

[Day 2] Unit 2 - "embeddings & vector stores / Pathbases"



Embeddings

- "cheat sheet" for machines to understand meaning of text, images & other data.
- help make sense of massive data.
- low dimensional numerical representation of "anything"
- designed to capture meaning & relationships between different pieces of data.
 - ↳ compact vectors preserve semantic information
 - maps complex data into a simpler space
- * Incredibly efficient for all types of data, save storage space, makes possible to find patterns/cls.

Semantic Rel.

- embeddings capture how similar / diff two pieces of data are (ex: king & queen → bike)
- encode meaning, large scale processing
- Retrieval & Recommendation: precompute embedding, query is converted to embedding, find nearest embedding
 - retrieves & rec. semantically similar embeddings.
 - recommend: find closest similar embeddings to what user prov. liked/intrst.

Joint-Embedding

- Handle multimodal data (pics, videos, text)
- allow for more nuanced understanding

- see how well they return relevant data, & filter out irrelevant
- Precision: amount of retrieved items that are relevant.
- Recall: proportion of all relevant items you managed to find.
@ $K=5$; how many top 5 pred. are correct.

Normalized discount. • gives higher scores when relevant items are at top of results.

Cumulative gain.

- higher = better

BERT / MTFB

- standardized collections of data sets & tasks, allow evaluate & compare different emb. models (fair & consistent)

Hugging Emb. Model

- key factors: model size, embedding dimensionality, latency/exit.
- Removal avg. gen. • using emb. to make lang. models "smarter", use embedding to find info from know. base, use info to boost prompts gen. to lang.
- ↑ accurate & ↑ relevant responses

2-stage process

- 1st
 - Create index, generate emb. for chunks using a doc encoder, storing emb. in vector DB.
- 2nd
 - Query processing stage: take vec q, turn into emb., perform similarity search in vector DB to find chunks closest to q
 - Need for labeled data for training (\$ + time cons)
→ ex. sd: run LLM to generate synthetic question + doc pairs

Challenges

Types of embeddings

- Text: allow to np. text as dense numerical vectors
 - * powerful for NLP tasks
- Token/Word: Breakdown text string into smaller meaningful units (each unique gets unique # ID)
- One-Hot embedding: np token IDs as Binary vectors
 - * don't capture semantics
- Word emb: dense, fixed length vectors aim to capture meaning of individ. words
- Word2vec: meaning of word is determined by words around it
 - architech: CBOW & skip gram → better for infrequent words & smaller DS
 - limits: small window of context
- Co-occurrence matrix: conta how often each word appears in context of every other word in DS
- Doc emb: emb. large pieces of text
 - bag of words model (LSA: look for underlying themes, LDA: probab. approach) (TFIDF: stat. methods that weigh words on how often they appear)
- Visual emb: training convolutional neural networks or vision trans.
- Multi-modal emb: combining text & image emb.
- Graph emb: np. objects and relationships within a network.
 - algos: deepwalk, node2vec

Training emb.

- on dual encoder architech, trained w/ contrastive loss
 - pulls emb. of similar data closer together & repelling diff.
 - 2 stages: pretraining & fine tuning
- finding items based on their incinity not just matching kw
 - in algo like approx. nearest neighbor.

Vector search



ANN tricks

- locality sensitive hashing, tree based methods
 - HNSW go-to algo. handles new insertion & vectors efficient.
 - SCAN: partitioning
- data emb. → index vector → embed. query → search for similar
- Info retrieval, recommend., semantic text similarity, classification, clustering

[Live Stream]

- Emb. advancement: using LLMs to speed up emb. training

- LLMs: filter training data sets. & data curation.

- matrixKA emb:

• Alloy DB

- Efficient hash for vector search: benchmarks for vector

→ RAG search: combining keyword + embeddings (hybrid search)

→ filtering and re-ranking.

- Upgrading all embedding models when using new ones.

→ Upgrade embeddings as well.

* AI search: efficient way to match prev. emb to new.

AI pair to Upgrade DB.

- RAG + measuring drift, focus more on recall > retrieval

* bridge gap between index and new context (new data)

- Freshness of data as a new eval. metric.

Pop quiz:

- 1) Application emb are NOT best suited for

→ simple rule based system

- 2) Adv. or scan > ANN search algos?

→ designed for high-dim. data, & excellent acc / latency tradeoffs

- 3) weakness of BOW for gen. embeddings?

→ ignore word ordering + sem. meaning

- 4) challenge when using emb for search? address?

→ emb might not capture literal meaning

- 5) Adv. of LSH for vector search

→ reduces search space by reducing dims to buckets

- [code labs]
- Doc. QA w/ RAG
- wing chroma
- limits of LMs: only "know" info - trained on, limited input context window → address by: RAG
 - RAG (Retrieval augmented generation) stages:
 - 1) Indexing
 - 2) Retrieval
 - 3) Generation
 * allows for providing info that model hasn't seen b4
 - Data → chm emb. DB w/ ChromaDB → retrieval (find relevant docs) → Augmented gen (answer Q's)

Embeddings and Similarity Scores

- Calc. similarity scores of 2 emb. vectors: $u^T v$
- Range: 0 = dissimilar, Range = 1 = comp. similar
- calc by matrix scl-mult: $\text{sim} = \text{df} @ \text{df.T}$

Classifying emb. w/ Keras + Hu Gemini API

- Dataset + Embeddings (comp. on @ a time) → Build classification model. → train model → evaluate model performance