
Оглавление

Отзывы о книге «Google BigQuery. Всё о хранилищах данных, аналитике и машинном обучении»	13
Предисловие	15
Для кого написана эта книга?.....	16
Условные обозначения	16
Использование примеров программного кода	17
Благодарности	17
От издательства.....	18
Глава 1. Что такое Google BigQuery?	19
Архитектуры обработки данных	19
Система управления реляционными базами данных.....	20
Фреймворк MapReduce.....	22
BigQuery: бессерверный распределенный движок SQL.....	23
Работа с BigQuery.....	25
Анализ наборов данных.....	25
ETL, EL и ELT	26
Эффективная аналитика	28
Простота управления.....	30
История появления BigQuery.....	31
Что позволило создать BigQuery?	34
Отделение вычислений от хранилища	35
Хранилище и сетевая инфраструктура.....	36
Управляемое хранилище.....	37
Интеграция с платформой Google Cloud	39
Безопасность и соответствие требованиям	40
Выводы	41

Глава 2. Основы запросов	42
Простые запросы.....	44
Извлечение записей с помощью SELECT.....	44
Создание псевдонимов столбцов с помощью AS	46
Фильтрация с WHERE	48
SELECT *, EXCEPT, REPLACE.....	49
Подзапросы с WITH	50
Сортировка с ORDER BY.....	50
Агрегирование.....	51
Агрегирование с GROUP BY	51
Подсчет записей с COUNT	52
Фильтрация сгруппированных значений с HAVING.....	53
Поиск уникальных значений с DISTINCT	54
Краткое руководство по массивам и структурам.....	55
Создание массивов с помощью ARRAY_AGG.....	57
Массив структур STRUCT	59
Кортежи	60
Работа с массивами	60
Развертывание массива	61
Соединение таблиц	62
Основы соединения таблиц	63
Оператор внутреннего соединения INNER JOIN.....	66
Оператор перекрестного соединения CROSS JOIN.....	67
Оператор внешнего соединения OUTER JOIN.....	69
Сохранение и совместное использование	70
История запросов и кеширование	70
Сохранение запросов	72
Представления и общедоступные запросы.....	73
Выводы	74
Глава 3. Типы данных, функции и операторы.....	75
Числовые типы и функции.....	76
Математические функции	77
Стандартное вещественное деление.....	78

Функции SAFE	78
Сравнение	79
Точные десятичные вычисления с NUMERIC	80
Тип BOOL	81
Логические операции	81
Условные выражения	83
Обработка NULL с помощью COALESCE	84
Явное и неявное приведение типов	85
Использование COUNTIF, чтобы избежать приведения логических значений	87
Строковые функции	88
Интернационализация	89
Формирование и парсинг строк	91
Функции для обработки строк	91
Функции преобразования	92
Регулярные выражения	92
Краткие итоги по строковым функциям	94
Операции со значениями TIMESTAMP	95
Парсинг и форматирование отметок времени	95
Извлечение календарных данных	97
Арифметические операции с временными метками	98
DATE, TIME и DATETIME	99
Функции для работы с географическими координатами	100
Выводы	101
Глава 4. Загрузка данных в BigQuery	104
Основы	104
Загрузка из локального источника	105
Корректировка схемы	112
Копирование в новую таблицу	116
Управление данными (DDL и DML)	116
Эффективная загрузка данных	118
Федеративные запросы и внешние источники данных	121
Как использовать федеративные запросы	121
Когда использовать федеративные запросы и внешние источники данных	125

Интерактивное исследование и запрос данных из Google Sheets	132
Запросы SQL для выборки данных из Cloud Bigtable	141
Передача и экспорт данных.....	147
Служба передачи данных Data Transfer Service.....	147
Экспортирование журналов Stackdriver	153
Использование Cloud Dataflow для чтения/записи в BigQuery	154
Перемещение локальных данных.....	159
Методы миграции данных	159
Выводы	162
Глава 5. Разработка с BigQuery	163
Программный доступ	163
Доступ к BigQuery через REST API.....	163
Google Cloud Client Library	171
Доступ к BigQuery из инструментов исследования данных	188
Блокноты в Google Cloud Platform	188
Работа с BigQuery, pandas и Jupyter	193
Работа с BigQuery из R.....	198
Cloud Dataflow	199
Драйверы JDBC/ODBC.....	202
Внедрение данных из BigQuery в Google Slides (в G Suite)	203
Vash-скрипты для BigQuery	205
Создание наборов данных и таблиц	205
Выполнение запросов	208
Объекты BigQuery	210
Выводы	212
Глава 6. Архитектура BigQuery	213
Архитектура высокого уровня	213
Жизненный цикл запроса	213
Обновление BigQuery	218
Система обработки запросов (Dremel).....	219
Архитектура Dremel	221
Выполнение запроса	226
Хранилище.....	241
Хранение данных	241

Метаданные	248
Выводы	258
Глава 7. Оптимизация производительности и затрат	259
Принципы производительности	259
Ключевые составляющие производительности.....	260
Управление затратами.....	260
Измерение производительности и поиск проблем.....	262
Определение скорости выполнения запроса с помощью REST API.....	263
Определение скорости выполнения запроса с помощью BigQuery Workload Tester	265
Выявление проблем в рабочих нагрузках с помощью Stackdriver	267
Чтение плана запроса	269
Увеличение скорости выполнения запросов	274
Минимизация ввода/вывода.....	275
Кеширование результатов предыдущих запросов	280
Эффективное выполнение соединений	284
Исключение перегрузки рабочих серверов	293
Использование приближенных функций агрегирования.....	296
Оптимизация хранения данных и доступа к ним.....	299
Минимизация сетевых издержек	300
Выбор эффективного формата хранения	303
Секционирование таблиц для уменьшения объема сканирования	313
Кластеризация таблиц на основе ключей с большой мощностью множества	316
Случаи использования, нечувствительные ко времени	321
Пакетные запросы	321
Загрузка файлов	323
Выводы	324
Контрольный список	324
Глава 8. Продвинутые запросы	326
Множественные запросы	326
Параметризованные запросы.....	327
Пользовательские функции SQL	332
Повторное использование частей запросов	337

Продвинутый SQL	341
Работа с массивами	342
Оконные функции	351
Метаданные таблиц	356
Язык определения данных и язык манипулирования данными	360
За пределами SQL	365
Пользовательские функции на JavaScript	366
Скрипты	367
Продвинутые функции	375
Геоинформационная система BigQuery	375
Полезные статистические функции	383
Алгоритмы хеширования	385
Выводы	389
Глава 9. Машинное обучение в BigQuery	390
Что такое машинное обучение?	390
Формулировка задачи машинного обучения	391
Типы задач машинного обучения	392
Построение регрессионной модели	396
Выбор метки	396
Выбор признаков в наборе данных	397
Создание обучающего набора данных	401
Обучение и оценка модели	402
Получение прогнозов с помощью модели	404
Исследование весов модели	407
Более сложные регрессионные модели	409
Создание модели классификации	414
Обучение	415
Оценка	416
Прогнозирование	417
Выбор порога	418
Настройка механизма машинного обучения в BigQuery	420
Управление делением данных	420
Балансировка классов	422
Регуляризация	422

Кластеризация методом k-средних.....	423
Выбор признаков для кластеризации.....	424
Кластеризация пунктов проката велосипедов	425
Кластеризация.....	426
Исследование кластеров.....	427
Принятие решений на основе данных.....	429
Рекомендательные системы	430
Набор данных MovieLens.....	430
Разложение матрицы	432
Получение рекомендаций	434
Включение информации о пользователях и фильмах.....	436
Нестандартные модели машинного обучения в GCP.....	443
Настройка гиперпараметров	444
AutoML	448
Поддержка TensorFlow.....	450
Выводы	453
Глава 10. Администрирование и безопасность BigQuery	455
Защищенность инфраструктуры.....	455
Управление идентификацией и доступом	457
Идентификация.....	457
Роль	458
Ресурс.....	461
Администрирование BigQuery.....	462
Управление заданиями	462
Авторизация пользователей	463
Восстановление удаленных записей и таблиц	463
Непрерывная интеграция/непрерывное развертывание	464
Экспорт биллинга — получение информации о расходах.....	467
Оперативные панели, мониторинг и журналы аудита.....	470
Доступность, восстановление после отказа и шифрование	471
Зоны, регионы и объединения регионов.....	471
BigQuery и обработка отказов	472
Сохранность, резервное копирование и восстановление после аварий.....	476
Конфиденциальность и шифрование.....	477

Соответствие требованиям законодательств	478
Местоположение данных.....	478
Ограничение доступа к подмножествам данных.....	480
Удаление всех сделок, связанных с конкретным физическим лицом	483
Предотвращение потери данных.....	487
СМЕК.....	488
Защита от утечки данных.....	490
Выводы	491
Об авторах	493
Об обложке	494

Формулировка задачи машинного обучения

Допустим, вы управляете несколькими сотнями кинотеатров по всей стране и хотите спрогнозировать, сколько билетов будет продано за определенное время показа в конкретном кинотеатре, — это может очень пригодиться при планировании графика показа фильмов. При наличии статистических данных о показе фильмов задачу машинного обучения можно сформулировать следующим образом: на основе информации о прокате, имеющейся в наборе данных, узнать, сколько билетов было продано на каждый сеанс в каждом кинотеатре, затем применить полученную модель к фильму-кандидату и определите, каким будет спрос на этот фильм в конкретное время.

Атрибуты фильма, которые мы используем в качестве входных данных в модели машинного обучения, называются *признаками* модели. Метка — это прогнозное значение, которое требуется узнать, и в данном случае метка — это количество проданных билетов. Ниже приводятся некоторые примеры признаков, которые вы можете включить в свою модель:

- Оценка содержимого кинофильма¹ (например, PG-13 означает, что фильм рекомендуется для детей младше 13 лет).
- Фильм будет показан в рабочие или выходные дни?
- В какое время дня будет демонстрироваться фильм (днем, вечером, ночью)?
- Жанр фильма (комедия, триллер и т. д.).
- Как давно снят фильм (в днях).
- Средний рейтинг фильма по оценкам критиков (по шкале от 1 до 10).
- Общая сумма кассовых сборов за предыдущий фильм этого же режиссера, если применимо.
- Общая сумма кассовых сборов за предыдущий фильм с этим же актером в главной роли, если применимо.
- Местоположение кинотеатра.
- Тип кинотеатра (например, большой кинотеатр с несколькими залами, кинотеатр под открытым небом для автомобилистов, кинозал в торговом центре и т. д.).

Обратите внимание, что название фильма само по себе не является хорошим обучающим признаком.² Даже если фильм «Шпион, выйди вон!» («Tinker Tailor

¹ См. https://en.wikipedia.org/wiki/Motion_picture_content_rating_system.

² Отдельные слова в названии фильма могут оказаться неплохими признаками, если применить к названиям фильмов методы обработки естественного языка, такие как лексемизация, морфологический поиск и получение векторных представлений слов. Также могут пригодиться признаки, вычисленные из названий фильмов; например,

Soldier Spy») 2011 года будет присутствовать в наборе обучающих данных, нам будет неинтересно прогнозировать прокат именно этого фильма (потому что он уже был показан в наших кинотеатрах). Гораздо интереснее спрогнозировать сборы, скажем, от проката фильма «Глубоководный горизонт» («Deep Water Horizon»), еще одного триллера с похожими отзывами критиков, выпущенного в 2016 году.

То есть модель машинного обучения должна опираться на особенности фильма (характеризующие фильм), а не на признаки, однозначно его идентифицирующие. В этом случае модель сможет догадаться, что при одинаковых условиях показа фильм «Глубоководный горизонт» принесет такую же прибыль, как и «Шпион, выйди вон!», потому что они относятся к одному жанру и имеют похожий рейтинг.

Первые четыре признака (рейтинг, дни показа, время показа, жанр) являются категориальными, то есть принимают одно из конечного числа возможных значений. В BigQuery категориальным считается любой признак, являющийся строкой. Если в базе данных категориальные признаки представлены значениями другого типа (например, признак времени показа может быть числом, например 1430, или отметкой времени), в запросе они должны преобразовываться в строки. Следующие четыре признака (время с момента выпуска, рейтинги критиков, кассовые сборы за предыдущий фильм этого же режиссера и за предыдущий фильм с этим же актером в главной роли) — числовые, то есть имеют значимое числовое выражение. Последние два признака (тип и местоположение кинотеатра) должны быть представлены особым образом; возможные варианты мы обсудим далее в этой главе.

Меткой, или правильным прогнозом, является количество проданных билетов. Во время обучения модели машинного обучения BigQuery передаются входные признаки и соответствующие метки, на основе которых создается модель, обобщающая эту информацию (рис. 9.1). Затем, во время прогнозирования, обученную модель можно применить к новому набору входных признаков и получить оценку количества билетов, которое можно продать, если запланировать показ фильма в определенное время и в определенном месте.

Типы задач машинного обучения

Мы часто используем разные модели и методы машинного обучения в зависимости от характера входных признаков и меток. В этом подразделе мы кратко перечислим поддерживаемые типы задач машинного обучения, а в оставшейся части главы подробно рассмотрим способы их решения.

хорошим прогнозирующим элементом может послужить длина названия или наличие в нем слова «шпион».

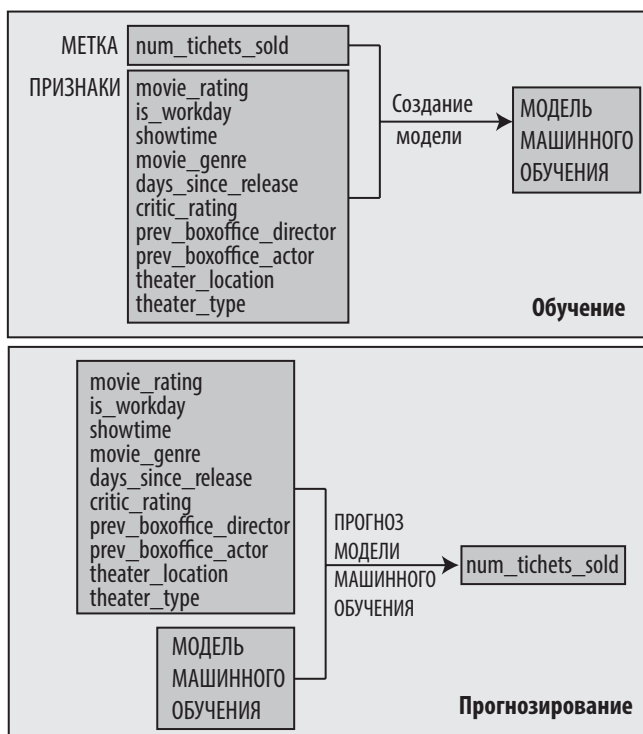


Рис. 9.1. Во время обучения модели машинного обучения передаются входные признаки и соответствующие метки. Затем обученную модель можно использовать для прогнозирования. На основе набора входных признаков модель способна спрогнозировать метку

Регрессия

В примере, описанном в предыдущем разделе, требовалось спрогнозировать количество проданных билетов на конкретный фильм. В данном случае метка является числом, и задачи машинного обучения подобного типа называют *регрессией*.

Классификация

Если результатом (меткой) задачи машинного обучения является категориальная переменная, такие задачи относятся к разряду задач *классификации*. Модели классификации возвращают вероятность принадлежности записи к категории (классу), определяемой меткой. Например, чтобы получить модель машинного обучения, прогнозирующую успех шоу, вы могли бы создать модель классификации, возвращающую вероятность успеха.

Многие задачи классификации имеют два класса, например: успех или провал шоу, высокий или низкий спрос на товар, прибытие рейса вовремя или с опозданием. Это задачи так называемой *бинарной классификации*. В таких случаях столбец метки должен иметь значение True или False либо 1 или 0. Прогнозом такой модели является вероятность, что метка будет иметь значение True. Чтобы определить наиболее вероятный класс, мы обычно устанавливаем пороговую вероятность равной 0.5.

Задача классификации может иметь несколько классов. Например, возвращаясь к нашему сценарию проката велосипедов, может понадобиться предсказать пункт проката, куда будет возвращен велосипед, и поскольку для этой категориальной метки существуют сотни возможных значений, эта задача относится к разряду задач *мультиклассовой (множественной) классификации*. Результатом такой модели машинного обучения является набор вероятностей, по одной для каждого пункта проката, и сумма этих вероятностей будет равна 1.0. В задачах множественной классификации, как правило, наибольший интерес представляют три или пять наибольших предсказаний, а не фактические значения вероятностей.

Рекомендательная система

Особый случай многоклассовой классификации, когда цель состоит в том, чтобы рекомендовать «следующий» продукт на основании рейтингов или предыдущих покупок, называется *рекомендательной системой*. Как и любую другую задачу множественной классификации, задачу подбора рекомендаций можно решить с помощью стандартного подхода. Тем не менее для таких задач были созданы специальные типы моделей машинного обучения, и для реализации рекомендаций лучше использовать именно их. Рекомендательные системы также лучше подходят для решения задач таргетинга — поиска клиентов, которым понравится продукт или рекламное предложение.

Кластеризация

Если в наборе исходных данных вообще нет меток, к ним не получится применить методы машинного обучения с учителем. Зато можно попробовать выявить сложившиеся группы в данных; задачи этого вида называются *кластеризацией (группировкой)*. Например, клиентов можно разделить на кластеры (группы) в зависимости от характеризующих их признаков. Также можно воспользоваться службой Cloud Data Labeling Service и прибегнуть к помощи людей, осуществляющих маркировку данных, перед проведением обучения с учителем.

Неструктурированные данные

До сих пор мы исходили из того, что данные могут быть структурированными или полуструктурированными. Если какие-то входные признаки не структу-

рированы (например, являются изображениями или текстом на естественном языке), их можно обработать с помощью готовых моделей, таких как Cloud Vision API или Cloud Natural Language, и использовать полученные результаты в виде числовых или категориальных входных признаков. Например, с помощью Cloud Natural Language API можно определить ключевые сущности в электронных письмах клиентов или эмоциональную окраску их отзывов и использовать полученные сущности в качестве категориальных переменных, а оценки эмоциональной окраски — в качестве числовых признаков.

Преобразовать неструктурированные данные в структурированные можно с помощью строковых функций или машинного обучения. На практике часто используется прием под названием «*мешок слов*», суть которого состоит в том, чтобы разделить текстовое поле на отдельные слова и интерпретировать наличие/отсутствие отдельных слов как признаки. Например, из названия фильма «Шпион, который меня любил» («The Spy Who Loved Me») можно выделить два признака со значением True: `has_spy` и `has_love`. Все остальные признаки будут иметь значение False (слова «the», «Who» и «Me» лучше отбросить, потому что они слишком часто встречаются в тексте и не имеют большого значения для прогноза). Также в качестве признака можно использовать количество слов в названии (фильмы с большим количеством слов в названии, скорее всего, окажутся авторским кино и могут быть привлекательными для разных аудиторий).

Если не структурирована сама метка (например, требуется, чтобы модель генерировала идеальные ответы на вопросы клиентов, опираясь на набор архивных ответов), такие задачи относят к разряду задач генерирования естественного языка — они не поддерживаются в BigQuery.

Краткая сводка по типам моделей

В табл. 9.1 перечисляются типы задач машинного обучения. В следующем разделе мы обсудим типы моделей, поддерживаемые в BigQuery.

Таблица 9.1. Типы моделей машинного обучения и как они реализуются в BigQuery

Описание задачи	Тип задачи машинного обучения	Тип модели в BigQuery
Метки отсутствуют, данные нельзя обеспечить метками	Кластеризация	kmeans
Числовые метки	Регрессия	linear_reg dnn_regressor boosted_tree_regressor
Рекомендация продуктов пользователям	Рекомендации	matrix_factorization

Таблица 9.1 (окончание)

Описание задачи	Тип задачи машинного обучения	Тип модели в BigQuery
Выбор целевой аудитории для предложения продукта	Таргетинг клиентов	<code>matrix_factorization</code>
Метки имеют значения 1/0, True/False (две категории)	Бинарная классификация	<code>logistic_reg</code> <code>dnn_classifier</code> <code>boosted_tree_classifier</code>
Метки представлены фиксированным набором строк	Многоклассовая классификация	<code>logistic_reg</code> <code>dnn_classifier</code> <code>boosted_tree_classifier</code>
Входные признаки не структурированы	Классификация изображений, классификация текста, анализ эмоциональной окраски, извлечение сущностей	Используйте выходные данные Cloud Vision API или Cloud Natural Language API в качестве входных данных для любой из стандартных моделей BigQuery, перечисленных выше
Метки не структурированы	Ответы на вопросы, аннотирование текста, генерирование подписей к изображениям	Используйте продукты Cloud AutoML

Построение регрессионной модели

В примере построения регрессионной модели мы используем набор данных `london_bicycles`. Предположим, что у нас есть два типа велосипедов: тяжелые и надежные городские велосипеды и быстрые, но менее надежные шоссейные велосипеды. Для клиентов, совершающих длительные поездки, у нас должны быть в наличии шоссейные велосипеды, а для коротких поездок — обычные городские. Чтобы организовать правильное распределение велосипедов, необходимо спрогнозировать продолжительность поездок.

Выбор метки

Первым шагом к решению задачи машинного обучения является ее формулирование — определение признаков и меток для модели. Поскольку задача первой модели состоит в том, чтобы на основе набора архивных данных предсказать продолжительность поездки, меткой будет служить продолжительность поездки.

Но насколько верно мы определили цель задачи? Должна ли модель прогнозировать продолжительность каждой поездки или общую продолжительность

всех поездок для каждого пункта проката, например, в течение часа? Если мы выберем последний вариант, тогда меткой должна быть сумма продолжительностей всех поездок за определенный час. Допустим, по опыту ведения бизнеса мы знаем, что пункт проката, выдающий 1000 велосипедов на 20 минут каждый, должен выдавать городские велосипеды, тогда как пункт проката, выдающий ежедневно 100 велосипедов на 200 минут, должен выдавать шоссейные велосипеды. То есть прогнозирование общей продолжительности не поможет принять правильное решение, а вот предсказание продолжительности каждой отдельной поездки будет весьма кстати.

Другой вариант — оценка вероятности, что поездка продлится меньше 30 минут. В этом случае метка будет иметь два значения — True/False — в зависимости от длительности поездки (больше или меньше 30 минут). Это еще больше поможет бизнесу, потому что вероятность может определять относительную пропорцию городских и шоссейных велосипедов, имеющихся в каждом пункте проката.

На практике довольно часто приходится выбирать между несколькими целями. Иногда можно создать метку в виде взвешенной комбинации из целей и обучить единственную модель. Иногда полезнее обучить несколько моделей, по одной для каждой цели, и использовать разные модели в разных сценариях. А иногда лучше представить конечному пользователю результаты всех моделей и дать ему возможность выбора. Все зависит от того, чем вы занимаетесь.

В этом примере мы построим две модели: одну для прогнозирования продолжительности поездки, а другую для прогнозирования вероятности того, что поездка займет больше 30 минут. Затем мы дадим конечному пользователю возможность принять решение на основе двух прогнозов.

Выбор признаков в наборе данных

Если предположить, что продолжительность поездок зависит от пункта проката, дня недели и времени суток, эти параметры могут послужить нам входными признаками. Прежде чем продолжить и создать модель с этими тремя признаками, желательно убедиться, что эти факторы действительно влияют на метку.

Выбор и формирование признаков для модели машинного обучения называется *конструированием признаков* (feature engineering). Конструирование признаков часто является наиболее важным условием создания точных моделей машинного обучения, и этот шаг может гораздо сильнее повлиять на точность прогнозирования, чем выбор алгоритма или настройка гиперпараметров. Чтобы получить хороший набор признаков, необходимо глубоко понимать данные и предметную область. Часто на этом этапе проверяется множество гипотез; у вас есть идея относительно признака, вы проверяете ее обоснованность (влияние признака на метку), а затем добавляете этот признак в модель. Если эта идея не подтверждается, вы проверяете следующую.