

2) What is Stemming? Explain Porter Stemming Algorithm.

Stemming:

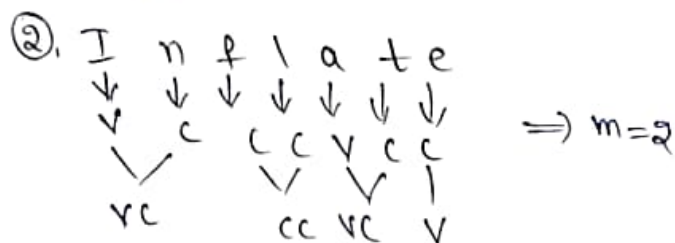
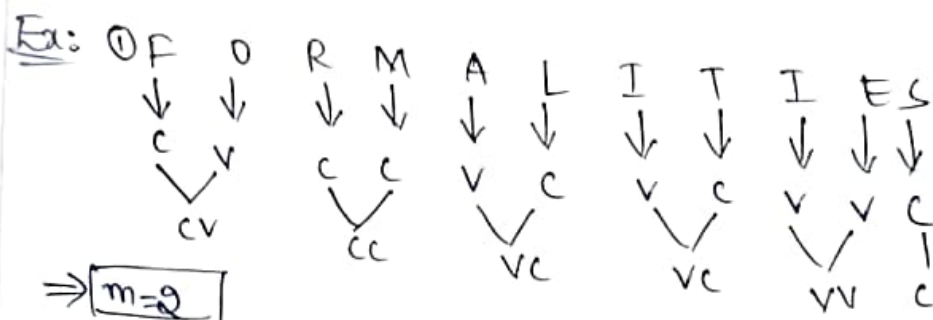
- It is a simple rule based approach to remove suffixes.
- It reduces words to their base form.
- It eliminates suffixes from words.
- It is used to improve recall.
- Also used to reduce index size.
- It is used to enhance query matching.

3

Porter Stemming

1. $\text{measure}(m)$: counts Vowel-Consonant (VC) Sequences
2. $\langle x \rangle$: stem ends with x .
3. v : stem contains vowel
4. d : stem ends in double Consonant.
5. 0 : stem ends in consonant - Vowel Sequence (excluding w, x, y)

$\text{measure}(m)$.

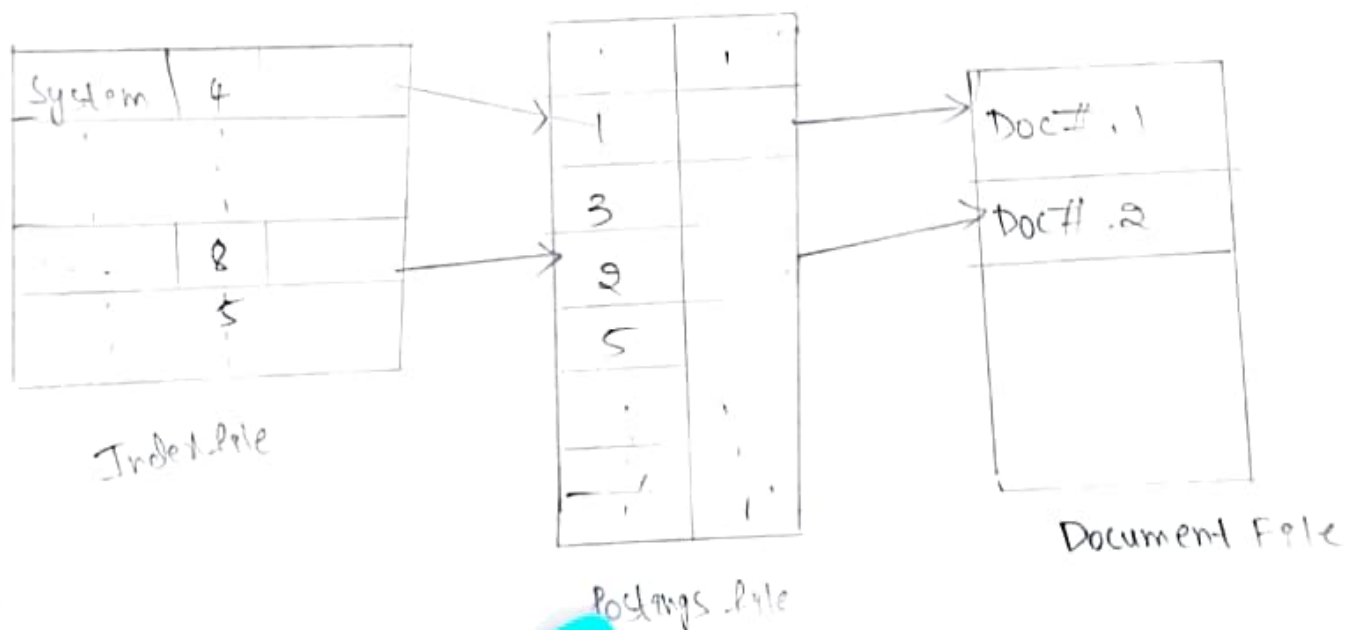


Step	Condition	Suffix	Replacement	Example
1a.	NULL	sses	ss	Stresses → stress
1b.	*v*	ing	NULL	making → mak
1b1.	NULL	ate	at	Inflate → Inflat
1c.	*v*	v	I	happy → happi
2.	$m > 0$	aliti	al	Formalities → Formal
3.	$m > 0$	Icate	IC	Duplicate → Duplic
4.	$m > 1$	Able	NULL	Adjustable → Adjust
5.	$m > 1$	e	NULL	Inflate → Inflat

5) Explain about Inverted File Structure with example.

Inverted File Structure

- It is a data structure, it allows efficient, full-text searches in the database.
- It stores a mapping of words to their locations in the database table or document.
- Inverted file based on methodology of storing an inversion of documents.
- For each word a list of documents in which the word is found is stored.
- Each document is given a unique & numeric identifier that is stored in inversion list.
- It is a data structure used in IRS to organize data and allow for efficient full-text searches.
- Each term is associated with a list of document identifiers.



Inverted File Structure

- An inverted file structure is a data organization method.
- This data organization method used in IRS
- Also known as an inverted index

Key components

1. Index Terms: Unique words or phrases
2. Posting lists: Documents containing each term.

Benefits

- Fast querying
- Reduced storage
- Improved scalability

Types

- Simple Inverted Index
- Compressed inverted index
- Levelled inverted index

Applications

- Search Engines
- Digital Libraries
- Database Systems

Example:

Document Collection

Doc ID	Text
1	"The quick brown fox jumps"
2	"Brown cats are sleepy"
3	"Foxes are quick animals"

ram)

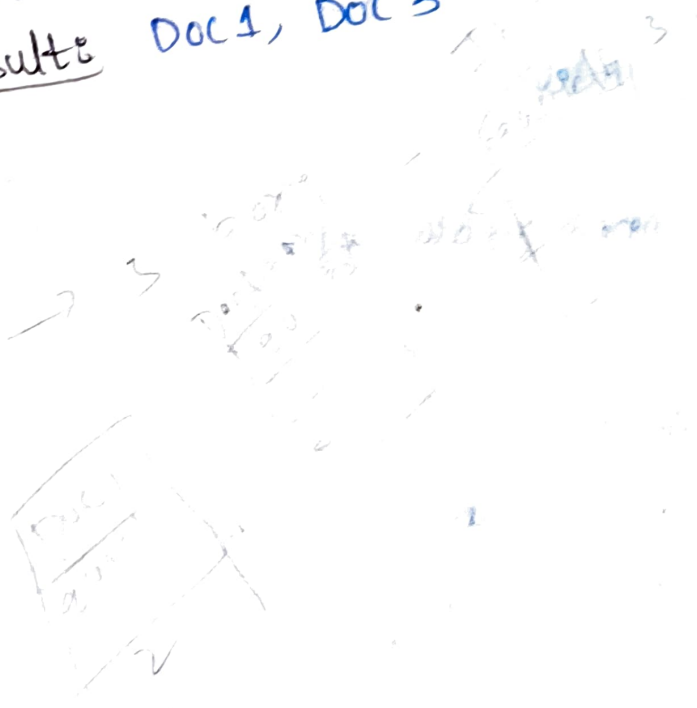
am?

Inverted File Structure:

Term	Posting list
quick	1,3
brown	1,2,3
fox	1,3
jumps	1
Cats	2
Sleepy	2
animals	3

Query: "quick brown fox"

Results: Doc 1, Doc 3



Finally, T_5 is selected as a new seed because it is
for viz., cluster 3.

Unit-2 * N-gram Data Structure

- N-gram is one of the data structure special technique for
- It is considered as a special technique for stemming.
- It is not concerned about the semantics (meaning) of the word.
- ~~Instead, they depend upon the fixed consecutive series of "n" characters.~~
- N-gram data structure can be viewed as Unique.
- N-grams are fixed length consecutive series of "n" characters.

Specializations:

- Special Data structure
 - Ignore words
 - Ignore sentences
 - N-gram = N-Length
 - Input as continuous data
 - Logical linkages
- Special Data Structures:
- It is one of the special data ~~structure~~ structure.
- It is unique.
- Ignore words:
- It ignores words, repeating a word once or twice.
- N-gram:
- Indicating as N-gram is equal to N-Length.

For $n=1$, 1-gram (Unigram)
 For $n=2$, 2-gram (Bigram)
 For $n=3$, 3-gram (Trigram)
 ...
 n-gram

Example: Hello How are you today

1-gram \rightarrow "Hello" "How" "are" "you" "today"

2-gram \rightarrow "Hello How" "are you" "you today"

3-gram \rightarrow "Hello How are" ~~How are you~~
 "are you today"

\rightarrow 1-gram (no word history)

\rightarrow 2-gram (one word history)

\rightarrow 3-gram (two word history)

$$\text{Prediction}(w_i/w_{i-1}) = \frac{\text{Count}(w_{i-1}, w_i)}{\text{Count}(w_{i-1})}$$

$w_{i-1} = \text{do}$
 $w_i = \text{Am}$

Ex: / Do

I am Amet

I like computer

Do Amet like computer

Amet I am

Do I like Amet

Do I like computer

I do like Amet

Next word	Frequency	Probability
I	6	2/4
Am	2	0/4
Amet	5	1/4
Like	5	1/4
Computer	3	0/4
do	4	0/4

$$= \frac{\text{count}(\text{do}, \text{I})}{\text{count}(\text{do})} = \frac{2}{4} = \frac{\text{do, Am}}{\text{do}}$$

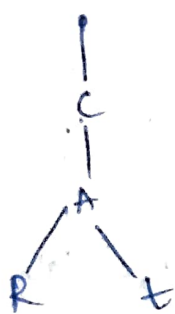
$$P = (w_1/w_{q-2}, w_{q-1}) = \frac{\text{count}(w_{q-2}, w_{q-1})}{\text{count}(w_{q-2}, w_{q-1})}$$

~~P = f(w)~~

* PAT Data Structure

- The name PAT comes from PATRICIA trees, which are used to search text.
- PAT is one of the data structure.
- It is also known as PAT tree or PAT array.
- Input data is transformed into searchable data structure.

Ex: ① Cat
Car

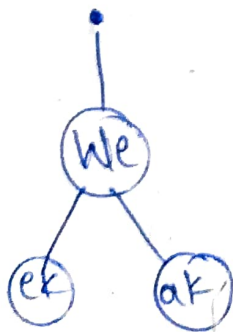


Regular PAT Tree



Compact PAT Tree

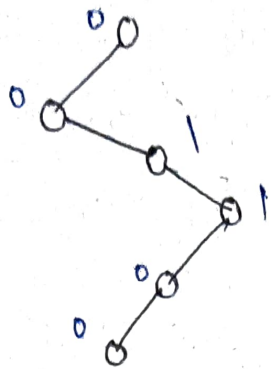
②. Week, weak



Binary Representation

ex: 001100

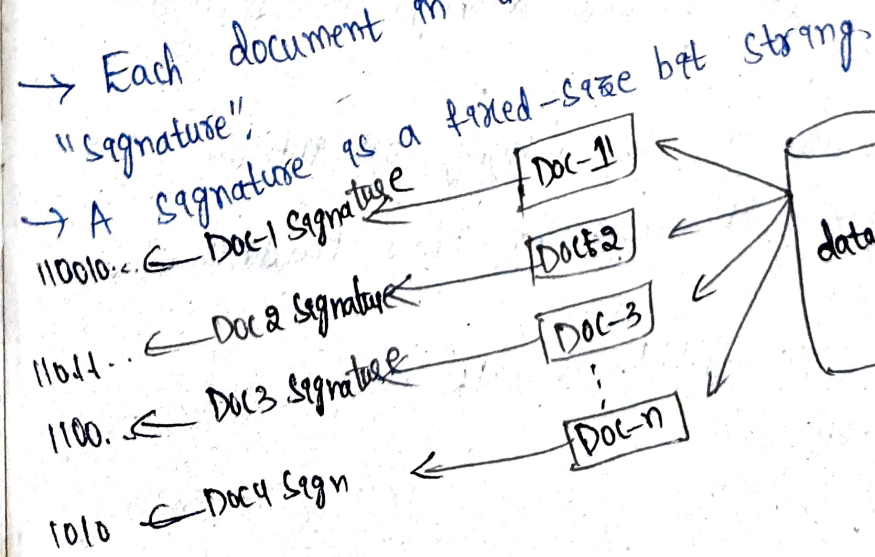
- * 0 - left side
- * 1 - right side



- PAT trees are created from the given input string.
- PAT trees supports strings and binary numbers.

* Signature File Structure

- The main aim of Signature file structure is to perform a fast test
- So that the unnecessary items that are not related to a query are removed.
- File structure is highly compresses & unordered.
- It requires significantly less space than an Inverted file structure.
- It is a method used for indexing & searching text documents
- Each document in the database is represented by a "signature".



→ By using hash

Input
(Document)

Hash
function

Output
(Document Signature)

Ex: this is a string → 0001000

~~Ex: Computer~~

User →

Query

Computed

1110000110

Doc

1010100110

AND
1110000110
1010100110

10100001100

→ match found

→ If we get all zero, then match not found.

* Hypertext & XML Data Structure

→ Hypertext is a data structure.

→ Basically, Internet is a global information network - which introduced a new storage data structure called Hypertext.

→ This Hypertext data structure display information on a WEB Page.

→ this structure is used largely on internet.

→ It is different from the other traditional data structures.

→ Languages like HTML & XML are used to store hypertext.

HTML - Hyper Text Markup Language
XML - Extensible Markup Language

1. HTML

- HTML is a standard markup language.
- It is used for creating web pages.

Key features

- Tags (<P>,)
- Attributes (href, src)
- Hyperlinks

Ex:

```
<html>  
<head>  
  <title> Example Page </title>  
</head>  
  <body>  
    <p> Welcome all </p>  
  </body>  
</html>
```

2. XML

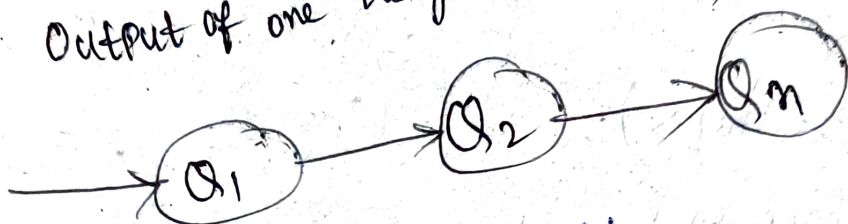
- XML stands for "Extensible Markup Language"
- It is a flexible markup language.
- Used to store & exchange data.

Ex:

```
<Person>  
  <name> abc </name>  
  <address>  
    <street> ... </street>  
    <city> ... </city>  
    <state> ... </state>  
  </address>  
</Person>
```


* Hidden Markov models (HMM)

- It is used for searching as Textual queries
- The output of one query is supplied/given to ~~the~~ another query as input
- Output of one query = Input of another query



- It is a chain process.
- Q_1 takes input & produce output.
- Output of Q_1 is given as input Q_2 .
- Q_2 produce some output.
- It is repeated until end of the process.
- Development for HMM approach
 - begins with ^{by} applying Bayes Rule to Conditional Probability.

$$P(A/B) = \frac{P(B/A) P(A)}{P(B)}$$

* Indexing

- It is a process of organizing & structuring data.
- It is the oldest technique for finding the items.
- Originally, indexing is called as "Cataloging".
- The evolution of IRS have changed the objectives of indexing.
- It is an important process in IRS.
- Indexing process can be manual or automatic.

