этой книги и моего курса как раз и побудило меня прочитать ее. И в полном согласии с теорией параллельной эволюции я узрела в авторе единомышленника, терзаемого по ночам отсутствием ресурсов по прикладному машинному обучению, одной из потенциально самых полезных и вместе с тем самых недопонятых отраслей инженерной практики, – терзаемого с такой силой, что он захотел как-то исправить ситуацию. Так что если вы вознамерились захлопнуть книгу, то, пожалуйста, сделайте мне одолжение – задумайтесь на минутку, почему главы следуют именно в таком порядке. Обещаю, уже это принесет полезные плоды.

Так что же будет в книге дальше? Машинное обучение можно сравнить с руководством по приготовлению еды в массовом масштабе. Поскольку вы

еще не прочитали книгу, опишу содержание книги в кулинарных терминах. Вам нужно определить, что имеет смысл приготовить и каковы цели (*принятие решений и управление продуктом*), понять, кто является поставщиками и клиентами (*знакомство с предметной областью и деловое чутье*), как обрабатывать ингредиенты в большом количестве (*инженерия и анализ данных*), как быстро оценивать много разных сочетаний ингредиентов, пригодных для создания потенциальных рецептов блюд (*инженерная разработка прототипа*), как проверить качество рецепта (*статистика*), как эффективно превратить потенциальный рецепт в миллионы порций (*организация производства*) и как гарантировать, что блюда останутся высококачественными, даже если грузовик привез тонну картошки вместо заказанного риса (*обеспечение надежности*). Эта книга – одна из немногих, где рассматриваются все этапы технологического процесса.

Теперь самое время подпустить немного дегтя в бочку меда. Это хорошая книга. Правда. Но не идеальная. Иногда автор срезает углы – как часто делают инженеры, профессионально занимающиеся машинным обучением, – хотя в целом все изложено правильно. А поскольку предмет книги быстро эволюционирует, автор и не пытается всегда быть на переднем крае. Но даже в отсутствие совершенства книгу все равно стоит прочитать. Учитывая, как мало на рынке подробных руководств по прикладному машинному обучению, ясное и последовательное введение в эту тему дорогого стоит. Я очень рада, что такая книга появилась!

Что мне в ней очень нравится, так это полнота, с которой излагается самая важная вещь, которую нужно знать о машинном обучении: ошибки возможны... и иногда болезненные ошибки. Как любят говорить мои коллеги, занимающиеся надежностью: «Надежда на лучшее – это не стратегия». Надеяться, что ошибок не будет, – худший из возможных подходов. Эта книга устроена иначе – и лучше. Она быстро заставляет расстаться с ложным чувством безопасности, которое вы могли лелеять, надеясь построить систему ИИ, более «умную», чем вы сами. (Забудьте об этом, ничего не получится.) Затем автор становится вашим заботливым проводником по пути, на котором можно надеяться на самых разнообразных ошибок, и объясняет, как их предотвратить, обнаружить и исправить. В книге прекрасно подчеркнута важность мониторинга, описаны подходы к сопровождению модели, рассказано, что делать, если что-то ломается, как спланировать резервные стратегии на случай разного рода отказов и как управлять ожиданиями пользователей (есть также раздел о том, что делать, если пользователями являются машины). Эти вопросы очень важны при практическом применении машинного обучения, но в других книгах для них часто не находится места. Впрочем, эта книга не из их числа.

Если вы собираетесь использовать машинное обучение для решения крупномасштабных практических задач, могу только порадоваться, что вы наткнулись на эту книгу. Читайте с удовольствием!

Кэсси Козырьков,
главный специалист по теории принятия решений в Google,
автор курса «MakingFriends with Machine Learning»
на облачной платформе Google

Сентябрь 2020

Предисловие

За последние несколько лет машинное обучение (МО) для многих стало синонимом искусственного интеллекта. И хотя машинное обучение как научная дисциплина существует уже несколько десятков лет, в мире найдется лишь горстка организаций, в полной мере осознавших его потенциал.

Несмотря на доступность современных библиотек, пакетов и каркасов машинного обучения с открытым исходным кодом, которые поддерживаются ведущими организациями и широким сообществом ученых и программистов, большинство компаний все еще испытывают трудности с применением машинного обучения к решению практических деловых задач.

Одна из проблем – нехватка кадров. Но, даже располагая талантливыми специалистами по машинному обучению и анализу данных, в 2020 году большинство организаций все еще тратят от 31 до 90 дней на развертывание всего одной модели, а 18 % – более 90 дней, у некоторых на доведение идеи до продукта уходит даже больше года. Основные проблемы, с которыми приходится сталкиваться при освоении потенциала МО, например управление версиями модели, воспроизводимость результатов и масштабирование, имеют скорее инженерную, чем научную природу.

Существует много книг по машинному обучению – как теоретических, так и практических. Из типичного учебника вы можете узнать о разных типах машинного обучения, об основных семействах алгоритмов, о том, как они работают и как с их помощью создавать модели из данных.

В типичном учебнике меньше внимания уделяется инженерным аспектам реализации проектов машинного обучения. Такие вопросы, как сбор, хранение и предобработка данных, конструирование признаков, а также тестирование и отладка моделей, развертывание и вывод из эксплуатации, сопровождение на этапе выполнения и в процессе эксплуатации, зачастую остаются за кадром.

Эта книга призвана заполнить пробел.

НА КОГО РАССЧИТАНА ЭТА КНИГА

Я предполагаю, что читатель знаком с основами машинного обучения и способен построить модель при наличии подходящим образом отформатированного набора данных, применяя свой любимый язык программирования или библиотеку. Если вы неуверенно владеете применением алгоритмов машинного обучения к данным и не видите разницы между логистической регрессией, методом опорных векторов и случайным лесом, то я рекомендую предварительно прочитать мою книгу «The Hundred-Page Machine Learning Book»¹.

¹ Бурков А. Машинное обучение без лишних слов. СПб.: Питер, 2020.

Целевая аудитория этой книги – аналитики данных, стремящиеся стать инженерами по машинному обучению, инженеры по машинному обучению, стремящиеся привести больше порядка в свою работу, студенты, изучающие машинное обучение, а также архитекторы программных систем, которым приходится иметь дело с моделями, разработанными аналитиками данных и инженерами по машинному обучению.

КАК ИСПОЛЬЗОВАТЬ ЭТУ КНИГУ

Эта книга представляет собой подробный обзор передовых практик и паттернов проектирования в области прикладного машинного обучения. Я рекомендую читать ее с начала до конца. Но можете читать главы в любом порядке, поскольку в них рассматриваются различные аспекты жизненного цикла проекта с использованием машинного обучения и прямых зависимостей между главами нет.

Андрей Бурков

Глава 1

Введение

Хотя предполагается, что читатель этой книги знаком с основами машинного обучения, все же важно начать с определений, чтобы у нас было общее понимание используемых в книге терминов.

Ниже я повторяю некоторые определения из главы 2 книги «Машинное обучение без лишних слов», а также даю несколько новых определений. Если вы читали мою первую книгу, то некоторые части этой главы могут показаться вам знакомыми.

После этой главы мы будем одинаково понимать такие понятия, как обучение с учителем и без учителя. У нас будет общий взгляд на прямо и косвенно используемые данные, первичные и аккуратные данные, обучающие и резервированные данные.

Мы будем знать о том, когда следует использовать машинное обучение, а когда этого делать не стоит, а также о различных формах машинного обучения, например: на основе модели и на основе экземпляров, глубокое и поверхностное, классификация и регрессия и т. д.

Наконец, мы определим предмет инженерии машинного обучения и представим жизненный цикл проекта машинного обучения.

1.1. ОБОЗНАЧЕНИЯ И ОПРЕДЕЛЕНИЯ

Начнем с базовых математических обозначений и определим термины и понятия, к которым часто будем обращаться в этой книге.

1.1.1. Структуры данных

Скаляром¹ называется простое числовое значение, например 15 или -3.25 . Переменные или константы, принимающие скалярные значения, обозначаются курсивом, например x или a .

Вектором называется упорядоченный список скалярных значений, именуемых атрибутами. Мы обозначаем вектор полужирным шрифтом, например \mathbf{x} или \mathbf{w} . Векторы изображаются в виде направленных стрелок, а также

¹ Если термин выделен полужирным шрифтом, значит, он присутствует в алфавитном указателе в конце книги.

точек в многомерном пространстве. Иллюстрации трех двумерных векторов $\mathbf{a} = [2, 3]$, $\mathbf{b} = [-2, 5]$ и $\mathbf{c} = [1, 0]$ приведены на рис. 1.1. Атрибут вектора обозначается курсивной буквой с верхним индексом, например: $w^{(j)}$ или $x^{(j)}$. Индекс j обозначает конкретное **измерение** вектора, т. е. позицию атрибута в списке. Например, в векторе \mathbf{a} , показанном красным цветом на рис. 1.1, $a^{(1)} = 2$ и $a^{(2)} = 3$.

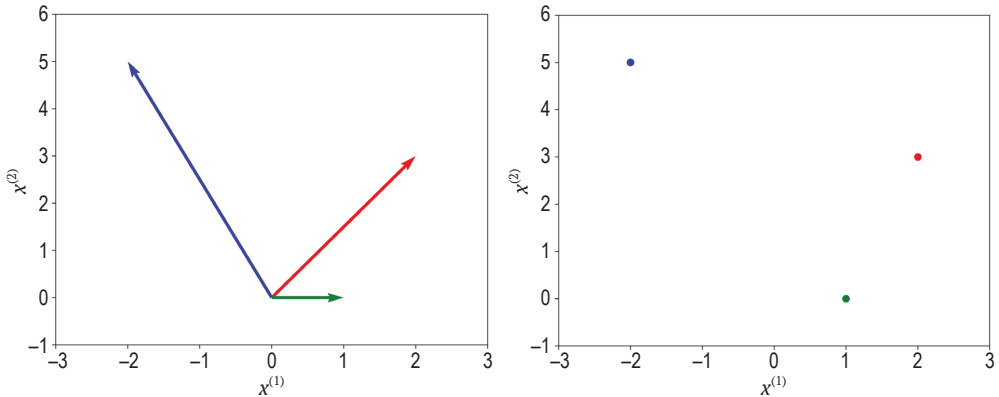


Рис. 1.1 ❖ Три вектора, представленных направленными стрелками и точками

Обозначение $x^{(j)}$ не следует путать с оператором возведения в степень, например 2 в x^2 (возведение в квадрат) или 3 в x^3 (возведение в куб). Если мы хотим применить оператор возведения в степень, скажем, в квадрат, к индексированному атрибуту вектора, то пишем: $(x^{(j)})^2$.

Переменная может иметь два и более индексов, например: $x_i^{(j)}$ или $x_{i,j}^{(k)}$. Так, в нейронных сетях $x_{i,u}^{(j)}$ обозначает j -й входной признак u -го блока в l -м слое.

Матрица – это прямоугольный массив чисел, организованный по строкам и столбцам. Ниже приведен пример матрицы с двумя строками и тремя столбцами:

$$\mathbf{A} = \begin{bmatrix} 2 & -2 & 1 \\ 3 & 5 & 0 \end{bmatrix}.$$

Матрицы обозначаются полужирными заглавными буквами, например \mathbf{A} или \mathbf{W} . Из примера матрицы \mathbf{A} выше видно, что матрицы можно трактовать как регулярные структуры, состоящие из векторов. И действительно, столбцами приведенной выше матрицы \mathbf{A} являются векторы \mathbf{a} , \mathbf{b} и \mathbf{c} , показанные на рис. 1.1.

Множеством называется неупорядоченная коллекция неповторяющихся элементов. Множество обозначается каллиграфической заглавной буквой, например \mathcal{S} . Множество чисел может быть конечным (содержать фиксированное количество значений). В этом случае его элементы перечисляются в фигурных скобках, например $\{1, 3, 18, 23, 235\}$ или $\{x_1, x_2, x_3, x_4, \dots, x_n\}$. Множество также может быть бесконечным и включать все значения в некотором

интервале. Множество, содержащее все значения между a и b включительно, обозначается $[a, b]$. Если же множество не включает граничные значения a и b , то оно обозначается (a, b) . Например, множество $[0, 1]$ включает среди прочего значения 0, 0.0001, 0.25, 0.784, 0.9995 и 1.0. Буквой \mathbb{R} обозначается множество всех чисел от минус бесконечности до плюс бесконечности.

Если элемент x принадлежит множеству \mathcal{S} , то мы пишем $x \in \mathcal{S}$. Новое множество \mathcal{S}_3 можно получить как **пересечение** двух множеств \mathcal{S}_1 и \mathcal{S}_2 . В этом случае мы пишем $\mathcal{S}_3 \leftarrow \mathcal{S}_1 \cap \mathcal{S}_2$. Например, $\{1, 3, 5, 8\} \cap \{1, 8, 4\}$ дает новое множество $\{1, 8\}$.

Новое множество \mathcal{S}_3 можно получить как **объединение** двух множеств \mathcal{S}_1 и \mathcal{S}_2 . В этом случае мы пишем $\mathcal{S}_3 \leftarrow \mathcal{S}_1 \cup \mathcal{S}_2$. Например, $\{1, 3, 5, 8\} \cup \{1, 8, 4\}$ дает новое множество $\{1, 3, 5, 8, 4\}$.

$|\mathcal{S}|$ обозначает размер множества \mathcal{S} , то есть число элементов в нем.

1.1.2. Заглавная сигма

Сумма множества чисел $\mathcal{X} = \{x_1, x_2, \dots, x_{n-1}, x_n\}$ или атрибутов вектора $\mathbf{x} = [x^{(1)}, x^{(2)}, \dots, x^{(m-1)}, x^{(m)}]$ обозначается следующим образом:

$$\sum_{i=1}^n x_i \stackrel{\text{def}}{=} x_1 + x_2 + \dots + x_{n-1} + x_n,$$

или

$$\sum_{i=1}^n x^{(j)} \stackrel{\text{def}}{=} x^{(1)} + x^{(2)} + \dots + x^{(m-1)} + x^{(m)}.$$

Здесь $\stackrel{\text{def}}{=}$ означает «по определению равно».

Евклидова норма вектора \mathbf{x} , обозначаемая $\|\mathbf{x}\|$, характеризует «размер», или «длину», вектора. Она определяется как $\sqrt{\sum_{j=1}^D (x^{(j)})^2}$.

В качестве расстояния между векторами \mathbf{a} и \mathbf{b} берется **евклидово расстояние**:

$$\|\mathbf{a} - \mathbf{b}\| \stackrel{\text{def}}{=} \sqrt{\sum_{i=1}^N (a^{(i)} - b^{(i)})^2}.$$

1.2. ЧТО ТАКОЕ МАШИННОЕ ОБУЧЕНИЕ

Машинное обучение – это раздел информатики, посвященный построению алгоритмов, которые работают с набором примеров, описывающих какое-то явление. Примеры могут поступать из природы, создаваться людьми или генерироваться другим алгоритмом.

Машинное обучение также можно определить как процесс решения практической задачи путем

- 1) сбора набора данных и
- 2) алгоритмического обучения **статистической модели** на этом наборе.

Предполагается, что эта статистическая модель каким-то образом используется для решения практической задачи. Для краткости я буду использовать термины «обучение» и «машинное обучение» как синонимы. По той же причине я буду часто говорить «модель», имея в виду статистическую модель.

Обучение бывает с учителем, без учителя, с частичным привлечением учителя и с подкреплением.

1.2.1. Обучение с учителем

В случае **обучения с учителем** аналитик работает с набором **помеченных примеров** $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$. Каждый элемент x_i называется **вектором признаков**. В информатике вектор – это одномерный массив. Одномерный массив, в свою очередь, является упорядоченной и индексированной последовательностью значений. Длина этой последовательности значений, D , называется **размерностью** вектора.

Вектор признаков – это вектор, в котором каждое измерение $j = 1 \dots D$ содержит значение, описывающее пример. Каждое такое значение называется **признаком** и обозначается $x^{(j)}$. Например, если каждый пример x представляет человека, то первый признак, $x^{(1)}$, может содержать рост в сантиметрах, второй признак, $x^{(2)}$, – вес в килограммах, $x^{(5)}$ – пол и т. д. Для всех примеров в наборе данных значение в позиции j вектора признаков всегда имеет один и тот же тип. Это означает, что если $x_i^{(2)}$ содержит вес в килограммах для некоторого i , то и для любого k от 1 до N $x_k^{(2)}$ будет содержать вес в килограммах. **Метка** y_i может быть либо элементом конечного множества классов $\{1, 2, \dots, C\}$, либо вещественным числом, либо более сложной структурой, такой как вектор, матрица, дерево или граф. Если явно не оговорено противное, то в этой книге y_i является либо элементом конечного множества классов, либо вещественным числом¹. Можно считать, что класс – это категория, к которой относится пример.

Скажем, если примерами являются сообщения электронной почты, а наша задача заключается в обнаружении спама, то есть два класса: спам и не спам. В случае обучения с учителем задача предсказания класса называется **классификацией**, а задача предсказания вещественного числа называется **регрессией**. Значение, которое должно быть предсказано моделью, обученной с учителем, называется **целевым показателем**, или целью. Примером регрессии является задача предсказания заработной платы сотрудника с учетом его опыта работы и знаний. Примером классификации является ситуация, когда врач вводит характеристики пациента в приложение, а то возвращает диагноз.

Различие между классификацией и регрессией показано на рис. 1.2. В случае классификации алгоритм обучения ищет линию (или, в более общем слу-

¹ Вещественное число – это величина, представляющая расстояние вдоль прямой от некоторой начальной точки на ней. Примеры: 0, -256.34, 1000, 1000.2.

чае, гиперповерхность), которая разделяет примеры разных классов. С другой стороны, в случае регрессии алгоритм обучения стремится отыскать линию или гиперповерхность, которая хорошо соответствует обучающим примерам.

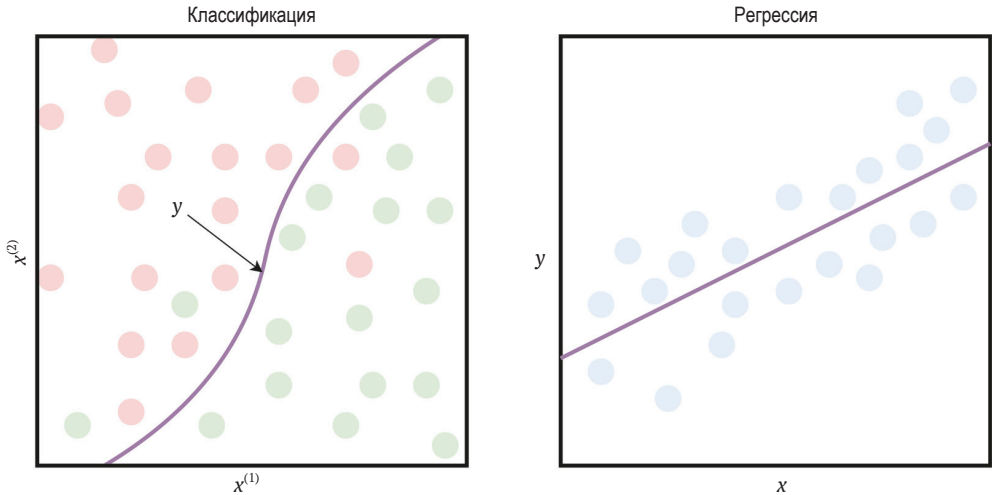


Рис. 1.2 ❖ Разница между классификацией и регрессией

Цель **алгоритма обучения с учителем** – использовать набор данных для порождения модели, которая на входе принимает вектор признаков \mathbf{x} , а на выходе выдает информацию, позволяющую вывести метку для этого вектора признаков. Например, модель, созданная с использованием набора данных о пациентах, может на входе принимать вектор признаков, описывающий пациента, и на выходе выдавать вероятность наличия y у пациента рака.

Даже если модель традиционно представляет собой математическую функцию, размышляя о действиях модели с входными данными, удобно думать, что модель «смотрит» на значения некоторых признаков на входе и, основываясь на опыте работы с аналогичными примерами, выводит значение. Это выходное значение представляет собой число или класс, «наиболее похожий» на метки, которые мы видели в прошлом в примерах с похожими значениями признаков. Такое описание выглядит упрощенно, но именно так работает модель решающего дерева и алгоритм k ближайших соседей.

1.2.2. Обучение без учителя

В случае **обучения без учителя** набор данных содержит **непомеченные примеры** $\{x_1, x_2, \dots, x_N\}$. Как и раньше, \mathbf{x} – вектор признаков, а цель **алгоритма обучения без учителя** заключается в порождении модели, которая на входе принимает вектор признаков \mathbf{x} и преобразовывает его либо в другой вектор, либо в значение, используемое для решения практической задачи. Например, в случае **кластеризации** для каждого вектора признаков из на-

бора данных модель возвращает ИД кластера. Кластеризация применяется для отыскания групп похожих объектов в большом наборе объектов, например изображений или текстовых документов. Используя кластеризацию, аналитик может, например, выбрать достаточно репрезентативное, но малое подмножество непомеченных примеров из большого набора, чтобы потом пометить их вручную: из каждого кластера выбирается всего несколько примеров, вместо того чтобы отбирать непосредственно из большого набора с риском выбрать очень похожие примеры.

В задаче **понижения размерности** выходом модели является вектор признаков с меньшим числом измерений, чем на входе. Например, у исследователя имеется вектор признаков, слишком сложный для визуализации (поскольку число измерений больше трех). Модель понижения размерности может преобразовать этот вектор в другой (сохраняя часть информации) – двумерный или трехмерный. Этот новый вектор признаков можно изобразить на графике.

В задаче **обнаружения выбросов** выходом является действительное число, показывающее, насколько входной вектор признаков отличается от «типичного» примера в наборе данных. Обнаружение выбросов применяется для решения задачи о проникновении в сеть (путем обнаружения аномальных сетевых пакетов, которые отличаются от типичного пакета в «нормальном» трафике) или обнаружения новизны (например, документа, отличающегося от других документов в наборе).

1.2.3. Обучение с частичным привлечением учителя

В случае **обучения с частичным привлечением учителя** набор данных содержит как помеченные, так и непомеченные примеры. Обычно число непомеченных примеров намного превышает число помеченных. Цель **алгоритма с частичным привлечением учителя** такая же, как у алгоритма обучения с учителем. Мы надеемся, что, располагая большим числом непомеченных примеров, алгоритм обучения сможет отыскать (можно было бы сказать «породить» или «вычислить») лучшую модель.

1.2.4. Обучение с подкреплением

Обучение с подкреплением – это раздел машинного обучения, в котором рассматривается ситуация, когда машина (именуемая агентом) «обитает» в окружающей среде и способна воспринимать состояние этой среды как вектор признаков. Машина может выполнять действия в незаключительных состояниях. Разные действия приносят разные вознаграждения, а также могут переводить машину в другое состояние окружающей среды. Типичная цель алгоритма обучения с подкреплением – обучиться оптимальной **политике**.

Оптимальная политика – это функция (аналогичная модели в обучении с учителем), которая на входе принимает вектор признаков состояния, а на

Конец ознакомительного фрагмента.

Приобрести книгу можно

в интернет-магазине

«Электронный универс»

e-Univers.ru