

Term Paper Performance Review Research

Language model optimized with least amount of data
Child language acquisition

But, Roberta — in traditional sense, they are not LLMs
(at least a billion parameter) ↩

Chatgpt — General task solvers

Benchmark of chatgpt in summarization datasets

General task solvers are better than specific task solvers

Performance of LLM in low resource language.

only good at previously seen data

Task contamination

SOTA

Class Notes

Edit distance : way to quantify string similarity

spelling mistake:

graffe

↙ giraffe graf grail

spell checking

Harvard President Claudine

Harvard university President Claudine

co reference



if two strings
contain same
entity

Define cost — LEVENSHTEIN

The apple was red

1. The apple was ^{2\$}big ^{2\$}and red. → 4\$

2. The ^{2\$}fish was ^{2\$}blue. → 8\$
apple 2\$ red 2\$

task	cost
i insertion	2\$
d deletion	2\$
s substitution	

1st was more similar to the string.

→ minimum Edit distance (check test similarity)

d s s - i s

cost = 8

INTENTION

→ EXECUTION

↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
x E x E C U T I O N

min Edit Distance?

→ Can start from anywhere

i = 1\$
d = 1\$
s = d + i

Use DP (Dynamic Programming)

— search problem

Find most cost efficient approach

Edit Distance Algo

- ✓ 1. Levenshtein
- ✓ 2. Hamming Distance
- 3. Jaro-Winkler

Token-based

- ✓ 1. Cosine-Similarity (similarity in text/may)
- ✓ 2. Jaccard Index
- 3. Tversky Index

Sequence-based (programming)

- 1. LCS (Longest common substring)
- 2. Ratcliff-Obershelp similarity

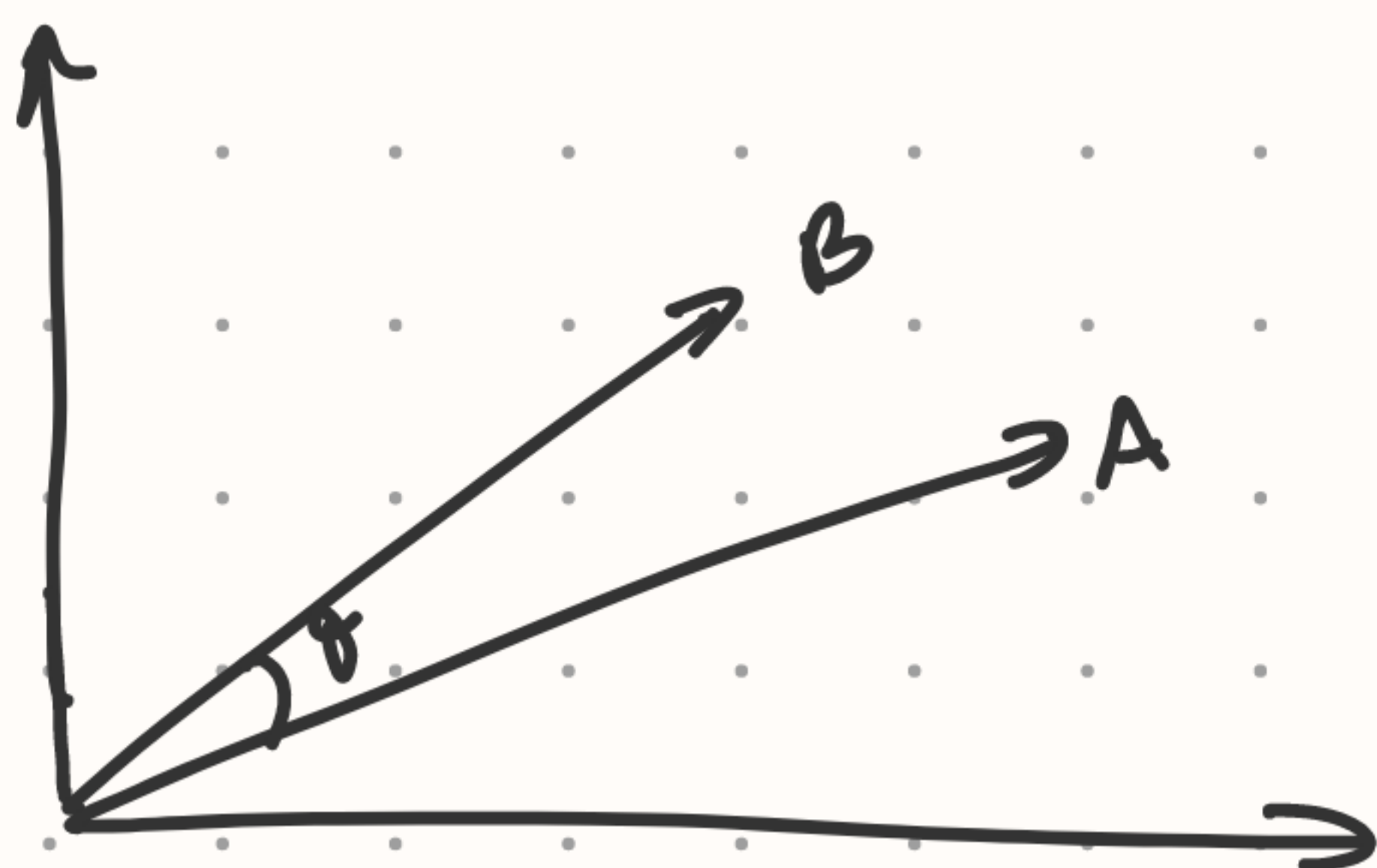
Cosine-Similarity

converts the whole corpus into same dim. vectors and projects it.

similar sentence \rightarrow similar vectors

more similar \rightarrow less distance

2. same sentence \rightarrow project on each other
 $\theta = 0$



cosine similarity (doc1, doc2)

Document 1: The cat in the hat.

Document 2: A black cat in a black hat.

Doc \rightarrow vectors : Find all unique tokens from two docs.

	The	cat	in	hat	a	black	
Doc 1	2	1	1	1	0	0	Frequency
Doc 2	0	1	1	1	2	2	

$$A = \begin{bmatrix} 2 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 2 \\ 2 \end{bmatrix}$$

cosine similarity (doc1, doc2)

$$= \frac{A \cdot B}{\|A\| \cdot \|B\|} = \frac{2 \cdot 0 + 1 \cdot 1 + 1 \cdot 1 + 1 \cdot 1 + 0 \cdot 2 + 0 \cdot 2}{\sqrt{7} \cdot \sqrt{11}}$$

$$= 0.342$$

similar \rightarrow close to 1

diff \rightarrow close to -1

\perp vectors \rightarrow close to 0
(no similarity / dissimilarity)

— complete ly diff set of words
(0 frequency)

There should be a (-ve) value theoretically.

Insight

- 1) The patient died before he doc came.
- 2) The doc came after the patient died.

Cosine similarity cannot detect semantic meaning of sentence.

(Cons)

computationally strong (Pros)

- Sensitive to sentence long.
long vector (not necessary) | Cons

sparsity - having a lot of 0s in vectors/matrices
cosine similarity create sparse vectors that will dominate the entire value.

Jaccard Index is more sensitive to context in finding text similarity.

$$J.I = \frac{A \cap B}{A \cup B}$$

$$\text{Jaccard similarity}(\text{doc1}, \text{doc2}) = \frac{\text{doc1} \cap \text{doc2}}{\text{doc1} \cup \text{doc2}}$$

if 2 docs talking about same context

docs \rightarrow tokens

$$\text{J.S.} = \frac{\{ \text{'The', 'cat', 'in', 'the', 'hat'} \} \cap \{ \text{'A', 'black', 'cat', 'in', 'a', 'black', 'hat'} \}}{\{ \text{'The', 'cat', 'in', 'the', 'hat', 'A', 'black', 'cat', 'in', 'a', 'black', 'hat'} \}}$$

$$= \frac{3}{6}$$

$$\{ \text{'the', 'cat', 'in', 'hat', 'a', 'black'} \}$$

$$= \frac{3}{6} \approx 50\% \text{ similarity.}$$

Imp. of tokenization/lemmatization.

throw, threw, thrown \rightarrow w/o lemmatization, consider them diff.
(verb forms)