

vector Database [Backbone of RAG]

①

DataSet: three sentences, each ~~has~~^{as} 3 words (or tokens)

how are you

who are you

who am I

(In practice, a dataset may contain millions or billions of sentences. The max number of tokens may be tens of thousand / billions. For instance Mistral Language model includes 7B parameters.

process

process

process

① word Embedding

② Encoding

③ Mean Pooling

④ Indexing

①-④

①-④

After process (①-④) for all sentences, now we have Index our dataset in the vector database

Query: "am I you"

process

①-④

a 2-d query vector

estimate similarity

Dot Products between query vector and database vectors.

They are all 2-d

(transposing the query vector for matrix multiplication)

Nearest Neighbor

Find the largest dot product by linear scan

In practice, because scanning billions of vectors is slow, we use an Approximate Nearest Neighbor (ANN) algorithm like the Hierarchical Navigable Small words (HNSW)

Data

how are you

who are you

who am I

Query

am I you

Vocabulary size

1
n

Word Embedding

a	an	the	how	why	who	what	are	is	am	be	was	you	we	I	they	...
0	-1	0	1	0	1	0	0	-1	1	0	0	0	3	1	0	
2	0	2	0	0	0	-1	1	0	0	0	2	1	0	2	0	
-1	0	-1	1	2	0	0	1	0	1	-1	0	0	-1	0	3	
0	1	0	0	1	0	1	0	1	0	1	-2	0	0	0	1	

Data:

1	0	0
0	1	1
1	1	0
0	0	0

1	0	0
0	1	1
0	1	0
0	0	0

1	1	1
0	0	2
0	1	0
0	0	0

Query:

1	1	0
0	2	1
1	0	0
0	0	0

Text Embedding includes ① & ②

- ① Encoding: Feed the sequence of words embedding to an encoder to obtain a sequence of feature vectors, one per word.
- ② Mean Pooling: Merge the sequence of feature vectors into a vector using "mean pooling" which is to average across the columns.

Reluc

- * In this Example **encoder** is a simple one layer perceptron (linear layer + Reluc)
- * In practice, the encoder is a transformer or one of its many variants (LLMs)
- * The result of text Embedding or sentence embedding is a single vector.
- * other pooling techniques are: CLS - mean pooling is the most common.

1	0	0
0	1	1
1	1	0
0	0	0

1	0	0
0	1	1
0	1	0
0	0	0

1	1	1
0	0	2
0	1	0
0	0	0

1	1	1
0	2	1
1	0	0
0	0	0

Query

DataSet after word Embedding (page 2)

Text Embedding

here: a simple one layer perceptron / In practice: LLMs

Encoder

Transformer-based Architecture

$$Z = Wn + b$$

$$\text{ReLU}(Z) = \max\{0, z\}$$

1	1	0	0	0
0	1	0	1	0
1	0	1	0	-1
1	-1	0	0	0

Weight bias

DataSet:

1	1	1
0	1	1
1	0	X
1	X	X

1	1	1
0	1	1
0	0	X
0	X	X

1	1	3
0	0	2
0	1	0
0	1	X

Query:

1	3	1
0	2	1
1	0	X
1	X	X

$$\begin{aligned}
 1 &= [(1 \times 0) + (1 \times 1) + (0 \times 0) + (0 \times 0)] + 0 \\
 1 &= [(0 \times 0) + (1 \times 1) + (0 \times 0) + (1 \times 0)] + 0 \\
 -1 &= [(1 \times 0) + (0 \times 1) + (1 \times 0) + (0 \times 0)] + -1 \\
 -1 &= [(1 \times 0) + (-1 \times 1) + (0 \times 0) + (0 \times 0)] + 0
 \end{aligned}$$

bias

$$-1 = [(1 \times 0) + (-1 \times 2) + (0 \times 0) + (0 \times 0)] + 0$$

Mean Pooling

\sum over columns
DataSet

3/3
2/3
1/3
1/3

3/3
2/3
0
0

5/3
2/3
1/3
1/3

Query:

5/3
3/3
1/3
1/3

$\frac{3}{3}$
$\frac{2}{3}$
$\frac{1}{3}$
$\frac{1}{3}$

$\frac{3}{3}$
$\frac{2}{3}$
0
0

$\frac{5}{3}$
$\frac{2}{3}$
$\frac{1}{3}$
$\frac{1}{3}$

$\frac{5}{3}$
$\frac{3}{3}$
$\frac{1}{3}$
$\frac{1}{3}$

(4)

Query

Dataset after Mean Pooling & Encoding
Text embedding of Dataset

Encoding



projection

1	1	0	0
0	0	1	1

Reducing the text embedding vector dimensions using projection matrix

a short representation of dataset sentences which allows faster comparison & retrieval

PxV

vector storage

$\frac{5}{3}$	$\frac{5}{3}$	$\frac{7}{3}$
$\frac{2}{3}$	0	$\frac{2}{3}$

PxV

$\frac{8}{3}$
$\frac{2}{3}$

Dot Product (Finding similarity)

$$\frac{40}{9} + \frac{4}{9}$$

$$\frac{40}{9} + 0$$

$$\frac{56}{9} + \frac{4}{9}$$

$$\frac{44}{9}$$

$$\frac{40}{9}$$

$$\frac{60}{9}$$

Retrieval

Performing linear scan to find the largest dot product. This process is slow for billions of values: using ANN / HNSW Algorithms

$\frac{8}{3}$	$\frac{2}{3}$
---------------	---------------

T
(Transpose)