

МАШИННОЕ



ОБУЧЕНИЕ



Применение в бизнесе

Д ж е й д К а р т е р

Джейд Картер

Машинное обучение

http://www.litres.ru/pages/biblio_book/?art=69356998

SelfPub; 2023

Аннотация

Книга представляет комплексное руководство по применению МО в сфере бизнеса. Автор исследует различные аспекты МО и его роль в современных бизнес-процессах, а также предлагают практические рекомендации по использованию этих технологий для достижения конкурентных преимуществ и улучшения результатов. В книге рассматриваются алгоритмы МО и объясняется, как они могут быть применены в различных сферах бизнеса, включая маркетинг, финансы, производство, здравоохранение и другие. Автор предлагает практические примеры и сценарии использования МО и как оно может быть внедрено в организациях. Особое внимание уделяется вопросам предобработки и анализу данных. Методы работы с Big Data и подходы к обработке неструктурированных данных. Этические и юридические аспекты МО в бизнесе, включая вопросы конфиденциальности и защиты данных. Книга полезна для менеджеров, аналитиков, предпринимателей и всех, кто заинтересован в использовании МО для оптимизации бизнес-процессов и принятия обоснованных решений.

Содержание

Глава 1: Введение в машинное обучение и его роль в бизнесе	6
Глава 2: Типы задач машинного обучения в бизнесе	28
Глава 3: Подготовка данных для машинного обучения	65
Конец ознакомительного фрагмента.	77

Джейд Картер

Машинное обучение

Список сокращений

1. МО – машинное обучение
2. ИИ – искусственный интеллект
3. СЗ – супервизированное обучение
4. БЗ – безнадзорное обучение
5. ПЗ – полузаданные обучение
6. НС – нейронная сеть
7. SVM – метод опорных векторов
8. RF – случайный лес
9. CNN – сверточная нейронная сеть
10. RNN – рекуррентная нейронная сеть
11. MLP – многослойный персептрон
12. SGD – стохастический градиентный спуск
13. NLP – обработка естественного языка
14. CV – компьютерное зрение
15. DL – глубокое обучение
16. ROI – возврат инвестиций
17. KPI – ключевые показатели эффективности
18. CRM – управление взаимоотношениями с клиентами
19. ERP – система планирования ресурсов предприятия

20. BI – бизнес-аналитика

Глава 1: Введение в машинное обучение и его роль в бизнесе

1.1. Основные понятия и термины в машинном обучении

Обучение с учителем – форма машинного обучения, где системе предоставляется обучающая выборка с входными данными и соответствующими выходными значениями.

Признаки – характеристики или свойства объектов, которые описывают данные.

Метки (выходные значения, целевые переменные) – значения, которые система должна предсказывать или классифицировать на основе входных данных.

Модель – математическая функция, которая принимает входные данные и выдает предсказания или классификации.

Обучение – процесс, в ходе которого модель настраивается на основе обучающей выборки для минимизации ошибки предсказания.

Тестирование – процесс оценки производительности модели на новых данных, не участвующих в обучении, с целью оценки ее обобщающей способности.

Переобучение – состояние модели, когда она становится

слишком сложной и настраивается на шум в данных, в результате чего ее способность обобщения страдает.

Недообучение – состояние модели, когда она слишком проста и не способна выявить сложные закономерности в данных, что приводит к низкой производительности на новых данных.

Гиперпараметры – параметры модели, которые задаются вручную перед началом обучения и влияют на ее поведение и производительность, например, скорость обучения, количество эпох и размер скрытых слоев в нейронной сети.

Алгоритмы обучения – методы и процедуры, используемые для обучения моделей на основе обучающих данных, например, линейная регрессия, метод опорных векторов (SVM), деревья решений, нейронные сети и другие.

Регуляризация – техника, используемая для предотвращения переобучения модели путем добавления штрафов или ограничений на значения параметров модели.

Кросс-валидация – метод оценки производительности модели, который заключается в разделении обучающей выборки на несколько подмножеств (фолдов) для обучения и тестирования модели, с последующим усреднением результатов.

Метрики оценки – числовые значения, используемые для измерения качества предсказаний модели, например, точность, полнота, F-мера, среднеквадратическая ошибка (MSE) и другие.

Разделение выборки – процесс разбиения общего набора данных на обучающую, тестовую и, иногда, валидационную выборки для обучения, тестирования и настройки модели соответственно.

Размер выборки – количество образцов данных, доступных для обучения модели.

Препроцессинг данных – этап подготовки данных перед обучением модели, включающий операции, такие как нормализация, масштабирование, заполнение пропущенных значений, кодирование категориальных признаков и другие.

Распределение данных – статистическая характеристика данных, которая описывает их вероятностные свойства, такие как среднее значение, дисперсия и форма распределения.

Ансамбли моделей – методы, которые объединяют предсказания нескольких моделей для получения более точного и устойчивого результата, например, бэггинг, случайный лес и градиентный бустинг.

Большие данные – наборы данных, которые характеризуются объемом, разнообразием и скоростью обновления, требующие специальных подходов и инструментов для их анализа и обработки.

Параметры модели – внутренние настраиваемые переменные, которые определяют ее поведение и способность предсказывать выходные значения. При обучении модели параметры настраиваются таким образом, чтобы минимизиро-

вать ошибку предсказания.

Функция потерь – математическая функция, которая измеряет расхождение между предсказанными и фактическими значениями модели. Цель обучения заключается в минимизации значения функции потерь.

Градиентный спуск – метод оптимизации, используемый для настройки параметров модели путем поиска оптимальных значений, исходя из градиента функции потерь. Градиентный спуск позволяет модели постепенно приближаться к минимуму функции потерь.

Регрессия – задача машинного обучения, которая связана с предсказанием непрерывных выходных значений на основе входных данных. Например, регрессионная модель может прогнозировать цену недвижимости на основе ее характеристик.

Классификация – задача машинного обучения, которая заключается в присвоении входным данным определенных категорий или классов. Классификационная модель может, например, определять, является ли электронное письмо спамом или не спамом.

Нейронные сети – модели машинного обучения, которые состоят из искусственных нейронов, объединенных в слои. Нейронные сети способны обрабатывать сложные входные данные и выявлять скрытые закономерности. Они широко используются в различных областях, таких как компьютерное зрение и естественный язык.

Сверточные нейронные сети – специализированный тип нейронных сетей, которые эффективно работают с входными данными в виде изображений. Они используют операцию свертки для извлечения локальных признаков из изображений и позволяют достигать высокой точности в задачах компьютерного зрения.

Рекуррентные нейронные сети – тип нейронных сетей, которые обладают памятью и могут обрабатывать последовательные данные, сохраняя информацию о предыдущих состояниях. Они часто применяются в задачах обработки естественного языка и временных рядов.

Безопасность и этика в машинном обучении – область, которая изучает вопросы связанные с надежностью, прозрачностью и справедливостью моделей машинного обучения. Включает в себя вопросы конфиденциальности данных, предвзятости моделей и этического использования искусственного интеллекта.

Андерсемплинг – метод сокращения преобладающего класса в несбалансированных данных путем удаления части образцов этого класса.

Оверсемплинг – метод увеличения меньшего класса в несбалансированных данных путем добавления дубликатов или синтетических образцов этого класса.

Автоэнкодеры – тип нейронных сетей, используемых для обучения представлений данных путем кодирования и декодирования входных сигналов. Они могут быть использованы

для извлечения скрытых признаков или снижения размерности данных.

Алгоритмы кластеризации – методы, используемые для разделения множества данных на группы или кластеры на основе их сходства. Примеры включают k-средних, иерархическую кластеризацию и DBSCAN.

Обратное распространение ошибки – алгоритм, используемый для обучения нейронных сетей путем вычисления и корректировки градиента функции потерь от выхода к входу сети.

Метод главных компонент (PCA) – метод снижения размерности данных путем преобразования их в новое пространство признаков, состоящее из линейных комбинаций исходных признаков с наибольшей дисперсией.

Рекомендательные системы – системы, используемые для предоставления рекомендаций пользователю на основе его предпочтений и поведения. Они широко применяются в электронной коммерции, музыкальных стриминговых сервисах и социальных сетях.

Генеративные модели – модели, которые могут генерировать новые данные, имитируя вероятностные распределения исходных данных. Примеры включают генеративные состязательные сети (GAN) и вариационные автоэнкодеры.

Понимание этих концепций является важным фундаментом для дальнейшего изучения и применения методов машинного обучения.

1.2. Преимущества и потенциал применения машинного обучения в бизнесе

В последние годы машинное обучение стало одной из самых обсуждаемых и востребованных областей в сфере бизнеса. Его способность анализировать данные, выявлять скрытые закономерности и делать предсказания делает его мощным инструментом для повышения эффективности и принятия обоснованных решений. В этой главе рассмотрим преимущества и потенциал применения машинного обучения в бизнесе.

1. Улучшение прогнозирования и планирования

Машинное обучение предоставляет бизнесу мощный инструмент для предсказания будущих событий и трендов на основе анализа больших объемов данных. Эта способность может быть особенно ценной для компаний, поскольку позволяет им получать ценную информацию, которая помогает принимать осознанные и стратегические решения.

Одной из ключевых преимуществ МО для бизнеса является его способность предсказывать спрос на товары и услуги. Алгоритмы машинного обучения могут анализировать исторические данные о покупках, предпочтениях клиентов, сезонных факторах и других факторах, чтобы определить вероятные тренды спроса в будущем. Это позволяет компаниям прогнозировать спрос и принимать меры заранее, чтобы

эффективно планировать производство, управлять запасами и оптимизировать бизнес-процессы.

Прогнозирование рыночных тенденций является еще одной сильной стороной машинного обучения в бизнесе. Алгоритмы машинного обучения могут анализировать данные о рынке, экономических показателях, конкурентной среде, социальных медиа и других источниках, чтобы выявить тенденции и понять, как они могут повлиять на бизнес. Это позволяет компаниям принимать основанные на фактах решения, адаптироваться к изменениям рынка и найти новые возможности для роста.

МО также играет важную роль в планировании производства и оптимизации цепей поставок. Алгоритмы машинного обучения могут анализировать данные о заказах, производственных мощностях, поставках и других факторах, чтобы оптимизировать процессы производства и распределение ресурсов. Это позволяет компаниям улучшить эффективность и гибкость производства, снизить затраты и улучшить обслуживание клиентов.

Благодаря алгоритмам машинного обучения, бизнес может принимать более точные и основанные на данных решения. Модели машинного обучения могут анализировать сложные взаимосвязи между различными переменными и выявлять скрытые паттерны, которые могут быть незаметны для человеческого анализа. Это помогает компаниям принимать обоснованные и обоснованные решения, основанные на

объективных фактах и статистических моделях.

2. Автоматизация и оптимизация бизнес-процессов

МО имеет потенциал автоматизировать рутинные задачи и процессы в бизнесе, что может привести к значительным выгодам. Автоматизация позволяет освободить время и ресурсы сотрудников, чтобы они могли сконцентрироваться на более стратегических и креативных задачах.

Одной из областей, где машинное обучение может быть применено для автоматизации, является клиентское обслуживание. Чат-боты, основанные на алгоритмах машинного обучения, могут быть использованы для автоматизации ответов на типовые вопросы и запросы клиентов. Они могут обрабатывать и анализировать текстовые данные, понимать намерения клиентов и предоставлять релевантные ответы. Это позволяет снизить нагрузку на сотрудников, освободить их время от рутинных запросов и улучшить общее качество обслуживания клиентов.

Другой пример автоматизации с помощью МО – системы распознавания речи. Они могут быть использованы для автоматической транскрипции аудио- или видеозаписей, распознавания команд голосового управления или анализа разговоров с клиентами. Это снижает необходимость в ручной обработке и анализе больших объемов аудио- или видеоданных и повышает эффективность работы сотрудников.

Оптимизация бизнес-процессов с помощью алгоритмов МО также позволяет более эффективно использовать ресур-

сы и сократить издержки. Например, алгоритмы МО могут быть применены для прогнозирования спроса на товары или услуги, что позволяет компаниям планировать закупки и производство более точно и эффективно. Также алгоритмы МО могут помочь в оптимизации логистических и поставочных цепочек, оптимальном планировании маршрутов доставки или управлении запасами.

МО имеет потенциал значительно улучшить автоматизацию рутинных задач и процессов в бизнесе. Это позволяет более эффективно использовать ресурсы, сократить издержки и освободить время для выполнения более важных и стратегических задач.

3. Улучшение клиентского опыта и персонализация

МО играет важную роль в понимании предпочтений и поведения клиентов в бизнесе. Анализ больших объемов данных с применением алгоритмов МО позволяет выявлять скрытые паттерны и тренды, которые могут указывать на предпочтения и интересы клиентов.

Алгоритмы рекомендаций, основанные на МО, способны анализировать исторические данные о покупках, предпочтениях, поведении и интересах клиентов. Они создают уникальные профили клиентов и используют эти данные для предложения персонализированных товаров и услуг. Например, на основе предыдущих покупок клиентов и сходных паттернов поведения, система рекомендаций может предложить товары, которые могут заинтересовать конкретного

клиента.

Это имеет большое значение для бизнеса, поскольку персонализированные предложения повышают удовлетворенность клиентов. Когда клиенты получают рекомендации, которые соответствуют их предпочтениям и потребностям, они чувствуются более важными и учтенными. Это может привести к увеличению частоты покупок, повышению лояльности клиентов и росту прибыли.

Более того, МО позволяет бизнесу применять индивидуальные рекомендации, учитывая контекст и ситуацию клиента. Например, алгоритмы машинного обучения могут учитывать данные о местоположении, времени суток, погодных условиях и других факторах, которые могут влиять на предпочтения клиента. Это позволяет бизнесу предлагать более релевантные и актуальные предложения, улучшая впечатление клиентов и повышая шансы на успешное завершение сделки.

МО помогает бизнесу лучше понимать клиентов и предлагать более персонализированные предложения и рекомендации. Это способствует повышению удовлетворенности клиентов, росту лояльности и увеличению прибыли компании.

4. Обнаружение мошенничества и анализ рисков

МО имеет значительный потенциал для выявления аномалий и обнаружения потенциальных случаев мошенничества в бизнесе. Алгоритмы машинного обучения могут обра-

батывать и анализировать огромные объемы данных, искать необычные паттерны и сигналы, которые могут указывать на наличие мошеннической активности.

Это особенно важно для финансовых учреждений и компаний, где безопасность и защита данных являются приоритетными задачами. МО может быть применено для обнаружения мошеннических транзакций, фальшивых идентификационных документов, несанкционированного доступа к системам и других видов мошенничества.

Алгоритмы МО могут быть обучены на основе исторических данных о мошеннической активности, что позволяет им распознавать подозрительные ситуации и сравнивать текущие события с ранее известными шаблонами мошенничества. Например, модель МО может выявить необычные транзакции с необычно высокими суммами, необычные паттерны поведения клиентов или несоответствие типичным сценариям использования продукта или услуги. При обнаружении подозрительных сигналов система может предпринять соответствующие меры, например, заблокировать транзакцию или оповещать службу безопасности для проведения дополнительной проверки.

Это позволяет бизнесу более эффективно бороться с мошенничеством, защищать своих клиентов и себя от потенциальных угроз. В результате, финансовые учреждения и компании могут сэкономить значительные суммы денег, предотвратив финансовые потери, и поддерживать свою репутацию.

цию, обеспечивая безопасность и надежность своих услуг.

Однако, важно отметить, что МО не является идеальным и может сталкиваться с ограничениями и вызовами при обнаружении мошенничества. Некоторые виды мошенничества могут быть сложными и изменчивыми, и могут быть неизвестны для моделей машинного обучения, обученных на исторических данных. Кроме того, существует риск ложноположительных и ложноотрицательных результатов, когда модель неправильно классифицирует транзакцию как мошенническую или не замечает реальную мошенническую активность.

Поэтому важно комбинировать применение алгоритмов МО с другими методами и инструментами для обеспечения безопасности бизнеса. Это может включать мониторинг и аудит систем, вовлечение специалистов в области безопасности, разработку политик и процедур для обработки потенциальных случаев мошенничества.

МО имеет большой потенциал для выявления аномалий и обнаружения мошенничества в бизнесе. Оно помогает бизнесу защищать своих клиентов, предотвращать финансовые потери и поддерживать высокий уровень безопасности и доверия. Однако, необходимо учитывать ограничения и вызовы при использовании машинного обучения и принимать дополнительные меры для обеспечения безопасности и эффективности системы.

5. Инновации и новые возможности

МО предоставляет бизнесу уникальные возможности исследования и инновации, открывая новые горизонты в анализе данных и принятии решений. Алгоритмы машинного обучения способны обрабатывать и анализировать огромные объемы данных, выявлять скрытые паттерны и взаимосвязи, которые могут остаться незамеченными человеческим взглядом.

Анализ данных с помощью МО может привести к открытию новых знаний и неожиданных выводов. Например, модель МО может обнаружить скрытые корреляции между различными переменными, выявить факторы, влияющие на спрос на продукты или предсказать тенденции и тренды на рынке. Это позволяет бизнесу принимать более информированные и основанные на данных решения.

Благодаря МО, бизнес может разрабатывать новые продукты и услуги, оптимизировать бизнес-модели и создавать инновационные решения. Например, на основе анализа данных о потребностях клиентов, предпочтениях и поведении, бизнес может разработать более персонализированные продукты и предлагать индивидуальные рекомендации. Это улучшает опыт клиентов, повышает их удовлетворенность и способствует повторным покупкам.

Кроме того, МО может помочь бизнесу открыть новые рыночные сегменты и идентифицировать потенциально прибыльные возможности. Алгоритмы машинного обучения могут анализировать данные о поведении клиентов, соци-

альных тенденциях и экономических факторах, чтобы выявить нишевые сегменты рынка или потенциальные рыночные разрывы. Это позволяет бизнесу адаптироваться к изменяющейся среде и идентифицировать новые возможности для роста и развития.

Таким образом, МО предоставляет бизнесу новые возможности для исследования данных, инноваций и развития. Анализ данных с помощью алгоритмов машинного обучения помогает выявить скрытые паттерны, прогнозировать тренды и создавать более эффективные стратегии. Это открывает двери для разработки новых продуктов и услуг, оптимизации бизнес-процессов и открытия новых рыночных возможностей.

В заключение, МО имеет огромный потенциал для применения в бизнесе. Оно способно улучшить прогнозирование, оптимизировать бизнес-процессы, повысить качество обслуживания клиентов, обнаружить мошенничество и создать новые возможности для инноваций. Понимание и использование этих преимуществ позволяют бизнесу оставаться конкурентоспособным в современной высокотехнологичной среде.

1.3. Ограничения и вызовы использования машинного обучения в бизнесе

В ходе использования МО в бизнесе, мы сталкиваемся с

определенными ограничениями и вызовами.

Одним из ключевых факторов, которые необходимо учитывать при использовании машинного обучения в бизнесе, является качество данных. Качество данных оказывает прямое влияние на точность и достоверность результатов моделей машинного обучения.

Для того чтобы модели МО могли предсказывать и принимать решения на основе данных, эти данные должны быть высокого качества. Качество данных включает в себя такие аспекты, как полнота, точность и отсутствие шума. Неполные данные могут содержать пропущенные значения или отсутствующие фрагменты, что может исказить общую картину и снизить эффективность моделей.

Точность данных также является важным аспектом. Если данные содержат ошибки или неточности, то модели МО могут давать неверные предсказания или рекомендации. Например, если данные о клиентах содержат неточную информацию о их предпочтениях или покупках, то модель может сделать неверные выводы о предпочтениях и поведении клиентов.

Шум в данных представляет собой случайные или нежелательные вариации, которые могут вносить дополнительные искажения в процесс обучения моделей. Наличие шума может привести к некорректным или несостоятельным выводам. Например, если данные о погоде содержат случайные выбросы или ошибки измерений, то модель, обученная на

таких данных, может давать непредсказуемые результаты.

Для достижения высокого качества данных, необходимо уделить должное внимание процессу сбора, обработки и очистки данных. Это может включать автоматизацию процессов, применение алгоритмов обработки данных, удаление выбросов и ошибок, а также проверку и верификацию данных.

Однако, несмотря на все усилия, полностью избавиться от проблем с качеством данных невозможно. Важно иметь реалистические ожидания относительно качества данных и принять меры для минимизации влияния возможных недочетов. Это может включать мониторинг качества данных, использование алгоритмов, устойчивых к шуму, и внесение корректировок в модели, если данные изменяются или ухудшаются со временем.

Другим вызовом, связанным с использованием моделей МО в бизнесе, является их интерпретируемость. Некоторые типы моделей, особенно сложные нейронные сети, могут быть непрозрачными в своих принятиях решений. Это означает, что для людей может быть сложно объяснить, почему модель приняла ту или иную решающую ставку.

Интерпретируемость моделей играет важную роль в бизнесе, особенно когда принимаются важные решения, такие как предсказания рыночных трендов, определение стратегии продаж или принятие инвестиционных решений. Компании и организации могут столкнуться с вызовом в том, что требу-

ется объяснить, почему модель сделала определенное предсказание или рекомендацию.

Непрозрачность моделей может вызывать сомнения и недоверие в их результаты. Бизнес-лидеры и заинтересованные стороны могут испытывать необходимость в понимании причин, которые привели к определенным решениям. В некоторых отраслях, таких как финансовый сектор или здравоохранение, требуется обоснование и объяснение решений, сделанных моделью.

Для решения этого вызова и повышения интерпретируемости моделей МО, проводится активное исследование в области алгоритмов "черного ящика" и методов объяснения моделей. Некоторые подходы включают визуализацию важных признаков, анализ вклада каждого признака в принятие решения, использование методов "линейной аппроксимации" для построения понятных моделей и др.

Однако, эти дополнительные усилия по объяснению моделей могут потребовать дополнительных ресурсов и времени. Компании должны внимательно рассмотреть баланс между точностью и интерпретируемостью моделей, и определить, насколько важно иметь понятные объяснения за счет некоторого снижения точности предсказаний.

Вопрос интерпретируемости моделей МО остается актуальным в бизнесе. Балансировка между сложностью модели и ее понятностью является одним из вызовов, с которыми компании сталкиваются при использовании машинного обу-

чения в своей деятельности.

Еще одним ограничением, с которым сталкиваются компании при использовании машинного обучения, является нехватка экспертизы и ресурсов. Внедрение МО требует глубоких знаний и опыта в области алгоритмов, моделей и технологий.

Компании, не обладающие достаточным количеством квалифицированных специалистов, могут столкнуться с ограничениями при внедрении и использовании МО. Необходимо иметь специалистов, которые обладают навыками в области обработки данных, анализа, выбора и оптимизации моделей, а также умеющих эффективно работать с соответствующими инструментами и программными средствами.

Кроме нехватки экспертизы, использование МО может требовать значительных ресурсов. Некоторые модели машинного обучения требуют высокопроизводительного оборудования и вычислительных мощностей для обучения и развертывания моделей. Это может быть финансово затратным для многих компаний, особенно для малых и средних предприятий.

Для преодоления этого ограничения компании могут искать способы повышения уровня экспертизы своих сотрудников через обучение и повышение квалификации. Это может включать обучение внутреннего персонала, привлечение внешних консультантов или партнерство с университетами и исследовательскими организациями.

Для снижения финансовой нагрузки, связанной с использованием МО, компании могут рассмотреть возможность использования облачных сервисов и платформ, которые предоставляют вычислительные ресурсы на арендной основе. Это позволяет снизить затраты на инфраструктуру и обеспечить гибкость в использовании вычислительных ресурсов в зависимости от потребностей.

Однако, несмотря на ограничения, недостаток экспертизы и ресурсов не должен отпугивать компании от применения МО в бизнесе. Существуют различные способы преодоления этих вызовов, и с течением времени и развитием технологий, доступность и доступность ресурсов и экспертизы в области машинного обучения продолжают улучшаться.

Безопасность и этика являются критическими аспектами, которые необходимо учитывать при использовании МО в бизнесе. Одним из важных вопросов является обеспечение безопасности данных. Некорректная обработка и использование данных может привести к нарушению конфиденциальности и приватности клиентов. Важно обеспечивать адекватные меры защиты данных, чтобы предотвратить несанкционированный доступ, утечку информации или злоупотребление данными. Это может включать применение криптографических методов, контроль доступа, анонимизацию данных и обеспечение соответствия нормам и правилам обработки персональных данных.

Кроме того, модели МО могут быть предвзятыми и

несправедливыми. Это может произойти, если данные, на которых модель обучалась, содержали предвзятость или нерепрезентативность. Например, если модель обучалась на данных, в которых преобладали определенные группы, это может привести к систематическому неравенству и несправедливому воздействию на другие группы. Важно учитывать эти этические аспекты и принимать меры для минимизации предвзятости моделей, такие как балансировка классов или справедливая выборка данных.

Другим аспектом этики является вопрос о социальной ответственности. Модели МО могут иметь значительное воздействие на общество и людей. Важно учитывать потенциальные негативные последствия и воздействие, которое модели могут оказывать на различные группы людей или общество в целом. Это может включать вопросы дискриминации, неравенства, прозрачности и объяснимости принимаемых моделью решений. Компании должны стремиться к разработке и использованию моделей, которые учитывают эти этические аспекты и способствуют положительному воздействию на общество.

В свете этих вопросов безопасности и этики, компании должны принимать соответствующие меры для защиты данных, обеспечения справедливости моделей и социальной ответственности. Это может включать проведение оценки воздействия на приватность, этический аудит моделей, установление принципов и политик в области безопасности и этики,

а также обучение сотрудников основным принципам и нормам в использовании МО.

Несмотря на эти ограничения и вызовы, машинное обучение все равно предоставляет бизнесу значительные преимущества и потенциал для роста и развития. Понимание и учет этих ограничений помогает бизнесам принимать обоснованные решения и разрабатывать соответствующие стратегии для успешного внедрения машинного обучения в своей деятельности.

Глава 2: Типы задач машинного обучения в бизнесе

2.1. Классификация и предсказание

В машинном обучении классификация и предсказание являются одними из основных задач. Классификация относится к процессу разделения данных на заранее определенные категории или классы на основе их характеристик. Это позволяет модели машинного обучения классифицировать новые данные, определяя, к какому классу они относятся. Примером классификации может быть определение электронного письма как спама или не спама, или определение изображения как кошки или собаки.

Предсказание, с другой стороны, связано с использованием модели машинного обучения для предсказания значений или результатов на основе имеющихся данных. Модель обучается на исторических данных и затем используется для предсказания будущих значений. Например, модель машинного обучения может быть обучена на данных о продажах и использована для предсказания продаж на следующий месяц или год.

Классификация и предсказание имеют широкий спектр

применений в бизнесе. Они могут помочь в определении спроса на товары и услуги, выявлении потенциальных клиентов, прогнозировании рыночных тенденций и анализе рисков. Например, на основе данных о клиентах, модель машинного обучения может классифицировать их по уровню лояльности или предсказывать вероятность их оттока. Это позволяет бизнесу принимать более информированные решения о маркетинговых стратегиях, управлении клиентским опытом и удержании клиентов.

Классификация и предсказание также могут быть использованы для обнаружения аномалий и предотвращения мошенничества. Например, модель машинного обучения может классифицировать финансовые транзакции как нормальные или подозрительные на основе их характеристик, помогая бизнесу выявить потенциальные случаи мошенничества.

Давайте рассмотрим пример использования классификации и предсказания на наборе данных о банковских клиентах для определения их вероятности дефолта. Предположим, что у нас есть набор данных, содержащий информацию о клиентах банка, такую как возраст, пол, доход, семейное положение, кредитная история и другие параметры.

Мы можем использовать модель МО, например, логистическую регрессию, для классификации клиентов на два класса: дефолтные и недефолтные. Модель будет обучаться на исторических данных, где для каждого клиента известно,

произошел ли дефолт или нет. Затем, используя эту модель, мы можем предсказывать вероятность дефолта для новых клиентов на основе их характеристик.

Такой анализ может быть полезен для банков в принятии решений о выдаче кредитов. Например, если модель предсказывает высокую вероятность дефолта для определенного клиента, банк может принять решение о отказе в выдаче кредита или установить более строгие условия. Это позволяет снизить риски и улучшить управление кредитным портфелем.

Этот пример демонстрирует, как классификация и предсказание на основе данных могут быть использованы для принятия решений в банковской сфере, анализе рисков и определении оптимальных стратегий предоставления услуг клиентам.

Пример программы на языке Python, использующей библиотеку `scikit-learn` для классификации с помощью модели логистической регрессии:

```
```python
Импортирование необходимых библиотек
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

Загрузка набора данных

Предположим, что у нас есть CSV-файл с данными о
банковских клиентах
```

# Содержащий столбцы: возраст, пол, доход, семейное положение, кредитная история и целевая переменная (дефолт/недефолт)

```
data = pd.read_csv("bank_clients.csv")
```

# Разделение данных на признаки (X) и целевую переменную (y)

```
X = data.drop("target", axis=1)
```

```
y = data["target"]
```

# Разделение данных на тренировочный и тестовый наборы

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)
```

# Создание модели логистической регрессии

```
model = LogisticRegression()
```

# Обучение модели на тренировочном наборе данных

```
model.fit(X_train, y_train)
```

# Прогнозирование классов для тестового набора данных

```
y_pred = model.predict(X_test)
```

# Вычисление точности модели

```
accuracy = accuracy_score(y_test, y_pred)
```

```
print("Точность модели: {:.2f}".format(accuracy))
```

```
...
```

В этом примере мы используем модель логистической регрессии для классификации банковских клиентов на дефолтные и недефолтные. Мы загружаем данные из CSV-файла, разделяем их на признаки и целевую переменную, а

затем разделяем их на тренировочный и тестовый наборы данных. Модель логистической регрессии обучается на тренировочном наборе, а затем используется для предсказания классов для тестового набора. Наконец, мы вычисляем точность модели с помощью метрики `assurasy_score`.

Обратите внимание, что этот пример является общим и требует наличия данных в соответствующем формате и установленных библиотек `scikit-learn` и `pandas` для работы.

Логистическая регрессия (Logistic Regression) является одним из методов бинарной классификации в машинном обучении. Она используется для предсказания вероятности принадлежности объекта к определенному классу.

Основная идея логистической регрессии состоит в том, чтобы использовать логистическую функцию (также известную как сигмоидная функция) для преобразования линейной комбинации признаков объекта в вероятность принадлежности к классу. Формула логистической регрессии выглядит следующим образом:

$$p(y=1|x) = \text{sigmoid}(w^T * x + b)$$

где:

- $p(y=1|x)$  представляет собой вероятность принадлежности объекта к классу 1 при условии значения признаков  $x$ ,
- $w$  – вектор весов, соответствующий признакам,
- $b$  – смещение (bias),
- `sigmoid` – логистическая функция, определенная как  $\text{sigmoid}(z) = 1 / (1 + \exp(-z))$ .



Для обучения модели логистической регрессии используется метод максимального правдоподобия, который позволяет настроить веса и смещение модели таким образом, чтобы максимизировать вероятность наблюдаемых данных.

После обучения модели логистической регрессии, для новых объектов можно использовать полученные веса для вычисления их вероятности принадлежности к классу 1. Затем можно применить пороговое значение для принятия решения о классификации объекта.

Логистическая регрессия является одним из наиболее широко используемых методов классификации в различных областях, включая медицину, финансы, маркетинг и другие. Ее популярность объясняется несколькими причинами.

Во-первых, логистическая регрессия отличается простотой в реализации и интерпретации. Модель основана на линейной комбинации признаков, что делает ее относительно простой для понимания. При этом полученные веса модели можно интерпретировать в контексте важности каждого признака для классификации. Это позволяет исследователям и экспертам в соответствующих областях использовать результаты модели для принятия решений и проведения анализа данных.

Во-вторых, логистическая регрессия обладает хорошей способностью к обобщению. Даже при наличии большого количества признаков она способна эффективно работать с относительно небольшим объемом данных. Это делает ее при-

менимой в случаях, когда доступные данные ограничены.

В-третьих, логистическая регрессия позволяет моделировать вероятности принадлежности к классу, а не только делать бинарные предсказания. Это может быть полезно в задачах, где важно не только определить класс объекта, но и оценить уверенность в этом предсказании.

## **2.2. Кластеризация и сегментация**

Кластеризация и сегментация – это важные методы анализа данных, которые позволяют группировать объекты в подобные кластеры или сегменты на основе их схожести или общих характеристик. Эти методы имеют широкое применение в различных областях, включая маркетинг, социальные исследования, медицину, географический анализ и многие другие.

Кластеризация – это процесс разделения объектов на группы (кластеры) таким образом, чтобы объекты внутри одного кластера были более схожи между собой, чем с объектами из других кластеров. Кластеризация может быть использована для выявления скрытых паттернов, структуры или типов объектов в данных. Например, в маркетинге кластеризация может помочь определить группы потребителей с общими предпочтениями или поведением, что позволит создать более эффективные стратегии маркетинга для каждой группы.

Сегментация – это процесс разделения группы объектов на более мелкие сегменты на основе их характеристик или поведения. Сегментация позволяет более детально изучать каждую группу и разрабатывать персонализированные стратегии для каждого сегмента. Например, в медицине сегментация пациентов может помочь выделить подгруппы с определенными медицинскими характеристиками или рисками заболеваний, что позволит проводить более точные и целевые лечебные мероприятия.

Кластеризация и сегментация основаны на алгоритмах машинного обучения, которые автоматически определяют схожесть или различия между объектами и формируют кластеры или сегменты. Эти алгоритмы могут использовать различные подходы, такие как методы иерархической кластеризации, методы на основе плотности, методы разделения, а также комбинации этих методов.

Рассмотрим пример кода для кластеризации данных в банковской сфере с использованием метода К-средних (K-means) в языке программирования Python:

```
```python
# Импорт необходимых библиотек
import pandas as pd
import numpy as np
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt
# Загрузка данных
```

```
data = pd.read_csv("bank_data.csv") # Предположим, у нас
есть файл с данными о клиентах банка
# Подготовка данных
X = data[['Age', 'Income']] # Выбираем признаки, по кото-
рым будем проводить кластеризацию
# Масштабирование данных
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
# Определение оптимального числа кластеров
inertia = []
for k in range(1, 10):
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(X_scaled)
    inertia.append(kmeans.inertia_)
# Визуализация графика локтя
plt.plot(range(1, 10), inertia, marker='o')
plt.xlabel('Number of clusters')
plt.ylabel('Inertia')
plt.title('Elbow Method')
plt.show()
# Выбор оптимального числа кластеров
k = 3 # По графику локтя видим, что оптимальное число
кластеров равно 3
# Применение метода К-средних
kmeans = KMeans(n_clusters=k, random_state=42)
```

```
kmeans.fit(X_scaled)
# Добавление меток кластеров в данные
data['Cluster'] = kmeans.labels_
# Вывод результатов
for cluster in range(k):
    cluster_data = data[data['Cluster'] == cluster]
    print(f"Cluster {cluster + 1}:\n{cluster_data.describe()}\n")
'''
```

Описание кода:

1. Импортируем необходимые библиотеки, такие как pandas для работы с данными, numpy для математических операций, sklearn для использования алгоритма К-средних и matplotlib для визуализации.
2. Загружаем данные из файла "bank_data.csv". Предполагается, что у нас есть файл с данными о клиентах банка, включающими возраст (Age), доход (Income) и другие признаки.
3. Выбираем признаки (Age и Income) для проведения кластеризации и создаем новый DataFrame X.
4. Масштабируем данные с помощью стандартизации с помощью объекта StandardScaler.
5. Определяем оптимальное число кластеров с помощью метода локтя (Elbow Method) и визуализируем график.
6. Выбираем оптимальное число кластеров (в данном случае равно 3).
7. Применяем метод К-средних с выбранным числом кла-

стеров.

8. Добавляем метки кластеров в исходные данные.

9. Выводим описательную статистику для каждого кластера.

Примечание: В приведенном коде предполагается, что у вас есть файл "bank_data.csv" с соответствующими данными о клиентах банка.

Метод К-средних (K-means) – это один из наиболее популярных алгоритмов кластеризации в машинном обучении. Он используется для разделения набора данных на заданное число кластеров.

Процесс работы метода К-средних выглядит следующим образом:

1. Определение числа кластеров (K): Сначала необходимо определить, сколько кластеров требуется создать. Это может быть заранее известное число или выбор на основе анализа данных и целей задачи.

2. Инициализация центроидов: Центроиды представляют собой точки в пространстве данных, которые инициализируются случайным образом или на основе предварительных оценок. Их количество соответствует числу кластеров K.

3. Присвоение точек к кластерам: Каждая точка данных присваивается к ближайшему центроиду на основе некоторой меры расстояния, чаще всего используется Евклидово расстояние.

4. Пересчет центроидов: После присвоения всех точек

кластерам пересчитываются новые центроиды. Это делается путем вычисления среднего значения координат точек в каждом кластере.

5. Повторение шагов 3 и 4: Процессы присвоения точек к кластерам и пересчета центроидов повторяются до тех пор, пока не будет достигнуто определенное условие остановки. Обычно это ограничение числа итераций или малая изменчивость центроидов.

6. Вывод результатов: По завершении алгоритма получаем набор кластеров, где каждая точка данных относится к определенному кластеру.

Формула, используемая в методе К-средних для определения принадлежности точки кластеру, выглядит следующим образом:

$$d(x, c) = \sqrt{(x_1 - c_1)^2 + (x_2 - c_2)^2 + \dots + (x_n - c_n)^2}$$

где:

– $d(x, c)$ представляет собой расстояние между точкой данных x и центроидом c ,

– x_1, x_2, \dots, x_n представляют координаты точки данных x ,

– c_1, c_2, \dots, c_n представляют координаты центроида c .

Формула использует Евклидово расстояние для вычисления расстояния между точкой данных и центроидом. Она измеряет разницу между каждой координатой точки данных и соответствующей координатой центроида, затем суммирует квадраты этих разностей и извлекает квадратный корень из суммы.

Это расстояние помогает определить, к какому кластеру должна быть отнесена точка данных. Чем ближе точка к центроиду, тем меньше значение расстояния, и она будет отнесена к этому кластеру.

Метод К-средних использует эту формулу для вычисления расстояния между каждой точкой данных и всеми центроидами, а затем выбирает ближайший центроид для каждой точки данных в качестве принадлежности к кластеру.

Метод К-средних является итеративным алгоритмом, который стремится минимизировать сумму квадратов расстояний между точками данных и центроидами. Он обладает простотой реализации и хорошей масштабируемостью, что делает его популярным методом для кластеризации данных в различных областях, включая бизнес, науку, медицину и другие.

Рассмотрим пример кода сегментации клиентов в банковской сфере с использованием метода К-средних (K-means). Этот метод может помочь выявить группы клиентов с общими характеристиками и поведением, что позволит банку адаптировать свои продукты и услуги под каждую группу более эффективно.

```
```python
import pandas as pd
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
Загрузка данных о клиентах банка
```



```
data = pd.read_csv('customer_data.csv')
Предобработка данных: масштабирование числовых признаков
scaler = StandardScaler()
scaled_data = scaler.fit_transform(data[['Age', 'Income', 'Balance']])
Определение количества кластеров
k = 3
Создание и обучение модели К-средних
kmeans = KMeans(n_clusters=k, random_state=42)
kmeans.fit(scaled_data)
Получение меток кластеров для каждого клиента
cluster_labels = kmeans.labels_
Добавление меток кластеров в исходные данные
data['Cluster'] = cluster_labels
Вывод результатов сегментации
for cluster in range(k):
 cluster_data = data[data['Cluster'] == cluster]
 print(f"Cluster {cluster}:")
 print(cluster_data.describe())
 print("\n")
Описание каждого кластера:
– Можно проанализировать средний возраст, доход и баланс по каждому кластеру
– Определить основные характеристики и поведение клиентов в каждом кластере
```

В данном примере мы используем библиотеки `pandas` и `scikit-learn` для загрузки данных о клиентах банка, предобработки данных и применения метода К-средних. Сначала данные подвергаются масштабированию с помощью `StandardScaler`, чтобы привести числовые признаки к одному масштабу.

Затем мы задаем количество кластеров (в данном случае  $k = 3$ ) и создаем экземпляр модели `KMeans`. Обучение модели происходит методом `fit`, где модель вычисляет центроиды кластеров, чтобы минимизировать сумму квадратов расстояний до точек данных внутри каждого кластера.

Полученные метки кластеров добавляются в исходные данные. Мы выводим описание каждого кластера, анализируя средние значения возраста, дохода и баланса для клиентов в каждом кластере. Это позволяет нам понять основные характеристики и поведение клиентов в каждой группе.

Используя результаты сегментации, банк может адаптировать свою стратегию продаж, маркетинга и обслуживания для каждого кластера клиентов, что поможет улучшить удовлетворенность клиентов и повысить эффективность работы банка.

## 2.3. Регрессия и прогнозирование

Регрессия и прогнозирование являются важными инстру-

ментами в области машинного обучения и анализа данных. Они позволяют бизнесу строить математические модели, которые могут предсказывать значения зависимой переменной на основе входных данных и обученных параметров модели. Это полезно для прогнозирования будущих событий, трендов и результатов на основе имеющихся данных.

Одним из наиболее распространенных методов регрессии является линейная регрессия. В линейной регрессии строится линейная модель, которая аппроксимирует зависимость между независимыми переменными и зависимой переменной. Модель представляет собой уравнение прямой линии, которая наилучшим образом соответствует данным. При помощи этой модели можно делать прогнозы и анализировать влияние различных факторов на зависимую переменную.

В случае, когда зависимая переменная является категориальной или дискретной, используется логистическая регрессия. Логистическая регрессия позволяет предсказывать вероятность отнесения наблюдения к определенному классу или категории. Модель использует логистическую функцию для преобразования линейной комбинации независимых переменных в вероятность.

Для регрессии и прогнозирования необходимо иметь набор данных, включающий значения зависимой переменной и соответствующие значения независимых переменных. Эти данные используются для обучения модели, то есть оценки параметров модели на основе имеющихся данных. Затем мо-

дель может быть использована для прогнозирования значений зависимой переменной для новых наблюдений или для анализа и интерпретации влияния независимых переменных на зависимую переменную.

Применение регрессии и прогнозирования в бизнесе может быть разнообразным. Например, в финансовой сфере регрессия может использоваться для прогнозирования цен акций или доходности инвестиций на основе исторических данных. В маркетинге регрессия может помочь в определении факторов, влияющих на продажи или клиентскую активность. В медицине регрессия может быть применена для прогнозирования заболеваемости или оценки влияния факторов на здоровье пациентов.

Оценка качества модели регрессии и прогнозирования является важным шагом в анализе данных и принятии решений в бизнесе. Различные метрики используются для оценки точности модели и ее способности обобщаться на новые данные. Ниже рассмотрим основные метрики, которые применяются в регрессии и прогнозировании.

1. Среднеквадратичная ошибка (Mean Squared Error, MSE): Это одна из наиболее распространенных метрик оценки качества модели регрессии. Среднеквадратичная ошибка измеряет среднее квадратичное отклонение между предсказанными значениями модели и истинными значениями зависимой переменной. Чем меньше значение MSE, тем ближе предсказания модели к реальным значениям. Формула для

расчета MSE:

$$\text{MSE} = (1/n) * \sum (y - \hat{y})^2,$$

где  $n$  – количество наблюдений,  $y$  – истинное значение зависимой переменной,  $\hat{y}$  – предсказанное значение зависимой переменной.

2. Коэффициент детерминации (R-squared): Эта метрика оценивает, насколько хорошо модель соответствует данным. Коэффициент детерминации показывает долю дисперсии зависимой переменной, которая объясняется моделью. Значение коэффициента детерминации находится в диапазоне от 0 до 1, где 0 означает, что модель не объясняет вариацию данных, а 1 означает, что модель идеально соответствует данным. Формула для расчета коэффициента детерминации:

$$R^2 = 1 - (\text{SSR} / \text{SST}),$$

где SSR – сумма квадратов остатков, SST – общая сумма квадратов отклонений от среднего.

3. Средняя абсолютная ошибка (Mean Absolute Error, MAE): Эта метрика измеряет среднее абсолютное отклонение между предсказанными значениями модели и истинными значениями зависимой переменной. Она является более устойчивой к выбросам, чем среднеквадратичная ошибка. Формула для расчета MAE:

$$\text{MAE} = (1/n) * \sum |y - \hat{y}|.$$

4. Корень из среднеквадратичной ошибки (Root Mean Squared Error, RMSE): Эта метрика представляет собой квадратный корень из средне-

квадратичной ошибки и используется для измерения среднего отклонения предсказанных значений от реальных значений. RMSE также измеряется в тех же единицах, что и зависимая переменная, что облегчает интерпретацию. Формула для расчета RMSE:

$$\text{RMSE} = \sqrt{\text{MSE}}.$$

Кроме этих основных метрик, существуют и другие метрики оценки качества модели регрессии, такие как коэффициенты корреляции, коэффициенты эффективности и другие, которые могут быть применены в зависимости от конкретной задачи и требований бизнеса.

Важно выбирать подходящую метрику в соответствии с целями анализа и спецификой данных, чтобы получить объективную оценку качества модели регрессии и прогнозирования.

При выборе подходящей метрики для оценки качества модели регрессии и прогнозирования следует учитывать следующие факторы:

1. Цель анализа: Определите, какую информацию вы хотите получить из модели и какие вопросы вы хотите на них ответить. Например, если вам важно измерить точность предсказания, то среднеквадратичная ошибка (MSE) или корень из среднеквадратичной ошибки (RMSE) могут быть подходящими метриками. Если ваша цель заключается в понимании объясняющей способности модели, то коэффициент детерминации (R-squared) может быть полезной метри-

кой.

2. Специфика данных: Рассмотрите особенности ваших данных, такие как наличие выбросов, несбалансированность классов или другие аномалии. Некоторые метрики, такие как среднеквадратичная ошибка (MSE), могут быть чувствительны к выбросам, в то время как средняя абсолютная ошибка (MAE) более устойчива к ним. Также учтите, что некоторые метрики могут быть предназначены для специфических типов данных или задач, например, метрики оценки точности классификации.

3. Бизнес-контекст: Изучите требования вашего бизнеса и применение модели. Какие критерии важны для вашей организации? Например, если вы работаете в области финансов, то точность предсказаний может быть особенно важной. Если вы прогнозируете спрос на товары, то средняя абсолютная ошибка (MAE) может быть полезной для измерения ошибки в денежных единицах.

4. Сравнение моделей: Если у вас есть несколько моделей, которые вы хотите сравнить, убедитесь, что выбранная метрика позволяет справедливо оценить их производительность. Некоторые метрики могут быть более чувствительны к определенным типам моделей или данным.

В идеале, выбор метрики должен быть основан на комбинации этих факторов и отражать конкретные цели и требования вашей задачи. Важно также понимать интерпретацию выбранной метрики и уметь объяснить ее значение заказчи-

кам.

Регрессия и прогнозирование играют важную роль в принятии решений в бизнесе. Они позволяют предсказывать и анализировать будущие значения переменных на основе имеющихся данных. Это помогает бизнесу планировать и оптимизировать свою деятельность, принимать обоснованные решения и достигать своих целей.

## **2.4. Рекомендательные системы**

Рекомендательные системы являются важным инструментом в современном бизнесе, позволяющим предлагать пользователям персонализированные рекомендации товаров, услуг, контента и других элементов. Они основаны на алгоритмах машинного обучения, которые анализируют данные о предпочтениях и поведении пользователей для предсказания их предпочтений и предлагают соответствующие рекомендации.

Одной из основных целей рекомендательных систем является улучшение удовлетворенности пользователей и повышение конверсии. Путем предоставления релевантных и интересных рекомендаций, системы могут помочь пользователям находить нужные товары или контент, сэкономив их время и упростив процесс выбора. Также рекомендации способствуют удержанию пользователей и повторным покупкам, что в свою очередь может привести к увеличению вы-



ручки и прибыли компании.

Рекомендательные системы могут быть применены в различных отраслях, включая электронную коммерцию, медиа, социальные сети и другие. В электронной коммерции, например, они могут предлагать рекомендации товаров, основанные на истории покупок или просмотрах пользователей, а также использовать коллаборативную фильтрацию для нахождения схожих пользователей и предлагать им рекомендации, основанные на предпочтениях похожих пользователей.

### *Коллаборативная фильтрация*

Одним из наиболее распространенных алгоритмов, используемых в рекомендательных системах, является коллаборативная фильтрация. Этот метод основан на предположении, что если два пользователя проявили схожие предпочтения в прошлом, то они будут иметь схожие предпочтения и в будущем. Коллаборативная фильтрация использует матрицу оценок пользователей (например, оценки фильмов или товаров) для нахождения схожих пользователей или схожих товаров и рекомендует пользователю те элементы, которые оценили похожие пользователи.

Пример программы, реализующей коллаборативную фильтрацию для рекомендаций фильмов:

```
```python
import numpy as np
# Пример матрицы оценок пользователей
ratings = np.array([
```

```
[5, 4, 0, 0, 0, 0],  
[0, 0, 4, 0, 5, 0],  
[0, 0, 0, 2, 4, 5],  
[4, 0, 0, 0, 0, 4]  
)
```

Вычисление схожести пользователей на основе корреляции Пирсона

```
def compute_similarity(user1, user2):  
    mask = np.logical_and(user1 != 0, user2 != 0)  
    if np.sum(mask) == 0:  
        return 0  
    return np.corrcoef(user1[mask], user2[mask])[0, 1]
```

Функция рекомендации фильмов для пользователя

```
def recommend_movies(user_id, ratings,  
num_recommendations=5):
```

```
    num_users, num_movies = ratings.shape
```

Вычисление схожести пользователя с остальными пользователями

```
    similarities = []
```

```
    for i in range(num_users):
```

```
        if i != user_id:
```

```
            similarity = compute_similarity(ratings[user_id], ratings[i])
```

```
            similarities.append((i, similarity))
```

```
    similarities.sort(key=lambda x: x[1], reverse=True)
```

Выбор топ-N наиболее похожих пользователей

```
    top_similar_users = [similarity[0] for similarity in
```

```

similarities[:num_recommendations]]
# Получение рекомендаций на основе оценок похожих
пользователей
recommendations = np.zeros(num_movies)
for user in top_similar_users:
    recommendations += ratings[user]
    recommendations = np.where(ratings[user_id] == 0,
recommendations, 0)
    top_movies = np.argsort(recommendations)[::-1]
[:num_recommendations]
return top_movies
# Пример использования
user_id = 0
recommended_movies = recommend_movies(user_id,
ratings)
print(f"Рекомендованные фильмы для пользователя
{user_id}:")
for movie_id in recommended_movies:
    print(f"Фильм {movie_id}")
"""

```

В данном примере используется матрица оценок пользователей `ratings`, где каждая строка соответствует пользователю, а каждый столбец соответствует фильму. Оценки фильмов могут принимать значения от 0 до 5, где 0 обозначает отсутствие оценки.

Функция `compute_similarity` вычисляет схожесть поль-

зователей на основе корреляции Пирсона. Она сравнивает оценки двух пользователей, игнорируя нулевые значения, и вычисляет коэффициент корреляции.

Функция ``recommend_movies`` принимает идентификатор пользователя и матрицу оценок в качестве входных данных. Она вычисляет схожесть пользователя с остальными пользователями, выбирает топ-N наиболее похожих пользователей и выдает рекомендации на основе их оце-



Рекомендованные

Фильм 5

Фильм 4

Фильм 2

Фильм 3

Фильм 1

нок.

Пример использования демонстрирует, как получить рекомендации фильмов для определенного пользователя. Ре-

результатом программы является список идентификаторов фильмов, которые рекомендуется пользователю с указанным идентификатором.

Заметьте, что в данном примере использована простая реализация коллаборативной фильтрации. В реальных приложениях рекомендательных систем может потребоваться более сложные алгоритмы и обработка больших объемов данных.

Пример более сложной реализации коллаборативной фильтрации с использованием алгоритма Singular Value Decomposition (SVD) для рекомендаций фильмов:

```
import numpy as np
from scipy.sparse import csr_matrix
from scipy.sparse.linalg import svds
# Пример матрицы оценок пользователей
ratings = np.array([
    [5.0, 4.0, 0.0, 0.0, 0.0, 0.0],
    [0.0, 0.0, 4.0, 0.0, 5.0, 0.0],
    [0.0, 0.0, 0.0, 2.0, 4.0, 5.0],
    [4.0, 0.0, 0.0, 0.0, 0.0, 4.0]
])
# Выполнение сингулярного разложения (SVD)
def perform_svd(ratings, k):
    # Преобразование матрицы оценок в разреженную матрицу
    sparse_ratings = csr_matrix(ratings)
```

```

# Применение SVD для получения матриц U, Sigma и Vt
U, Sigma, Vt = svds(sparse_ratings, k)
# Построение диагональной матрицы Sigma
Sigma = np.diag(Sigma)
return U, Sigma, Vt

# Функция рекомендации фильмов для пользователя
def recommend_movies(user_id, ratings, U, Sigma, Vt,
num_recommendations=5):
    user_ratings = ratings[user_id]
    predicted_ratings = np.dot(np.dot(U[user_id, :], Sigma), Vt)
    # Исключение уже оцененных фильмов из рекомендаций
    predicted_ratings[user_ratings != 0] = -1
    top_movies      =      np.argsort(predicted_ratings)[::-1]
[:num_recommendations]
    return top_movies

# Пример использования
user_id = 0
k = 2 # Размерность скрытого пространства
U, Sigma, Vt = perform_svd(ratings, k)
recommended_movies      =      recommend_movies(user_id,
ratings, U, Sigma, Vt)
    print(f"Рекомендуемые      фильмы      для      пользователя
{user_id}:")
    for movie_id in recommended_movies:
        print(f"Фильм {movie_id}")
    """

```

В данном примере используется алгоритм Singular Value Decomposition (SVD) для выполнения сингулярного разложения матрицы оценок пользователей. Полученные матрицы U, Sigma и Vt представляют собой аппроксимацию исходной матрицы оценок с использованием латентного пространства низкой размерности.

Функция ``perform_svd`` выполняет сингулярное разложение матрицы оценок с помощью функции ``svds`` из модуля ``scipy.sparse.linalg``. Разложение возвращает матрицы U, Sigma и Vt.

Функция ``recommend_movies`` принимает идентификатор пользователя, матрицу оценок, а также матрицы U, Sigma и Vt в качестве аргументов. Она вычисляет предсказанные оценки для пользователя и рекомендует фильмы, имеющие наивысшие предсказанные оценки, исключая уже оцененные фильмы.

В приведенном примере выводится список рекомендованных фильмов для пользователя с идентификатором 0. Количество рекомендаций задается параметром

Рекомендуе

Фильм 5

Фильм 3

Фильм 2

Фильм 1

Фильм 0

`num_recommendations`.

Singular Value Decomposition (SVD), или Сингулярное разложение, является мощным алгоритмом линейной алгебры, который используется в различных областях, включая рекомендательные системы, сжатие данных, обработку изображений и многие другие.

Сингулярное разложение позволяет представить матрицу в виде произведения трех матриц: U , Σ и V^t . Формально, для матрицы A размерности $m \times n$ SVD определяется следующим образом:

$$A = U * \Sigma * V^t,$$

где U – матрица размерности $m \times m$, содержащая левые

сингулярные векторы,

Σ – диагональная матрица размерности $m \times n$, содержащая сингулярные значения,

V^t – транспонированная матрица размерности $n \times n$, содержащая правые сингулярные векторы.

Сингулярные значения в матрице Σ являются неотрицательными числами и упорядочены по убыванию. Они представляют собой меру важности каждого сингулярного вектора и определяют вклад каждого сингулярного вектора в исходную матрицу A .

При использовании SVD в рекомендательных системах, например, матрица A представляет собой матрицу оценок пользователей, где строки соответствуют пользователям, а столбцы – элементам (фильмам, продуктам и т.д.). SVD разделяет матрицу на скрытые факторы, представленные сингулярными векторами, и связывает их с пользователями и элементами. Это позволяет рекомендовать пользователям элементы, которые им могут понравиться, на основе сходства с другими пользователями или элементами.

Алгоритм SVD имеет несколько вариаций, которые могут быть использованы в зависимости от контекста и требований задачи. Некоторые из них включают Truncated SVD (SVD с ограниченным числом сингулярных значений), Implicit Matrix Factorization (IMF) и другие.

SVD является мощным инструментом для анализа данных и позволяет снизить размерность данных, извлекать

важные признаки и находить скрытые паттерны. Вместе с тем, алгоритм SVD требует значительных вычислительных ресурсов и может столкнуться с проблемами при обработке больших объемов данных. Поэтому для больших наборов данных используются приближенные методы SVD или альтернативные алгоритмы, такие как алгоритмы матричной факторизации.

Однако, SVD по-прежнему остается важным инструментом в области рекомендательных систем и других задач, где требуется анализ больших матриц данных.

Контекстная фильтрация

Еще одним распространенным методом является контентная фильтрация. Контентная фильтрация – это метод рекомендательных систем, который основывается на анализе характеристик элементов и предпочтений пользователей. В контексте контентной фильтрации, каждый элемент (товар, статья, фильм и т.д.) характеризуется набором признаков или характеристик, которые описывают его содержание или свойства.

Процесс контентной фильтрации начинается с анализа характеристик элементов и их значимости для пользователей. Характеристики элементов могут включать такие атрибуты, как автор, жанр, ключевые слова, рейтинги и другие свойства, которые могут быть извлечены из содержания элемента или предоставлены вручную.

Далее, на основе характеристик элементов, строится про-

фильм пользователя, который отражает его предпочтения и интересы. Профиль пользователя может быть создан путем анализа предыдущих взаимодействий пользователя с элементами, например, его рейтинги или история просмотров.

Затем, используя различные алгоритмы сходства, производится сравнение между профилем пользователя и характеристиками элементов. Целью является определение степени сходства между предпочтениями пользователя и характеристиками элементов.

На основе этого сравнения, система ранжирует и рекомендует пользователю элементы, которые наиболее соответствуют его предпочтениям. Например, если пользователь предпочитает фильмы определенного жанра, система может рекомендовать ему фильмы схожего жанра.

Преимуществом контентной фильтрации является то, что она не требует данных о предпочтениях других пользователей, так как она полностью основана на анализе характеристик элементов и предпочтениях пользователя. Это делает ее особенно полезной в случаях, когда у нас ограниченное количество данных о взаимодействиях пользователей.

Однако, контентная фильтрация также имеет свои ограничения. В частности, она может столкнуться с проблемой ограниченности характеристик элементов, особенно если характеристики не полностью охватывают аспекты предпочтений пользователя. Также возникает проблема обновления профиля пользователя и характеристик элементов с течени-

ем времени.

Метод является важным в рекомендательных систем, который позволяет рекомендовать пользователю элементы на основе их сходства с предпочтениями и характеристиками элементов. Она может быть эффективным инструментом в различных областях, таких как маркетинг, интернет-торговля, медиа и другие, где персонализированные рекомендации имеют важное значение для улучшения пользовательского опыта и увеличения продаж.

Рекомендательные системы также могут использовать гибридные подходы, комбинируя несколько методов для получения более точных и релевантных рекомендаций. Например, можно использовать коллаборативную фильтрацию для нахождения похожих пользователей и контентную фильтрацию для нахождения похожих элементов, и затем объединить результаты для формирования итоговых рекомендаций.

Рекомендательные системы являются мощным инструментом для улучшения пользовательского опыта, увеличения продаж и удержания клиентов. Они позволяют бизнесу создавать персонализированные рекомендации, основанные на данных и поведении пользователей, что способствует улучшению конкурентоспособности и достижению бизнес-целей.

Ниже приведен пример программы контентной фильтрации:

```

```python
Импорт необходимых библиотек
import pandas as pd
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity

Загрузка данных
data = pd.read_csv('movies.csv')

Создание матрицы TF-IDF на основе описаний фильмов
tfidf = TfidfVectorizer(stop_words='english')

tfidf_matrix =

tfidf.fit_transform(data['description'].fillna(""))

Вычисление матрицы сходства косинусной мерой
cosine_sim = cosine_similarity(tfidf_matrix, tfidf_matrix)

Функция для получения рекомендаций похожих филь-
MOV
def get_recommendations(title, cosine_sim, data, top_n=5):
 indices =
 pd.Series(data.index,
index=data['title']).drop_duplicates()
 idx = indices[title]
 sim_scores = list(enumerate(cosine_sim[idx]))
 sim_scores = sorted(sim_scores, key=lambda x: x[1],
reverse=True)
 sim_scores = sim_scores[1:top_n+1]
 movie_indices = [i[0] for i in sim_scores]
 return data['title'].iloc[movie_indices]

Пример использования функции для получения реко-

```

мендаций

```
movie_title = 'The Dark Knight Rises'
recommendations = get_recommendations(movie_title,
cosine_sim, data)
print(f"Рекомендации для фильма '{movie_title}':")
print(recommendations)
...
```

Программа выполняет следующие шаги:

1. Импортируются необходимые библиотеки. Библиотека ``pandas`` используется для работы с данными в виде таблицы, а библиотеки ``TfidfVectorizer`` и ``cosine_similarity`` из модуля ``sklearn.feature_extraction.text`` и ``sklearn.metrics.pairwise`` соответственно используются для работы с текстовыми данными и вычисления сходства между ними.

2. Загружаются данные о фильмах из файла `'movies.csv'` с помощью функции ``read_csv()`` из библиотеки ``pandas``. Данные обычно содержат информацию о фильмах, включая их названия, описания и другие атрибуты.

3. Создается объект ``TfidfVectorizer``, который преобразует текстовые описания фильмов в числовые векторы с использованием метода TF-IDF. TF-IDF (Term Frequency-Inverse Document Frequency) – это статистическая мера, используемая для оценки важности термина в документе. Он позволяет выделить ключевые слова и характеристики фильмов.

4. С помощью метода ``fit_transform()`` объекта

ТfidfVectorizer` преобразуется список описаний фильмов в матрицу TF-IDF. Эта матрица представляет собой числовое представление описаний фильмов, где каждый столбец соответствует определенному термину, а каждая строка – конкретному фильму.

5. Вычисляется матрица сходства между фильмами с использованием метода ``cosine_similarity()`` из модуля ``sklearn.metrics.pairwise``. Косинусное сходство измеряет угол между двумя векторами и предоставляет меру их сходства. В данном случае, матрица сходства показывает степень сходства между каждой парой фильмов на основе их описаний.

6. Определяется функция ``get_recommendations()``, которая принимает название фильма, матрицу сходства и данные о фильмах. Внутри функции происходит следующее:

- Создается объект ``pd.Series`` с индексами, соответствующими названиям фильмов и значениями, соответствующими их индексам в данных.

- Получается индекс выбранного фильма.

- Вычисляется список схожести выбранного фильма с остальными фильмами.

- Список сортируется по убыванию схожести.

- Выбираются топ-N фильмов на основе сходства.

- Возвращается список рекомендуемых фильмов.

7. Запрашивается у пользователя название фильма, для которого необходимо получить рекомендации.

8. Вызывается функция ``get_recommendations()``` с передачей ей названия фильма, матрицы сходства и данных о фильмах.

9. Выводятся на экран рекомендованные фильмы.

Программа использует алгоритм контентной фильтрации на основе TF-IDF и косинусного сходства для рекомендации фильмов на основе их текстовых описаний. Она преобразует текстовые данные в числовые векторы с использованием TF-IDF и затем вычисляет сходство между фильмами. Рекомендуемые фильмы выбираются на основе сходства с выбранным фильмом. Это позволяет предлагать пользователю фильмы, которые имеют схожие характеристики и описания с фильмами, которые он предпочитает.



# Глава 3: Подготовка данных для машинного обучения

*Качество данных определяет качество решений.  
Тщательная подготовка данных — залог успешного  
машинного обучения и эффективного бизнеса.*

В процессе применения машинного обучения в бизнесе подготовка данных играет важную роль. Качество данных определяет эффективность моделей машинного обучения и точность результатов, которые они предоставляют. В этой главе мы рассмотрим различные аспекты и задачи, связанные с подготовкой данных, и объясним, почему они важны для бизнеса.

Одной из причин, почему мы будем рассматривать подготовку данных, является достижение высокого качества прогнозов и решений. Чистые и точные данные являются основой для создания моделей машинного обучения, которые могут давать надежные прогнозы и принимать обоснованные решения. Подготовка данных помогает устранить шум, выбросы и другие аномалии, что повышает точность прогнозов и решений.

Другой важной ролью подготовки данных является оптимизация бизнес-процессов. Анализ данных, включенный в процесс подготовки, позволяет лучше понять структуру

и особенности данных. Это помогает оптимизировать бизнес-процессы и принимать обоснованные решения на основе данных. Например, анализ данных может выявить паттерны потребительского поведения, что позволит оптимизировать маркетинговые стратегии и улучшить взаимодействие с клиентами.

Также подготовка данных играет роль в персонализации и улучшении опыта клиента. Понимание предпочтений и потребностей клиентов на основе анализа данных позволяет создавать более персонализированные предложения и предлагать индивидуальные рекомендации. Это повышает уровень удовлетворенности клиентов и улучшает их опыт использования продуктов и услуг.

В данной главе мы рассмотрим различные задачи, связанные с подготовкой данных, включая сбор данных, очистку от шума и аномалий, анализ данных и обработку категориальных данных. Мы также рассмотрим методы и инструменты, которые помогут вам эффективно подготовить данные для использования в моделях машинного обучения.

### **3.1. Сбор, очистка и преобразование данных**

В мире машинного обучения и анализа данных сбор, очистка и преобразование данных играют ключевую роль. Эти этапы являются неотъемлемой частью подготовки данных перед применением алгоритмов машинного обучения.

В этой главе мы рассмотрим, почему эти действия важны и как они влияют на результаты анализа данных и принятие решений в бизнесе.

Сбор данных является первым и наиболее важным шагом. Для успешного машинного обучения необходимо иметь доступ к качественным и репрезентативным данным. Это может включать данные о клиентах, продуктах, транзакциях, рекламе и многом другом, в зависимости от конкретной задачи и области бизнеса. Сбор данных может осуществляться различными способами, включая опросы, сенсоры, базы данных, API и многое другое. Цель состоит в том, чтобы получить максимально полные и точные данные, которые позволят нам выявить закономерности и сделать правильные выводы.

Однако сырые данные не всегда готовы к использованию. Часто они содержат ошибки, пропуски, выбросы и другие неточности. Поэтому следующим шагом является очистка данных. Очистка данных включает в себя удаление или исправление ошибочных значений, заполнение пропущенных данных, удаление выбросов и приведение данных к единому формату. Цель состоит в том, чтобы убрать нежелательные влияния, которые могут исказить результаты анализа и прогнозирования.

После очистки данных часто требуется их преобразование. Преобразование данных может включать изменение формата, масштабирование, создание новых признаков и

многое другое. Например, числовые данные могут быть нормализованы, чтобы привести их к одному диапазону значений, или категориальные данные могут быть закодированы с использованием метода One-Hot Encoding для использования в алгоритмах машинного обучения. Преобразование данных позволяет создать более информативные и удобные для анализа наборы данных, а также улучшить производительность моделей машинного обучения.

Важно понимать, что сбор, очистка и преобразование данных являются итеративным процессом. В ходе анализа данных и разработки моделей могут возникать новые требования и потребности, которые потребуют обновления и доработки данных. Поэтому эти этапы являются непрерывным процессом, который требует внимания и усилий на протяжении всего жизненного цикла проекта. Понимание и умение применять эти методы позволит нам получить качественные данные и обеспечить надежные результаты анализа данных в бизнесе.

### **3.1.1. Сбор данных**

Раздел о сборе данных является важной частью подготовки данных для машинного обучения. Он занимается определением источников данных и разработкой методов их сбора.

Один из основных аспектов сбора данных – это определение необходимых данных для анализа и прогнозирования.

В бизнесе может быть множество различных типов данных, которые могут быть полезными для принятия решений, например, данные о клиентах, продажах, финансовых показателях или маркетинговых активностях. Важно определить, какие данные являются релевантными для вашей задачи и какие источники можно использовать для их получения.

Существует множество различных источников данных, которые можно использовать в бизнесе. Некоторые из них включают опросы и исследования, базы данных, внутренние системы и приложения, сенсоры и устройства интернета вещей (IoT), а также внешние источники данных через API (Application Programming Interface). Каждый источник данных имеет свои особенности и методы сбора.

При сборе данных необходимо обеспечить их качество и надежность. Это означает, что данные должны быть точными, полными, актуальными и соответствовать определенным стандартам. Во время сбора данных может возникнуть необходимость проверки и фильтрации данных, чтобы убедиться в их корректности. Также важно обеспечить безопасность данных и соблюдать соответствующие правила и регуляции в отношении конфиденциальности и защиты данных.

Для сбора данных могут использоваться различные методы и технологии. Например, для опросов и исследований можно применять онлайн-формы, телефонные интервью или личные встречи. Для сбора данных из баз данных можно использовать SQL-запросы или специальные ин-

струменты для извлечения данных. SQL (Structured Query Language) является стандартным языком для работы с реляционными базами данных. С помощью SQL-запросов можно выбирать, фильтровать и объединять данные из различных таблиц, а также проводить агрегацию и вычисления.

При работе с сенсорами и устройствами IoT (Internet of Things) может потребоваться настройка и мониторинг сенсоров для сбора нужной информации. Сенсоры могут собирать данные о различных параметрах, таких как температура, влажность, движение и другие. Для сбора данных от сенсоров могут использоваться специальные протоколы и средства связи, такие как Bluetooth, Wi-Fi или специальные сети передачи данных.

Использование API (Application Programming Interface) позволяет получать данные из сторонних сервисов или платформ. API предоставляют набор функций и методов, которые позволяют программно взаимодействовать с сервисами или приложениями. С помощью API можно получать данные о погоде, финансовых показателях, социальных медиа и других источниках. Это обеспечивает возможность интеграции с внешними системами и получения актуальной информации для анализа.

Каждый из этих методов сбора данных имеет свои особенности и требует соответствующей настройки и подготовки. Например, при использовании SQL-запросов необходимо быть знакомым с языком SQL и структурой базы дан-

ных. При работе с сенсорами и IoT-устройствами требуется установка и конфигурация сенсоров, а также обеспечение надежности и безопасности сети передачи данных. Использование API требует регистрации и получения ключа доступа, а также ознакомления с документацией и методами взаимодействия с сервисом.

Выбор конкретного метода сбора данных зависит от доступных ресурсов, специфики проекта и требований анализа данных. Каждый метод имеет свои преимущества и ограничения, поэтому важно выбрать наиболее подходящий для конкретной ситуации.

Определение необходимых данных является ключевым шагом в процессе сбора данных. Чтобы определить, какие данные нужны, следует учитывать цели и задачи анализа данных, а также специфику бизнеса или проекта. Важно начать с четкого определения целей анализа данных. Что именно вы хотите достичь с помощью анализа данных? Какие вопросы вы хотите ответить или какие решения вы хотите принять? Определите основные проблемы, которые вы хотите решить, и выделите ключевые метрики или показатели, которые помогут вам измерить успех.

Затем проанализируйте текущую ситуацию и ресурсы, которые у вас есть. Какие данные уже собираются или доступны в вашей компании или организации? Рассмотрите внутренние системы и базы данных, которые могут содержать полезную информацию. Определите, какие данные уже исполь-

зуются или собираются для других целей, и можно ли их переиспользовать или объединить.

Важно также рассмотреть внешние источники данных, которые могут быть полезны для ваших целей. Это могут быть открытые данные, сторонние сервисы или API, исследования и отчеты, данные от поставщиков или партнеров. Исследуйте, какие данные доступны в вашей отрасли или сфере деятельности, и определите, какие из них могут быть полезны для вашего анализа.

При определении необходимых данных следует также учитывать юридические и этические аспекты сбора данных. Обратите внимание на правила конфиденциальности и защиты данных, а также соответствие законодательству, связанному с сбором и использованием данных. Убедитесь, что вы имеете право собирать и использовать определенные данные и что вы принимаете меры для защиты приватности пользователей и конфиденциальности информации.

Важно также оценить качество данных, которые вы намерены собирать. Это включает проверку источников данных на достоверность и актуальность, а также обеспечение достаточной точности и полноты данных. Разработайте методы и процессы для контроля качества данных и фильтрации возможных ошибок или неточностей.

Корректный сбор данных является важным шагом для дальнейшего анализа и моделирования данных. От качества собранных данных зависит точность и надежность результа-



тов машинного обучения и прогнозирования. Поэтому внимательное и систематическое выполнение этого этапа является ключевым для успешной подготовки данных в бизнесе.

SQL-запросы и специальные инструменты для извлечения данных являются основными способами сбора данных из баз данных. Давайте рассмотрим каждый из них подробнее:

1. SQL-запросы: SQL (Structured Query Language) является стандартным языком для работы с реляционными базами данных. С помощью SQL-запросов можно выполнять различные операции, такие как выборка данных из таблиц, фильтрация, сортировка, объединение таблиц и другие. SQL предоставляет мощный и гибкий инструментарий для извлечения нужных данных из базы данных. Он позволяет составлять запросы на основе определенных условий и критериев, чтобы получить конкретные данные, необходимые для анализа или обработки.

2. Специальные инструменты для извлечения данных: Существуют различные инструменты, разработанные специально для удобного и эффективного извлечения данных из баз данных. Эти инструменты обычно предоставляют графический интерфейс и набор функций, которые облегчают выполнение запросов и работу с данными. Некоторые из популярных инструментов включают в себя MySQL Workbench, Microsoft SQL Server Management Studio, Oracle SQL Developer и другие. Они обеспечивают удобную среду

для написания SQL-запросов, просмотра и редактирования данных, а также визуализации результатов запросов.

Выбор между использованием SQL-запросов и специальных инструментов зависит от предпочтений и требований пользователя. SQL-запросы предоставляют более гибкий и мощный подход, позволяющий полностью контролировать процесс извлечения данных и выполнять сложные операции. С другой стороны, специальные инструменты обладают удобным интерфейсом и предоставляют дополнительные функции, такие как визуализация данных и автоматическая генерация запросов. Они особенно полезны для пользователей, не знакомых с языком SQL или предпочитающих более интуитивный способ работы с данными.

Важно отметить, что для использования SQL-запросов или специальных инструментов требуется доступ к базе данных и соответствующие привилегии. Также необходимо иметь понимание структуры базы данных, таблиц и связей между ними, чтобы эффективно составлять запросы и получать нужные данные.

В зависимости от конкретной ситуации и требований проекта, можно выбрать наиболее удобный и эффективный способ сбора данных из баз данных.

Рассмотрим несколько примеров использования SQL-запросов и специальных инструментов для извлечения данных из баз данных:

1. Пример использования SQL-запросов:

Предположим, у нас есть база данных с информацией о клиентах и их заказах в интернет-магазине. Мы можем написать SQL-запрос для извлечения данных о клиентах, сделавших заказы на определенную дату:

```
```sql
SELECT * FROM Customers
JOIN Orders ON Customers.CustomerID =
Orders.CustomerID
WHERE Orders.OrderDate = '2023-05-31';
```
```

В результате этого запроса мы получим все записи о клиентах и их заказах, сделанных 31 мая 2023 года.

## 2. Пример использования специального инструмента:

Допустим, у нас есть база данных с информацией о сотрудниках компании. Мы можем использовать инструмент MySQL Workbench для просмотра и редактирования данных. С помощью графического интерфейса инструмента мы можем выполнить запрос на выборку данных, например, для получения списка всех сотрудников определенного отдела:

Открываем MySQL Workbench и подключаемся к базе данных.

Выбираем нужную таблицу (например, "Employees").

Нажимаем кнопку "Execute SQL" и вводим запрос:

```
```sql
SELECT * FROM Employees WHERE Department =
'Marketing';
```
```

---

Нажимаем кнопку "Execute" или "Run" для выполнения запроса.

В результате мы увидим список всех сотрудников, работающих в отделе маркетинга.

### 3. Пример использования SQL-запросов:

Предположим, у нас есть база данных с информацией о студентах и их оценках. Мы можем написать SQL-запрос для извлечения среднего балла студентов по предмету:

```
```sql
```

```
SELECT Subject, AVG(Grade) AS AverageGrade
```

```
FROM Students
```

```
GROUP BY Subject;
```

```
```
```

В результате этого запроса мы получим список предметов и соответствующие средние оценки студентов по каждому предмету.

# Конец ознакомительного фрагмента.

Текст предоставлен ООО «ЛитРес».

Прочитайте эту книгу целиком, [купив полную легальную версию](#) на ЛитРес.

Безопасно оплатить книгу можно банковской картой Visa, MasterCard, Maestro, со счета мобильного телефона, с платежного терминала, в салоне МТС или Связной, через PayPal, WebMoney, Яндекс.Деньги, QIWI Кошелек, бонусными картами или другим удобным Вам способом.