## NLTK ?
└── natural language toolkit

→ Useful for working with human language data

- Tokenization    · Stopword removal    · Lemmatization
- POS tagging    · Parsing

<u>punkt</u> : Segments text into sentences using an unsupervised
          algorithm

[ <u>Working</u> :
└→ learns abbreviations, ~~collectio~~ collocations, and sentence
    boundaries from raw text.

Eg :

    $X$ = "ABC. He uses NLTK hamesha !"
    print (sent_tokenize ($X$))
    o/p: ["ABC.", "He uses NLTK hamesha !"]
                    ↓
         uses Punkt (Pre-trained model) to identify sentence
         boundaries

<u>Use Case</u> : preprocessing in tasks like sentiment analysis,
          summarization, or any sentence-level NLP

<u>Stopwords</u> : common words that are usually filtered out
            before processing

करी ECिरे ?

~~vork~~ words like "the", "is", "and" don't carry much
semantic weight and can introduce noise in models

Eg : from nltk. corpus import stopwords
hatao_bc = set (stopwords.words('english'))
tokens = ["this", "is", "a", "sample", "sentence"]
filter = [ word for word in tokens if word not in hatao_bc]
print (filter)        # o/p: ['sample', 'sentence']

**WordNet** : Used for lemmatization and semantic analysis

↳ A large English lexical db where:
- noun, verbs, adjectives are grouped into sets of synonyms

Eg: from nltk. stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
print ( lemmatizer. lemmatize ("running, pos="v"))
#o/p: run
↑ verb

**Q) Lemmatization क्या है ?**

→ Text normalization technique in NLP, ~~convert~~
→ Converts a word into its base or dictionary form (its lemma)
→ This will result to a meaningful root word like "good" from "better"

Eg: running → run (verb)
better → good
was → be

**Purpose** : reduce variations of words to a ~~comm~~ common base ~~for~~ form, making text data more uniform and reducing its complexity for analysis.

**Q) Semantic analysis** : Understanding the meaning, intent and emotional tone of text by examining the relationships and context b/w words and phrases.

अब दोनो को साथ मे use कैसे ~~करे~~ करते है ?

Lemmatization ── converts ──→ "running" to "run"

Now this will help semantic analysis model recognize them as the same concept. This consistency improves the model's ability to correctly interpret the overall meaning of sentence.

<u>punkt_tab</u> : internally used by Punkt tokenizer to store
learned sentence boundary parameters /
To train own Punkt tokenizer on a custom corpus

## TF-IDF (Term Frequency - Inverse Document Frequency)

→ Frequency of a word in a document.

→ Penalizes common words across documents.
o/p: Sparse matrix where each row is a document
and each column is a weighted word feature

$$TF(t,d) = \frac{(No.\ of\ occurance\ of\ term\ t\ in\ document\ d)}{(Total\ no.\ of\ terms\ in\ the\ document\ d)}$$

$$IDF(t,D) = \log_e \frac{(Total\ no.\ of\ documents\ in\ the\ corpus)}{(no.\ of\ documents\ with\ term\ t\ in\ them)}$$

$$TF-IDF(t,d,D) = \Big[TF(t,d)\Big] \times \Big[IDF(t,D)\Big]$$

Word indexes: $\{$ 'geeks': 1, 'for': 0, 'r2j': 2 $\}$
Doc index →

tf - idf value:

| | |
|---|---|
| (0,0) | 0.5493 — |
| (0,1) | 0.8355 — |
| (1,1) | 1.0 |
| (2,2) | 1.0 |

→ tf-idf value of a word
having index 1
i.e. geeks in
document index 0

Document   word
index      index

Eg: Doc1 : The cat sat on the mat.
Doc2 : The dog played in the park.  _/_/_
Doc3 : Cats and dogs are great pets.

In doc 1, word "cat" appears→ 1
total no. of terms in doc 1 → 6
So, TF (cat, Doc 1) = 1/6

In doc 2, TF (cat, Doc 2) = 0
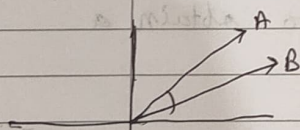In doc 3, TF (cat, Doc 3) = 1/6

(Doc 1 & 3)

Total no. of doc in the corpus = (D) ∴ 3
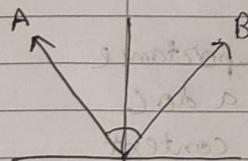no. of doc containing term "cat" : 2 ↗

$$IDF \ (cat, D) = \log\left(\frac{3}{2}\right) \approx 0.176$$

Cosine similarity : measure similarity b/w
two vectors using cosine of the angle b/w then

$$similarity = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^{m} A_i B_i}{\sqrt{\sum_{i=1}^{m} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}}$$



Similar

"hi, world!"

hi

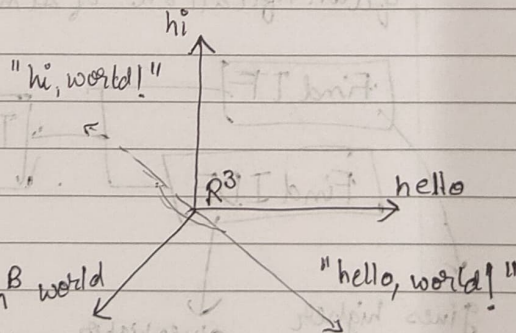R³        hello

A    B    B world        "hello, world!"
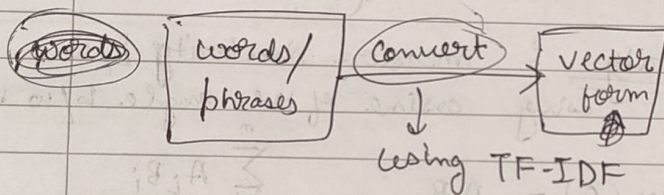
A

Unrelated        Opposite

Cosine similarity is a metric used to measure how similar the documents are irrespective of their size. Mathematically it calculates the cosine of the angle b/w two vectors
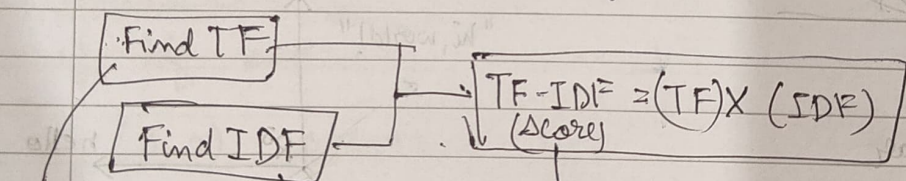
## Advantage

even if the two similar doc are far part by the Euclidean distance $\frac{t}{to}$ coy of the size (eg: word "yeah" could appear 50 times in one doc and 10 times in another) they could still have a smaller angle b/w them.

Smaller the angle = ↑ similarity

## Working



words → words/phrases → Convert → vector form

using TF-IDF

→ The vector presentation of the doc can then be used within cosine similarity formula to obtain a quantification of similarity

Find TF

Find IDF

$$TF-IDF = (TF) \times (IDF)$$
(score)

| gives higher score to words that are central to doc's theme | gives higher score to term that are uncommon across the whole dataset | indicates the importance of a term to a doc within the context to the entire corpus. Terms that are frequent in a doc but rare in the overall dataset recieve highest TF-IDF score |

→ For each doc, a vector is created where each dimension represents a unique word from the __/__/__ entire vocabulary of the corpus.

→ The value at each dimension (pos) in the vector is the TF-IDF score for that specific word in that doc

→ This process results in a numerical vector for each doc, allowing mathematical comparison.