

Natural Language Processing (NLP)

1. Search Engine

Data crawling (Nudge)

Spider, Robot.txt

DFS → Unstructured

Preprocessing

Tags, - Removed

Case -

Special character

Punctuation

Stop words

emoji

emoticons

Lemmatization, Stemming ✓

Tokenization

Notebook
Function
Framework

Structured Format

Indexing

Unique word

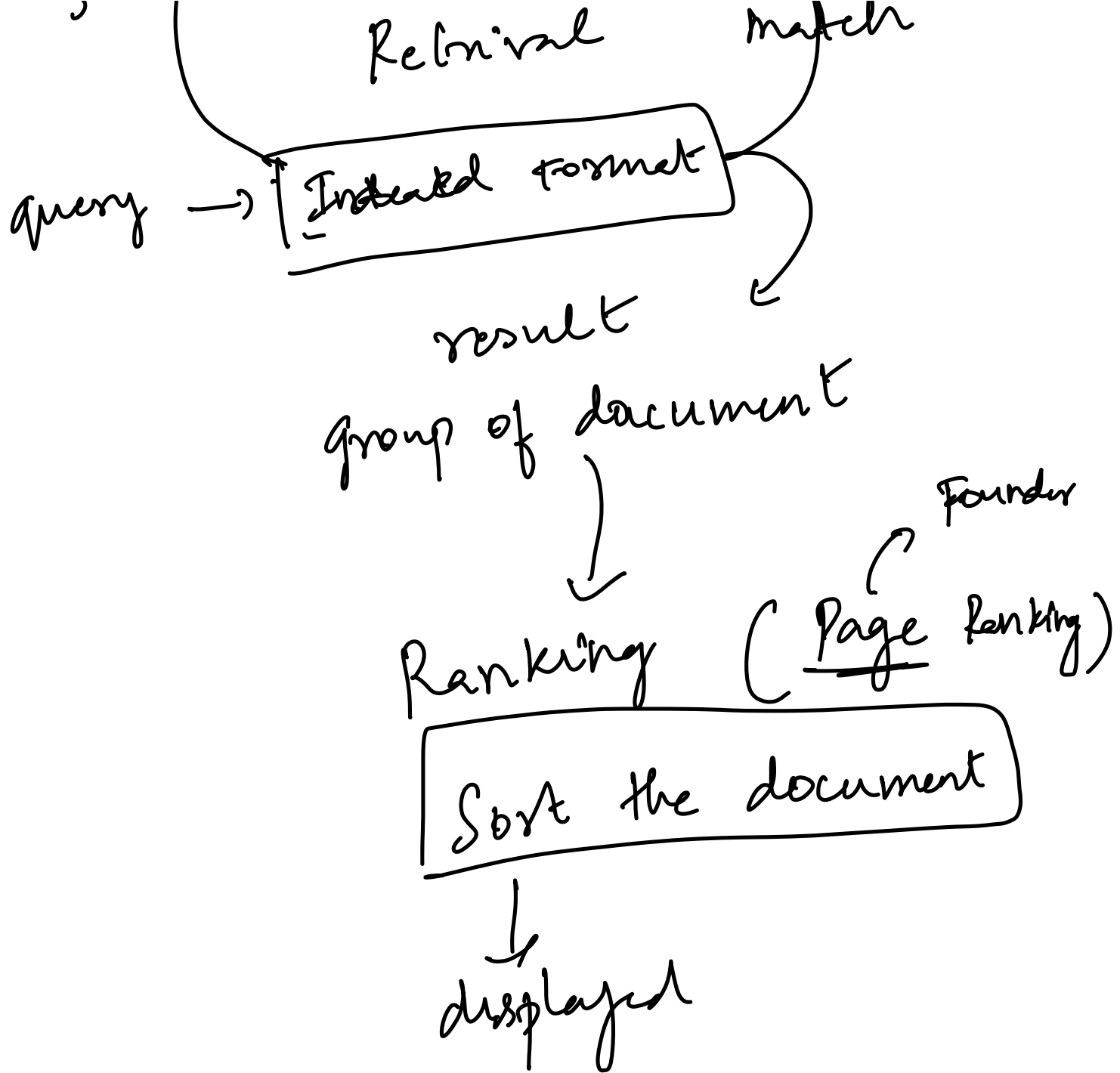
How many

Where it can

↓
NoSql (Cassandra, Hbase) RDBMS

n-grams
uni
bi
tri

Spider



2. Pre processing :

Textblob

NLTK

Spacy

Gensim

Beautiful Soup

3. Similarity

Euclidean distance
flaw \rightarrow Text

\rightarrow Magnitude

\rightarrow match is having
preference with
document word
count

Cosine similarity:

Angle

angle 0°

\Rightarrow

1

\Rightarrow

100% similar

angle 90°

\Rightarrow

0

\Rightarrow

0% similar

4. Recommendation Engine

Content based Recommender

Collaborative Filtering engine

Item based

User based

→ Apriori Algorithm

5. Text to Numbers.

1. One hot
2. Bag of words
3. TFIDF
4. Word Embedding

One hot

1. No Frequency of word in Documents corpus
2. No context
3. No order

Bag of words:

1. No Frequency of word in Corpus
2. No context
3. No order

TFIDF:

1. No context
2. No order

Word Embedding

word2vec

Skipgram

Continuous bag of words (CBOW)

disadvantage

→ generate → Train

1. Static embedding

word2vec

Glove

FastText

→ HW

6. Text Classification

Text to Number

ML model

DL model

LSTM algo

CNN algorithm

7. Sentiment Analysis.

Lexicon :

1. Textblob
2. Vader
3. Afinn

Text Classification

1. ML
2. DL
3. LSTM

4. Pretrained

8. Text Clustering

k-means → similar bucket

9. Topic Modelling

LDA → probabilistic Model

10. Seq2Seq

Encoder — Decoder

Machine Translation App

11. Transformer :

Drawback of Seq2Seq :

1. Input is taken Token by Token
2. The Decoder can only start after encoder completes the whole document

The encoder takes time to create the Context vector based on the size of document

3. If we have large token size Context Vector can miss the context of whole
4. Static embedding for token
bank, bank
track, track

Transformer Architecture

Attention is all you need

↳ Google Researchers

↳ Revolutionized

Encoder

BERT

Decoder

GPT

1 Decoder

GPT-1

GPT-2

GPT-3, 3.5

Creating the GPT's model

Prompt Engineer

Langchen

2202

Lambda index

Vector Space

pure love

DB

MOE \rightarrow Mixture of experts \rightarrow mutual

Lang Graph

Bedrock Agents

tools to use on Agents

Olama

Strands Agents

MLP

Optional based on time

Image Generation

GAI