

Введение: Проблема поиска по смыслу

Представьте, что вы ищете информацию на своем компьютере. Обычно вы используете поиск по ключевым словам. Если вы ищете документ со словом «отчет», система найдет все файлы, где это слово встречается в точности. Но что, если нужный вам документ называется «ежеквартальный итог», а слово «отчет» в нем не используется? Традиционный поиск его не найдет.

Здесь возникает идея **семантического поиска**, то есть поиска по *смыслу*, а не по точным словам. Вы хотите найти документы, которые *похожи по значению* на ваш запрос. Например, при поиске «грустные фильмы о собаках» вы бы хотели найти фильм «Хатико: Самый верный друг», даже если в его описании нет именно этих слов.

Чтобы научить компьютеры понимать смысл, нам нужен способ представить значение слов, предложений, изображений и других данных в числовом формате. Именно эту задачу решают **векторные представления (embeddings)**, и для работы с ними были созданы специализированные инструменты, такие как **векторная база данных Qdrant**.

Глава 1: Что такое векторы (Embeddings)?

1.1. От данных к числам

В основе современных систем искусственного интеллекта (AI), таких как большие языковые модели (Large Language Models, **LLM**), лежит способность преобразовывать сложные данные (текст, изображения, аудио) в математические объекты, которые компьютер может обрабатывать. Таким объектом является **вектор**.

В контексте AI, вектор — это просто упорядоченный список чисел. Каждое число в этом списке представляет какую-то характеристику или «измерение» исходных данных.

Простая аналогия: Представьте, что мы хотим описать фильм с помощью трех характеристик: «Комедия», «Экшен» и «Драма», оценивая каждую по шкале от 0 до 10.

- Фильм «Мстители» можно представить как вектор: [2, 9, 4] (мало комедии, много экшена, немного драмы).
- Романтическую комедию «Дневник Бриджит Джонс» — как вектор: [8, 1, 6] (много комедии, мало экшена, средняя драма).
- Драму «Список Шиндлера» — как вектор: [1, 2, 10] (почти нет комедии и экшена, максимум драмы).

Этот список чисел и есть вектор.

1.2. Векторные представления (Embeddings)

Современные AI-модели, такие как **GPT**, **Claude** или **Llama**, делают то же самое, но в гораздо большем масштабе. Они могут взять фрагмент текста (слово, предложение или целый документ) и преобразовать его в вектор, состоящий не из 3, а из сотен или даже тысяч чисел (измерений). Этот процесс называется **созданием векторного представления** или **embedding**.

Этот полученный вектор (embedding) — это числовой «слепок» или «ДНК» семантического значения исходных данных.

Ключевая идея заключается в следующем: **объекты с похожим смыслом будут иметь похожие векторы**.

Например, векторы для предложений «Какая сегодня погода?» и «Будет ли сегодня дождь?» будут находиться в многомерном пространстве очень близко друг к другу. А вектор для предложения «Я люблю пиццу» будет находиться далеко от них.

Эти векторы создаются с помощью сложных нейронных сетей, таких как **трансформеры (transformers)**, которые обучаются на огромных объемах текстовых и других данных, чтобы научиться улавливать нюансы языка и смысла. В основе архитектуры трансформеров лежит механизм **self-attention** (само-внимания), позволяющий модели взвешивать важность разных слов в предложении при создании его общего векторного представления.

Глава 2: Поиск похожих векторов

Теперь, когда у нас есть способ превращать данные в векторы, возникает следующий вопрос: как математически определить, что два вектора «близки» друг к другу? Для этого используются метрики расстояния или сходства.

2.1. Расстояние в многомерном пространстве

Представьте себе обычную двухмерную карту (пространство с двумя измерениями: X и Y). Расстояние между двумя точками можно измерить линейкой. В мире векторов мы делаем то же самое, но в пространстве с сотнями измерений.

Две самые популярные метрики:

1. **Евклидово расстояние (Euclidean Distance, L2)**
2. **Косинусное сходство (Cosine Similarity)**

2.2. Математика: Евклидово расстояние

Евклидово расстояние — это «прямолинейное» расстояние между двумя точками (концами векторов) в пространстве.

Формула: Для двух векторов \vec{A} и \vec{B} в n-мерном пространстве, где $\vec{A} = (a_1, a_2, \dots, a_n)$ и $\vec{B} = (b_1, b_2, \dots, b_n)$, евклидово расстояние $d(\vec{A}, \vec{B})$ вычисляется так:

$$d(\vec{A}, \vec{B}) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \dots + (a_n - b_n)^2} = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

Пошаговое объяснение формулы:

1. $(a_i - b_i)$: Для каждого измерения (от 1 до n) мы находим разницу между соответствующими компонентами векторов \vec{A} и \vec{B} . Это говорит нам, насколько далеко друг от друга находятся векторы по этой конкретной оси.
2. $(a_i - b_i)^2$: Мы возводим каждую разницу в квадрат. Это делается по двум причинам: во-первых, чтобы все значения были положительными (расстояние не может быть отрицательным), и во-вторых, чтобы придать больший «вес» большим различиям.
3. $\sum_{i=1}^n$...: Знак суммы (сигма) означает, что мы складываем все эти квадраты разностей по всем измерениям. Мы получаем общую сумму квадратов расстояний по всем осям.
4. $\sqrt{\dots}$: Мы извлекаем квадратный корень из этой суммы. Это возвращает нас к исходным единицам измерения и является аналогом теоремы Пифагора для многомерного пространства.

Практическое применение: Чем меньше значение евклидова расстояния, тем ближе векторы друг к другу и, следовательно, тем более похожи исходные данные. Если расстояние равно 0, это означает, что векторы идентичны.

2.3. Математика: Косинусное сходство

В отличие от евклидова расстояния, которое измеряет дистанцию, косинусное сходство измеряет **угол** между двумя векторами. Оно показывает, насколько векторы «смотрят» в одном направлении, независимо от их длины (магнитуды). Это особенно полезно для текстовых данных, где длина документа может сильно варьироваться, но нас интересует именно его тематика (направление).

Формула: Косинусное сходство между векторами \vec{A} и \vec{B} вычисляется как косинус угла θ между ними:

$$\text{similarity}(\vec{A}, \vec{B}) = \cos(\theta) = \frac{\vec{A} \cdot \vec{B}}{\|\vec{A}\| \|\vec{B}\|} = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2} \sqrt{\sum_{i=1}^n b_i^2}}$$

Пошаговое объяснение формулы:

1. $\vec{A} \cdot \vec{B}$ (**Скалярное произведение, Dot Product**): Это числитель дроби. Он вычисляется как $\sum_{i=1}^n a_i b_i$. Мы перемножаем соответствующие компоненты векторов и складываем результаты. Скалярное произведение будет большим, если компоненты векторов имеют одинаковые знаки и большие значения, то есть если векторы «со-направлены».
2. $\|\vec{A}\|$ (**Магнитуда или норма вектора A**): Это знаменатель дроби. Магнитуда — это длина вектора, вычисляемая как $\sqrt{\sum_{i=1}^n a_i^2}$. Это, по сути, евклидово расстояние от начала координат (вектора $[0, 0, \dots, 0]$) до конца вектора \vec{A} . То же самое делается для вектора \vec{B} .
3. \dots : Мы делим скалярное произведение на произведение магнитуд векторов. Эта операция называется **нормализацией**. Она убирает влияние длины векторов и оставляет только информацию об их направлении.

Практическое применение: Результат косинусного сходства всегда находится в диапазоне от -1 до 1:

- **1:** Векторы указывают в одном и том же направлении (максимальное сходство). Угол между ними 0° .
- **0:** Векторы ортогональны (перпендикулярны) друг другу (нет сходства). Угол 90° .
- **-1:** Векторы указывают в противоположных направлениях (максимальная непохожесть). Угол 180° .

При поиске похожих текстов мы ищем векторы с косинусным сходством, близким к 1.

2.4. Проблема масштаба

Представьте, что у вас есть миллион документов. Когда поступает новый поисковый запрос, его нужно превратить в вектор, а затем вычислить его расстояние до *каждого* из миллиона векторов в вашей коллекции. Это называется **полным перебором (brute-force search)**. Для небольшого количества векторов это работает, но для миллионов или миллиардов — это невероятно медленно и непрактично.

Именно эту проблему решают векторные базы данных.

Глава 3: Что такое векторная база данных?

Векторная база данных (Vector Database) — это специализированная система, созданная для эффективного хранения, индексирования и поиска огромного количества векторов.

Ее главное отличие от традиционных баз данных (например, SQL):

- **Традиционная БД** работает со структурированными данными (таблицы, строки, колонки) и выполняет поиск по точным совпадениям (например, `WHERE user_id = 123`).
- **Векторная БД** работает с неструктурированными данными (в виде векторов) и выполняет поиск по сходству (например, «найди 10 векторов, наиболее близких к этому»).

3.1. Приближенный поиск ближайших соседей (ANN)

Ключевая технология, лежащая в основе векторных баз данных, — это алгоритмы **приближенного поиска ближайших соседей (Approximate Nearest Neighbor, ANN)**.

- **Точный поиск (KNN, K-Nearest Neighbors)**, как мы выяснили, требует сравнения со всеми векторами и гарантирует нахождение абсолютно самых близких соседей, но он медленный.
- **Приближенный поиск (ANN)** идет на компромисс: он жертвует незначительной долей точности ради огромного прироста в скорости. Вместо того чтобы гарантировать нахождение 100% самых близких соседей, он находит, например, 99% из них, но делает это в сотни или тысячи раз быстрее. Для большинства практических задач (поиск, рекомендации) этого более чем достаточно.

3.2. HNSW: Двигатель быстрого поиска

Один из самых популярных и эффективных ANN-алгоритмов, который используется в Qdrant, — это **HNSW (Hierarchical Navigable Small World)**, что переводится как «Иерархический навигационный мир тесных связей».

Структура: HNSW строит сложную структуру данных в виде многослойного графа. Представьте себе это как несколько карт одной и той же местности с разным уровнем детализации. Слой 2 (входной) — (точка входа) | Слой 1 — | | Слой 0 (детальный) — — — — —

- **Слой 0 (нижний):** Содержит все векторы вашей базы данных. Это самая детальная «карта».
- **Более высокие слои (1, 2, ...):** Содержат лишь небольшую часть векторов из нижних слоев, служа «скоростными магистралями» для навигации.

Процесс поиска:

1. Поиск начинается с точки входа на самом верхнем, самом разреженном слое.
2. Алгоритм находит на этом слое вектор, ближайший к вашему запросу.
3. Затем он «спускается» на слой ниже, используя найденный вектор как отправную точку, и повторяет поиск на более детальной «карте».
4. Этот процесс продолжается до тех пор, пока алгоритм не достигнет самого нижнего, детального слоя (Слой 0), где он и находит ближайших соседей для вашего запроса.

Это похоже на поиск адреса: вы сначала находите страну, потом город, потом улицу, а не проверяете каждый дом на планете. Благодаря такой структуре сложность поиска составляет примерно $O(\log n)$, что делает его чрезвычайно быстрым. Важным преимуществом является то, что добавление новых векторов не требует полной перестройки всей структуры.

Глава 4: Знакомство с Qdrant

Qdrant (произносится как «квадрант») — это высокопроизводительная **open-source** векторная база данных. Она написана на языке **Rust**, что обеспечивает высокую скорость работы и безопасность при работе с памятью. Qdrant создана для **production-ready** решений, способных работать с большими объемами данных.

4.1. Основные концепции Qdrant

Чтобы понять, как работает Qdrant, нужно знать несколько ключевых терминов:

- **Коллекция (Collection):** Это аналог таблицы в SQL-базе данных. Коллекция — это именованное хранилище для ваших векторов. Например, у вас может быть коллекция `product_images` для векторов изображений товаров и коллекция `user_manuals` для векторов текстовых документов. Все векторы в одной коллекции обычно имеют одинаковую размерность (например, 1536 измерений).
- **Точка (Point):** Это аналог строки в таблице. Каждая точка в коллекции представляет один объект данных. Точка состоит из трех основных частей:
 1. **ID:** Уникальный идентификатор точки (число или строка).
 2. **Вектор (Vector):** То самое векторное представление (embedding) вашего объекта.
 3. **Полезная нагрузка (Payload):** Это необязательные метаданные, связанные с вектором.
- **Полезная нагрузка (Payload):** Это одна из самых мощных возможностей Qdrant. Вектор — это просто набор чисел, он не несет в себе исходной информации. Payload позволяет хранить вместе с вектором любую полезную информацию в формате JSON.

Пример: Вы превратили описание фильма «Хатико» в вектор `[-0.1, 0.8, ...]`. Сам по себе этот вектор бесполезен для пользователя. В Payload вы можете сохранить: `json { "title": "Hachi: A Dog's Tale", "year": 2009, "genre": "Drama", "url": "https://example.com/hachi" }` Когда вы выполняете семантический поиск и находите этот вектор, Qdrant возвращает вам не только вектор, но и его Payload. Таким образом, вы можете показать пользователю название фильма, год выпуска и ссылку на него.

- **Фильтрация (Filtering):** Это ключевое преимущество Qdrant. База данных позволяет комбинировать семантический поиск по векторам с традиционной фильтрацией по полям в Payload. В отличие от многих других систем, Qdrant поддерживает **полноценную и гибкую фильтрацию**, включая **вложенные структуры JSON** и **сложные иерархии**. Вы можете выполнять сложные запросы, например:
 - «Найди товары, похожие по смыслу на *‘удобное офисное кресло’*, но только те, у которых цена < 300 , цвет = *‘черный’* И рейтинг.средний > 4.5 ».

Qdrant эффективно выполнит такой запрос, сначала отфильтровав точки по метаданным, а затем выполнив быстрый ANN-поиск среди оставшихся.

4.2. Qdrant в экосистеме: когда его выбирать?

Понимание, когда использовать Qdrant, а когда другие инструменты, очень важно. Рассмотрим два практических примера.

Пример 1: Приложение для заметок NimbleNote

- **Задача:** Создать быстрый прототип (MVP) приложения для поиска по личным заметкам.
- **Условия:** 1-10 млн заметок, низкая нагрузка (меньше 30 запросов в секунду), развертывание на одном сервере, важна простота и скорость разработки.
- **Начальный выбор:** Можно начать с более простой базы, как Chroma, которая очень легко устанавливается и не требует сложной настройки.
- **Эволюция:** Со временем приложению потребовалась более **гибкая фильтрация** (например, поиск по тегам, датам, вложенным категориям с логикой AND/OR) и **горизонтальное масштабирование** для роста числа пользователей.
- **Переход на Qdrant:** Qdrant становится идеальным следующим шагом, так как он предоставляет мощную фильтрацию и поддерживает **шардинг** (разделение данных на несколько серверов) для масштабирования.

Пример 2: Каталог товаров FinCommerce

- **Задача:** Production-система поиска по десяткам миллионов товаров.
- **Условия:** Требуются **сложные фильтры** (категории, диапазоны цен, теги, права доступа), высокая нагрузка (тысячи запросов в минуту).
- **Выбор:** Qdrant — идеальный кандидат. Он создан для production-нагрузок, его система фильтрации отлично справляется со сложными каталогами, а горизонтальное масштабирование коллекций позволяет выдерживать

нагрузку.

- **Эволюция (гипермасштабирование):** Если нагрузка вырастает до сотен миллионов векторов и десятков тысяч запросов в секунду, может потребоваться еще более сложная архитектура.
- **Следующий шаг:** Системы вроде Milvus с их микросервисной архитектурой позволяют **независимо масштабировать** узлы для записи и узлы для поиска. Это нужно в очень крупных, высоконагруженных проектах.

Таким образом, Qdrant занимает золотую середину: он достаточно прост для старта, но при этом мощный и масштабируемый для большинства production-сценариев.

Глава 5: Практические сценарии использования Qdrant

5.1. Семантический поиск

Это классический пример: создание «умного» поиска для вашего сайта, базы знаний, интернет-магазина или архива документов. Пользователи могут искать информацию, используя естественный язык, и получать релевантные результаты, даже если ключевые слова не совпадают.

5.2. Системы рекомендаций

Qdrant идеально подходит для создания рекомендательных движков.

- **«Пользователи, купившие товар X, также покупали...»:** Найдите вектор товара X. Затем найдите в Qdrant векторы других товаров, которые находятся к нему ближе всего. Это и будут ваши рекомендации.
- **Персональные рекомендации:** Создайте вектор, представляющий интересы пользователя (на основе его истории просмотров, лайков, покупок). Затем ищите в базе контента (статей, видео, товаров) векторы, наиболее близкие к вектору интересов пользователя.

5.3. Retrieval-Augmented Generation (RAG)

Это один из самых важных и современных сценариев использования векторных баз данных в связке с большими языковыми моделями (LLM), такими как **GPT-4** или **Claude 3**.

Проблема: LLM обладают огромными общими знаниями, но они не знают:

- Информацию, появившуюся после даты их обучения.
- Ваши внутренние, частные или конфиденциальные данные (документы компании, переписку в техподдержке, личные заметки).
- Они могут «галлюцинировать» — выдумывать факты, если не знают точного ответа.

Решение (RAG): Вместо того чтобы напрямую задавать вопрос LLM, мы сначала используем Qdrant, чтобы «подкормить» модель нужным контекстом.

Процесс RAG шаг за шагом:

1. Фаза индексации (заранее):

- Вы берете свои документы (например, внутреннюю документацию компании).
- Разбиваете их на небольшие, осмысленные части (чанки).
- С помощью embedding-модели превращаете каждый чанк в вектор.
- Сохраняете эти векторы и соответствующие им текстовые чанки (в Payload) в коллекцию Qdrant.

2. Фаза запроса (в реальном времени):

- **Шаг 1: Запрос пользователя.** Пользователь задает вопрос: «Какие у нас правила оформления отпуска?»
- **Шаг 2: Создание вектора запроса.** Ваш вопрос «Какие у нас правила оформления отпуска?» превращается в вектор с помощью той же embedding-модели.
- **Шаг 3: Поиск в Qdrant.** Этот вектор запроса отправляется в Qdrant для поиска наиболее похожих векторов (например, топ-3 самых близких).
- **Шаг 4: Извлечение контекста.** Qdrant возвращает наиболее релевантные чанки текста из вашей документации (например, «Сотрудник имеет право на 28 календарных дней отпуска...», «Заявление на отпуск необходимо подать за 2 недели...»).
- **Шаг 5: Аугментация (дополнение) промпта.** Вы формируете новый, расширенный промпт для LLM (этот процесс называется **prompt engineering**), который выглядит примерно так: Используя только приведенный ниже контекст, ответь на вопрос пользователя. Контекст: --- [Текст чанка 1 из Qdrant] [Текст чанка 2 из Qdrant] [Текст

чанка 3 из Qdrant] --- Вопрос: Какие у нас правила оформления отпуска?

- **Шаг 6: Генерация ответа.** Этот расширенный промпт отправляется в LLM (например, GPT-4). Теперь модель не пытается вспомнить общие правила отпуска, а генерирует точный ответ, основываясь *исключительно на предоставленном вами контексте* из ваших документов.

Таким образом, **RAG** с использованием Qdrant позволяет создавать чат-ботов и вопросно-ответные системы, которые могут точно и правдиво отвечать на вопросы по вашим собственным данным, значительно снижая риск галлюцинаций.

Заключение

Qdrant и другие векторные базы данных являются фундаментальным строительным блоком для нового поколения AI-приложений. Они решают критически важную задачу — организацию и быстрый поиск данных на основе их семантического значения.

Ключевые концепции, которые вы изучили:

- **Вектор (Embedding):** Числовое представление смысла данных.
- **Метрики сходства:** Способы измерить «близость» векторов (Евклидово расстояние, Косинусное сходство).
- **Векторная база данных:** Специализированная БД для хранения и быстрого поиска векторов.
- **ANN-поиск (HNSW):** Приближенный, но очень быстрый метод поиска похожих векторов.
- **Qdrant:** Конкретная реализация векторной БД, написанная на Rust, с коллекциями, точками и мощной фильтрацией по полезной нагрузке (Payload).
- **RAG:** Мощная техника для создания точных вопросно-ответных систем с помощью LLM и векторных баз данных.

Понимание этих принципов открывает дверь в мир создания интеллектуальных систем поиска, рекомендаций и взаимодействия с AI на основе ваших собственных данных. Следующим шагом для углубления знаний может стать изучение практических руководств по установке Qdrant и созданию своего первого RAG-приложения.