

Vector Databases

- Traditional Databases are not always good for finding "Similar" Entities

- Resolution: Vector Database

Vector Representations (Embedding)

An embedding of an Entity is a high-Dimensional vector of values (Real values)

$$t_i \longrightarrow \vec{v}_i \in \mathbb{R}^d$$

such that given a distance function $\Delta(t_i, t_j)$

$$\text{dist}(\vec{v}_i, \vec{v}_j) \sim \Delta(t_i, t_j)$$

$\cos(\vec{v}_i, \vec{v}_j)$: Similarity function

Vector Database:

A Collection of Tables
that every table T is
in form $\{A_1, \dots, A_d\}, \{v_1, \dots, v_d\}$

Traditional
attributes

vector representation

Example:

Animal

ID	name	Category	Embedding
1	Dog	1	1...0...1
2	Cat	2	0...0...1
	\vdots		
n	Lion	2	11...1001

Simplified Vector DB
 \vec{v}

1	- - - - -	1
	\vdots	

Every row is an Embedding.

Challenges (Compared to DB)

1 - Vague Similarity:

How to Compare two objects?

2 - Expensive Comparisons

$O(d)$

3 - High Dimension

$d \sim (100 - 10000)$

4 - Lack of Structure:

- The Columns do not have any meanings independently.

⇒ Indexing Strategies based on ordering / Hashing are not effective.

5 - Different Goals, Target Users

- Traditional queries are initiated & understood by Human beings

- VectorDB ~ ~ ~
~ ~ by machine

→ How to Combine these perspectives?

Traditional DBs:

Query Types

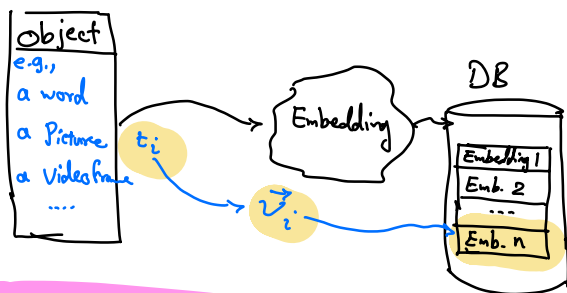
DDL	Create Table Create Schema
DML	Insert / Del / Update Select

Vector DB

Query Types

- Data Man. : Insert / delete
- Similarity Search
 - objects similar to a given query object
 - $f(D, q)$
 - ↳ query point
- Constrained (Predicated) queries
 - Find similar objects from the ones that satisfy a condition
 - ↳ Extended VDBs
- Multi-Vector Queries
 - Single vector → Multiple vectors
 - $M \rightarrow M$
 - $M \rightarrow S$

Insert:



Who Provides the details of the Embedding Transformation!




$$\vec{v}_i = f(t_i)$$

e.g.,

word Embedding: (word2vec)

Sentence	Embedding ($f(t_i)$)				
we	<table><tr><td>2.5</td><td>3</td><td>4.6</td><td>11</td></tr></table>	2.5	3	4.6	11
2.5	3	4.6	11		
love	<table><tr><td>1</td><td>1</td><td>2.5</td><td>9.2</td></tr></table>	1	1	2.5	9.2
1	1	2.5	9.2		
UIC	<table><tr><td>0.23</td><td>4.8</td><td>9.2</td><td>1.3</td></tr></table>	0.23	4.8	9.2	1.3
0.23	4.8	9.2	1.3		

$n \times d$ table of Embeddings

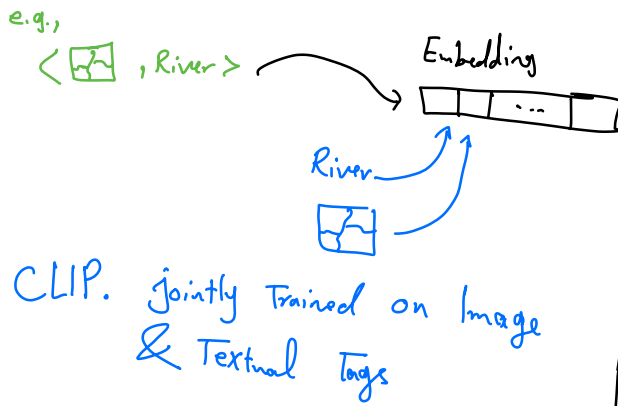
Image	Embedding ($f'(t_i)$)				
	<table><tr><td>1</td><td>2</td><td>...</td><td>12</td></tr></table>	1	2	...	12
1	2	...	12		
	<table><tr><td>3.8</td><td>9</td><td>...</td><td>10.9</td></tr></table>	3.8	9	...	10.9
3.8	9	...	10.9		
...	...				
	<table><tr><td>8.5</td><td>10.3</td><td>...</td><td>5.2</td></tr></table>	8.5	10.3	...	5.2
8.5	10.3	...	5.2		

$n \times d'$ table of embeddings

Embeddings in different "Embedding Spaces" are not comparable

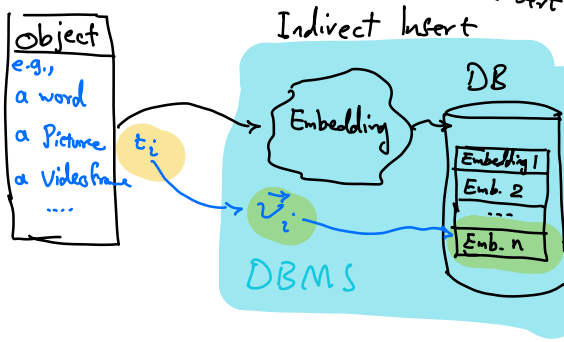
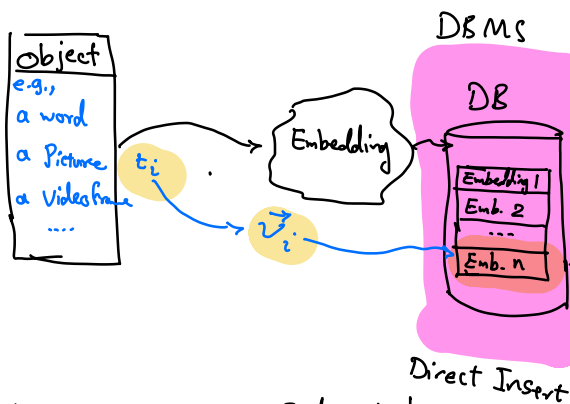
Jointly-Trained Embeddings

The Training Process is done on
 $\langle \text{Object}, \text{Tag} \rangle$
 ↑
 Textual



Insert Types

- Direct: The user is responsible for providing the Embedding
- Indirect: The DB is in charge of providing the Embedding



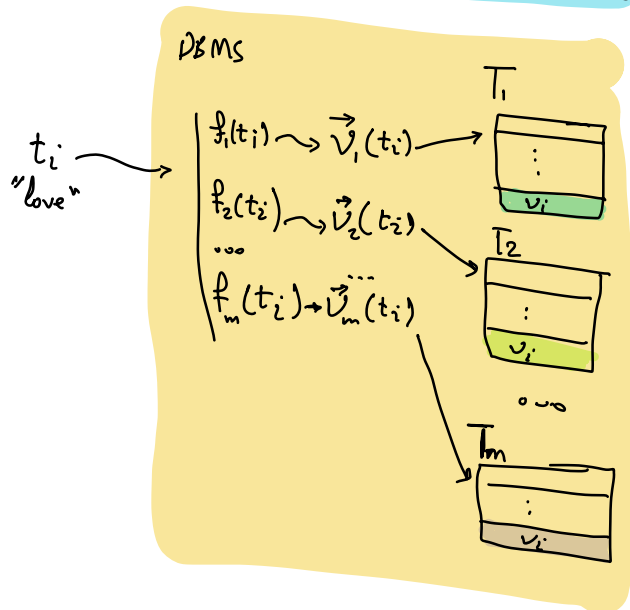
Direct Insert:

- issues
- ① Quality of the embedder left to the user.
 - ② The logical associations are Unknown.
 ↳ The DB cannot detect if the comparison is legit

- Benefits
- ① more freedom to user
 - ② Easier to adapt for DBMS

Indirect Insert:

Can use multiple Embedding



Similarity Metrics:

* I will use the words distance & similarity interchangeably

e.g., $\text{Sim}(a, b) = 1 - \text{dist}(a, b)$

Similarity Types

Metric:

Non-metric:

for a distance function Δ to be metric, it should satisfy the following properties

1- Identity: $\Delta(a, a) = 0$

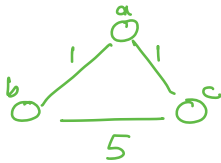
2- Positivity: $\Delta(a, b) \geq 0$

3- Symmetry: $\Delta(a, b) = \Delta(b, a)$

4- Triangular Ineq:

$$\Delta(a, b) \leq \Delta(a, c) + \Delta(b, c)$$

e.g.

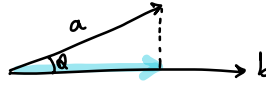


does not satisfy the Triangular Ineq

$$w(b, c) = 5 > 1 + 1$$

Distance / Similarity Measures

① Inner Product (Dot Product)



$$a = [1, 1, 3]$$

$$b = [2, 0, 1]$$

$$a \cdot b = \langle a, b \rangle = \sum a_i b_i$$

$$a \cdot b = 2 \times 1 + 1 \times 0 + 3 \times 1 = 5$$

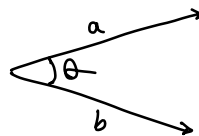
$$a \cdot b = \|a\| \cdot \|b\| \cos \theta$$

the projection of b on a

Not in metric

Identity	X
Positivity	X
Symmetry	✓
Triangular Ineq	X

Cosine Similarity



$$\text{Cosine}(a, b) = \cos(\theta)$$

$$\cos(a, b) = \frac{a}{\|a\|} \cdot \frac{b}{\|b\|}$$

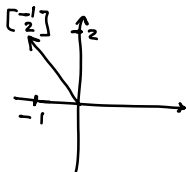
$$\|a\|_2 = \|a\| = \sqrt{a \cdot a}$$

$$\|a\|_p = (\underbrace{a \cdot a \dots a}_{p \text{ Times}})^{1/p}$$

e.g.,

$$a = [-1, 0, 3, 5]$$

$$b = [2, 1, -1, -3]$$



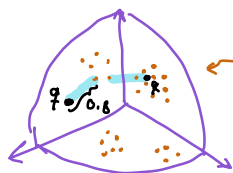
$$\cos(a, b) = ?$$

$$\|a\| = \sqrt{a \cdot a} = \sqrt{1+0+9+25} = \sqrt{35}$$

$$\|b\| = \sqrt{b \cdot b} = \sqrt{4+1+1+9} = \sqrt{15}$$

$$\frac{a}{\|a\|} = \frac{[-1, 0, 3, 5]}{\sqrt{35}}, \quad \frac{b}{\|b\|} = \frac{[2, 1, -1, -3]}{\sqrt{15}}$$

$$\begin{aligned} \Rightarrow \cos(a, b) &= \frac{[-1, 0, 3, 5]}{\sqrt{35}} \cdot \frac{[2, 1, -1, -3]}{\sqrt{15}} \\ &= \frac{(-2 + 0 + (-3) + (-15))}{\sqrt{525}} \\ &= \frac{-20}{\sqrt{525}} \approx -0.8 \end{aligned}$$



Points are on the surface of the sphere

Cosine is not Metric.

ℓ_p -norm measures

$$\|a-b\|_p = \left(\sum |a_i - b_i| \right)^{1/p}$$

ℓ_1 -norm

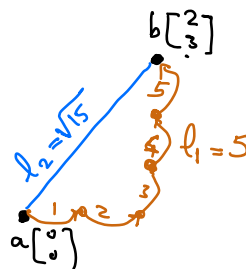
ℓ_2 -norm: Euclidean Distance

ℓ_0 -norm

ℓ_∞ -norm

$$\ell_1\text{-norm} = \sum_{i=1}^d |a_i - b_i|$$

$$\ell_2\text{-norm} = \sqrt{\sum (a_i - b_i)^2}$$



$$\ell_0\text{-norm} = \left(\sum |a_i - b_i|^0 \right)$$

= The number of non-zero values

ℓ_∞ -norm: The maximum absolute value of $|a_i - b_i|$

$$= \max |a_i - b_i|$$

e.g.,

$$a = [-1, 0, 3, 5]$$

$$b = [2, 1, -1, -3]$$

$$l_1 = \frac{(2+1)}{3} + \frac{(1-0)}{1} + \frac{(3+1)}{4} + \frac{(5-(-3))}{8} = 16$$

$$\begin{aligned} l_2 &= \sqrt{\frac{(2+1)^2}{9} + \frac{1^2}{1} + \frac{4^2}{16} + \frac{8^2}{64}} \\ &= \sqrt{90} = 3\sqrt{10} \end{aligned}$$

$$l_0 = \mathbb{I}(3, 1, 4, 8) = 4$$

$$l_\infty = \max(3, 1, 4, 8) = 8$$

L_p -norm is in Metric.