

Vector Models and Text processing

10.03.2023

Munich



Definitions

Vocabulary = set of all words

Corpus = Collection of writings (dataset)

N-gram = Sequence of N consecutive items

1 = Unigram 2 = Bigram 3 = Trigram

Bag of words

- Text is sequential

- Order doesn't matter
used in
Vector Models

Dog Toy = Toy Dog
in Bag of words

Count Vectorizer

- instance of Bag of words, order doesn't matter

Document - file, sentence, full book

→ convert sentences into numerical representation

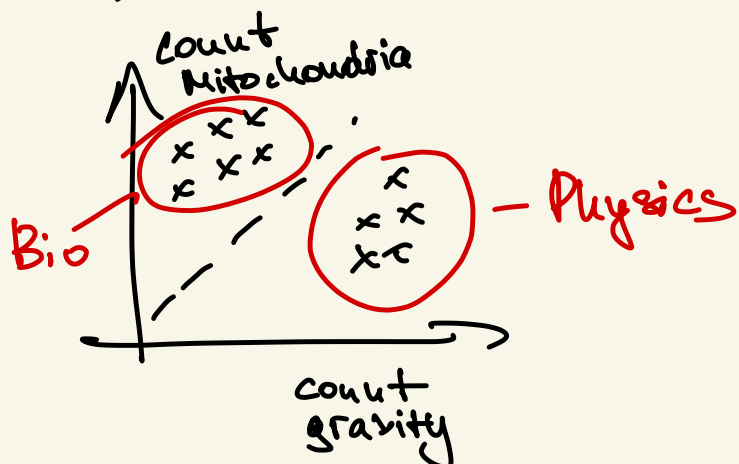
Vocabulary size → every word has a number

I like eggs

I hate cats

I like eggs and I like cats

and	cats	eggs	hate	I
0	0	1	0	1
0	1	0	1	1
1	1	1	0	2



Practical issues

Tokenization - convert string containing words into array of numbers

Normalization

long vs short documents $\hat{x} = \frac{x}{\|x\|_2}$ $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$
sparse matrices to be used

Tokenization

x old approach \rightarrow .split()

x String split works well for RNN classification

x Punctuation I hate cats. vs I hate cats?

CountVectorizer ignores punct. string split

x Casing Cat name vs cat animal

x Accents naïve vs naïve
count vectorizer strip accents

x Character Based

1M words = 1M embeddings = 1M last layer
not much information, but 26 letters + punctuation

Subword - Based Tokenization

walking \rightarrow walk + ing \rightarrow modifier

walking \leftrightarrow walks \leftrightarrow tree = equal weights

CountVectorizer has word/char based tokenization

Stop words

- very common words with less meaning (the, it, and, is)
- increasing vector dimensionality
- distance will be diminished (between vectors)

- Count Vectorizer has english stopwords library
- NLTK has stopwords from other languages

Stemming and Lemmatization

- Basic tokenization: *walk, walking, walks*
 ↳ very high dimensional
 all different tokens
- Search engine will treat all words independently
- Stemming → chop the end of the word (crude)
- Lemmatization → true root of the word (sophisticated)

Reduces vocabulary size

Stemming

sses → remove es: Bosses → Boss
 Replacement → Replace

Porter stemmer in NLTK

Lemmatization

- Look up table
 Better → Good
 Was → Be
 Is → Be
 Mice → Mice
 - Wordnet DB in NLTK
 - pos (parts of speech) signifies whether word is noun, adj, verb
 ↳ default is noun
- Trump has devoted following *noun*
 The cat is following me. *verb*

- NLTK can do POS tagging
- Application: search engine, doc. retrieval, social media tags, ads

- Search by converting terms in root form
- Online ads keyword \rightarrow match ads to keywords
Google Pixel \rightarrow show Apple iPhone

Vector Similarity

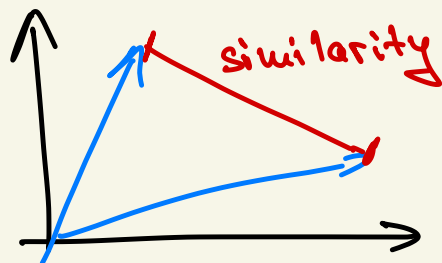
$$a = [0.5, 0.3, \dots] \rightarrow \text{Similarity Function} \rightarrow \text{score} = 0.12$$

$$b = [-0.1, 0.2, \dots] \rightarrow S(a, b)$$

- Find similar documents on specific topics

- If words were vectors: queen \leftrightarrow king
car \leftrightarrow automobile } high similarity
used in article spinning

- Euclidean Distance
similarity = distance



$$\|x - y\|_2 = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_D - y_D)^2}$$

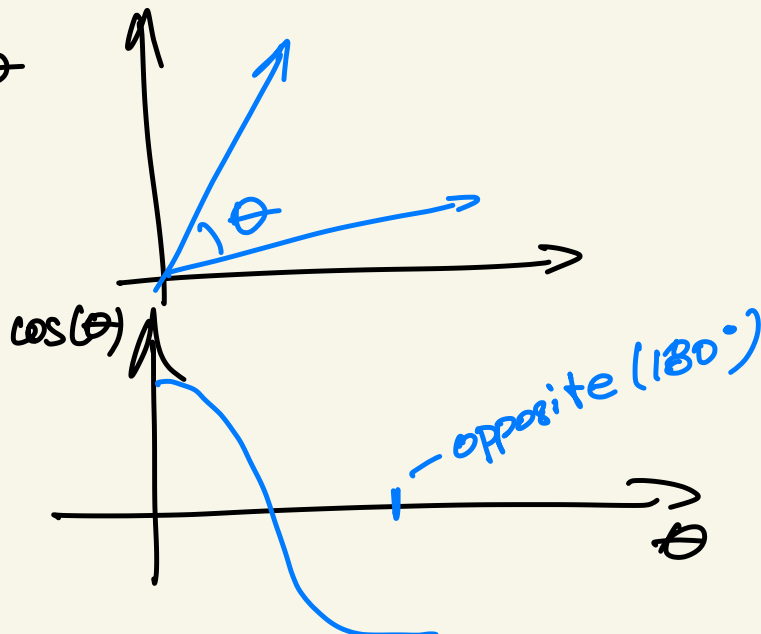
monotonically increasing

- Angle Between 2 Vectors (cosine similarity)

$$x^T y = \|x\|_2 \|y\|_2 \cos \theta$$

$$\cos \theta = \frac{x^T y}{\|x\|_2 \|y\|_2}$$

$$\|x\|_2 = \text{magnitude}$$

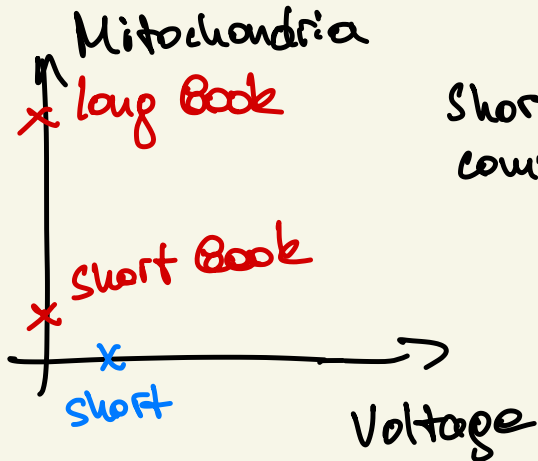


Cosine Distance = $1 - \text{Cosine Similarity}$

dist = $1 - (-1) = 2$ for 180°

dist = $1 - 1 = 0$ for 0°

Which to use?



Short Books have smaller eucl. distance compared to Books on Biology topics

- If vectors are normalized in L2 form, distance can be used

$$\|x - y\|_2^2 = (x - y)^T (x - y)$$

$$= x^T x - 2x^T y + y^T y \quad \text{But } x^T x = y^T y = 1 \text{ for L2}$$

$$= 2 - 2x^T y \quad x^T y = \frac{x^T y}{\|x\|_2 \|y\|_2} = \cos \theta$$

$$= 2 - 2 \cos \theta$$

$$= 2 \cdot \cos \text{ distance}$$

Word to Index Mapping

	word 1	word 2	word 3	
Doc 1	0	1	1
Doc 2	2	2	0	...
Doc 3	1	1	1	...

Vocabulary size $N = 5$

$0 \rightarrow I$ $1 \rightarrow \text{like}$ $2 \rightarrow \text{cats} \dots$

current_idx = 0

word2idx = {}

for doc in docs:

tokens = word_tokenize(doc)

for token in tokens:

if token not in word2idx:

word2idx[token] = current_idx

current_idx += 1 → map token to index
value is index

Word2Index → Count Vectorizer → TF IDF

- Words not appearing in test set?

if frequency in train set is low → ignore or assign $\frac{1}{\text{count token}}$

- Reverse Mapping

input # 345 → word

Neural Word embeddings

embedding = vector

convert words into vectors, doc = sequence of vectors

CNN / RNN / Transformers are built for sequences

Word2vec

neural network, weights are the embeddings

NN goal is to predict if a word is in the context of another word

... the fox jumped over the lazy dog ...

context window

Jumps → NN → Fox (✓)
→ Brown (✗)
→ Dog (✗) → weights are the neural embeddings vector

Glove

like recommender system

The quick brown fox jumps over the lazy dog.

score = $\frac{1}{2}$ score = $\frac{1}{4}$

1 word away \Rightarrow score = 1

Ratings are based on distance
then: how high a word rates another word

Word Vectors

- Embeddings are low-dimensional and dense
- Doc \rightarrow Tokenize \rightarrow word vectors \rightarrow Average(vectors)
 \hookrightarrow 1 vector of the same size for a Doc
- Word Analogies
arithmetic: $V(\text{king}) - V(\text{queen}) = V(\text{man}) - V(\text{woman})$
 $x = V(\text{king}) - V(\text{man}) + V(\text{woman})$
France : Paris Italy : Rome countries
Japan : Japanese languages
Miami : Florida cities
Dec : Nov months
- Learns meaning from word pos. in vector space
weights are close to optimal, no need for training