# Lecture 39

- job opportunity ...
  with learning rare ...

"Quality of response is matter"

- rare learning

job opportunity
↓
4th positions ✓
~ 24 positions
agenda → winning Team

## Vector Stores

- Vector store
- Cloud storage
- Local storage

- earlier data stored in rows & columns
  ↳ RDMS

- Vector → Representational
  ↳ Vector database / vector store

- different from conventional db

- vector store mainly used in recommendation.

"vector store are specialized databases optimized
for storing & querying high dimensional vector data".

- efficiently store & retriee vector

Feature ⌐ high speed similarity search
       └── scalability integration

⇒ each vector database store far 4 mei

① vector / vector embedding
② vector metadata
   ( information about data)
③ original information/data   (...)
④ unique id
n each vector


How vector store works
① cosine similarity
② dot product

ⅲ vector O و ـ ے ٹوٹتا کٹتا جنز چیز جس
کہ آپ کے چیز یہ @ ـتی ورژن سے کہ
especially in recommender system

for example 1 vector store a million embedding
              many algorithms   ←┘

              business point of view

• further research
• horror story ( know the vocabulary)
  "why"        why are you doing.

## Lang Chain & Vector Store →

→ ① fully manged service / cloud
(third party)
(api)

② locally manged database
(mange itself all by yourself)
→ resource cost ↑

↳ provide seamless integration

Benefits
- efficient data storage
- fast retrieval
- improved semantic understanding

why vector.

Textual data ——— embedding model → vector form
(computer can't understand) ···(numbers)

(computer can understand
by comparing vector value)

- Processing can't understand textual data
  - meaningful information

facebook ai semantic search



locally
mange

Vector Store → Faiss
→ qdrant

Chroma — Pinecone ↳ cloud / fully managed
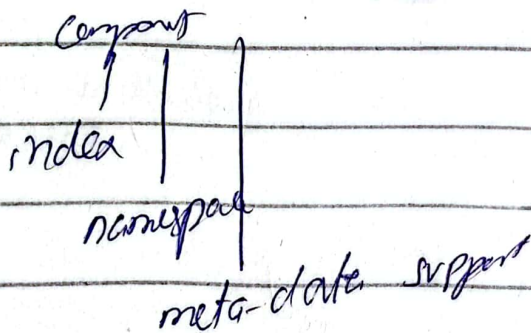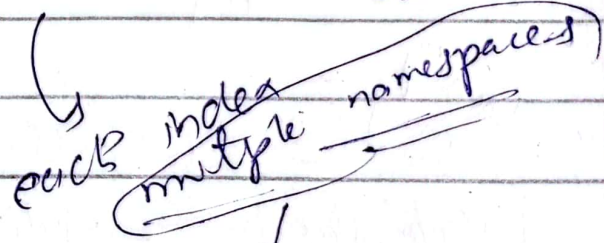
## FAISS

↳ library for effi———

## use cases

① u in RAG (Semantic Search)

② → recommendation system

(prefrences, queries,

③ → Anomaly detection

∴ can be multiple use cases. depend on critical thinking

## Pine cone architecture

Component

index

namespace

meta-data support

index ⟷ vector store

each index multiple namespaces

for concept of partition

index was 2.

partitions list (namespace) is

# Qdrant

cluster — multiple namespace
(vector stores)

- Pinecone, Quadrant ⟹ fully managed.

  login → project → vector store
  (different names on different platforms)

  Pinecone ⟹ index
  qdrant ⟹ cluster

## Practical

— review of previous 8 steps.

- embedding model    (∵ BAAI/bge small-en-v1.5)
  - load the model    (by PDF loader)
  - spliting data into chunks
  - meta data pre-processing    (edit meta data)
    (spliting page into chunks
  - updated meta data
    (you define of our own)

## Managed Services

- Qdrant credential
  url
  key
  collection name

  Automatically manages
  embeddings.

  ( directly create _____
   _____ )

- Do not repeat redundant time

  ┌─────────────────────────┐
  │ ∴ K ⇒ number of chunks  │
  └─────────────────────────┘
         ↓ important question

q. drant similarity search

  query → vector form → vector store
  cosine similarity
  → vector ے کم ہو یا vector الگ
  ے ( distance) up ہو تو الگ

- Pine cone

  same method

  vector similarity search

                        serless
                       (data store +
                          retrive ) cost
  Pinecone {
                       Pods
                       ( store + space
                         + retry)
                          ↓↑
                         cost

→ update your knowledge latest
→ necessary & critical step