WWW

Crawling

Spider

Robots.txt ✓

Unstructure | Structure

BigData HBASE → Preprocessing → Indexing → Transaction DB
Lookup Table

US
US
The the
hamlet Contradiction
ve

'1' Lower Case
10. Wrap ( space)
9. Contradiction
8. HTML, XML
1. Non-Alpha Numeric
2. Stop words ✓
3. Tagging POS, tag
4. Tokenization ✓
5. Stemming ✓
6. Lemmatization ✓
7. ascextend words

â ä ä

running, run, ran, (run) → Indexed tokens

User → Sort → Final Result
query → Page Ranking
Preprocess i/p → Result

Founder of Google
Larry (Page) + (Ranking) → Page Ranking

1. Hi how are you (!) bro        hi how you bro
2. I am good bro                 # good bro
3. I am very very good           | (very very) good

| | hi | how | you | bro | i | good | Vry | |
|---|---|---|---|---|---|---|---|---|
| (1) | 1 | 1 | 0 | (1) | 0 | 0 | 0 | indexing |
| (2) | 0 | 0 | (0) | (1) | 1 | 1 | 0 | |
| 3 | 0 | 0 | 0 | 0 | 1 | 1 | 2 | |

bro

Millions of website
tokens
Index

million rows

result optimized ↑ algorithm AI

Dhoni won
worldcup

sorted → Page
Ranking          query →

---

Limitation
          Frequency

1. Count is not captured ✓

2. Order is not captured ✓

word₁ : { 0, 2, 5, 6 } }
word₂ : { 1, 3, 5, 6 }
word₃ : { 2, 5, 8, 9 }

— | movie is very bad   ✓  ✓ —
— | movie is very very bad ✓  — ve

3. Context is not captured

4. Meaning is not captured

5.

1. Count vectorizer ✓ One hot

Doc1 = I like cats. do you like?
Doc2 = I like dogs to

Vocabulary = [ I, like, cats, dogs, too, do, you]

| | I | like | cats | dogs | to | do | you |
|---|---|---|---|---|---|---|---|
| doc1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 |
| doc2 | 1 | 1 | 0 | 1 | 1 | 0 | 0 |

Structured

2. Bag of words

| | I | like | cats | dogs | to | do | you |
|---|---|---|---|---|---|---|---|
| doc1 | 1 | 2 | 1 | 0 | 0 | 1 | 1 |
| doc2 | 1 | 1 | 0 | 1 | 1 | 0 | 0 |

Structured

# 3. Term Frequency - Inverse Document Frequency

## TF-IDF

$$\text{DOC1} = 1 \quad \text{I like cats. do you like?}$$
$$\text{DOC2} = 1 \quad \text{I like dogs to}$$

| Term Frequency | DOC1 ⇒ 5 | log |
|---|---|---|
| TF(DOC1) | DOC2 ⇒ 4 | IDF |
| | TF(DOC2) | |

| | TF(DOC1) | TF(DOC2) | IDF |
|---|---|---|---|
| I ✔ | 1/5 | 1/4 | 2/2 = 1 ⇒ 0 |
| like | 2/5 | 1/4 | 2/2 = 1 ⇒ 0 |
| cats ✔ | 1/5 | 0/4 | 2/1 = 2 = -3 |
| do | 1/5 | 0/4 | 2/1 = 2 = 0.3 |
| you | 1/5 | 0/4 | 2/1 = 2 = 0.3 |
| dogs | 0/5 | 1/4 | 2/1 = 2 = 0.3 |
| too | 0/5 | 1/4 | 2/1 = 2 = 0.3 |

## Inverse document Frequency

$$\left( \frac{\text{Total number of documents}}{\text{Number of doc which has word}} \right)$$

$$IDF \Rightarrow log \left(\text{Tot nos doc / Nos of doc which has word}\right)$$

(TF·IDF)

$$\Rightarrow \quad \frac{DOC1}{DOC2} = \begin{array}{l} I \\ I \end{array} \begin{array}{l} \text{like cats. do you like ?} \\ \text{like dogs to} \end{array}$$

6.3

$l^{th}$

TF·DF

DOC 1

TF·IDF (I') => $\boxed{1/5} * \boxed{0} = 0 \checkmark$

TF·IDF (like) = $2/5 \quad * \quad 0 = 0 \checkmark$

TF·IDF (cats) = $1/5 \quad * 0.3 = \boxed{0.06}$

TF·IDF C

✳. TF·IDF

$$= TF * log(IDF) \checkmark$$

$$= TF * \left(1 + log(IDF)\right)$$

why?

1000 document

10000 The

$$\left(\frac{1000}{1000}\right) = 1$$

$$\log(1) = 0$$