

IRS stands for Information Retrieval Systems.

Here

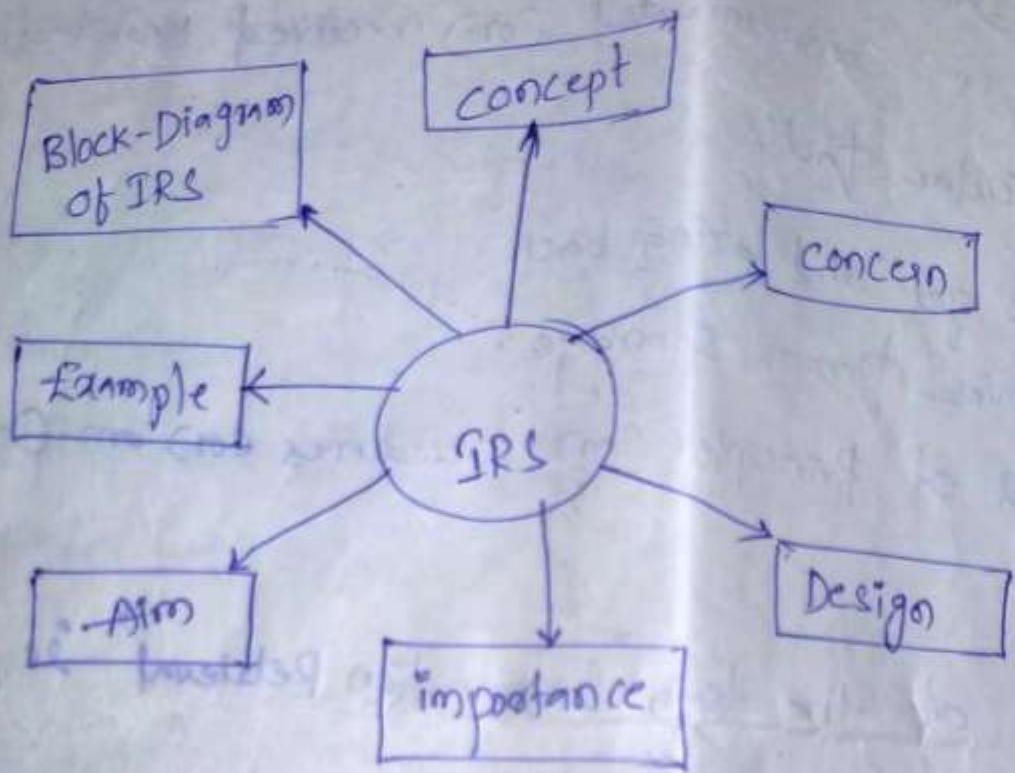
Information: Data & facts.

Information is communicated (or) received knowledge conveying a particular fact.

Retrieval:- To get and bring back i.e recover from storage.

System: A set of principles (or) procedures (or) an organized schema.

- ⇒ The meaning of the term "Information Retrieval" is "Very broad".
- ⇒ The term IR (Information Retrieval) was coined by "KELVIN MOORE" in 1950. It gained popularity in the research communication from 1961 onwards.
- ⇒ Goal of IRS is to provide the information needed to satisfy the user's question.



Introduction to Information Retrieval Systems:

→ The concept of Information Retrieval System (IRS) is "self-explanatory" from the "terminological point of view" and refers to a system which retrieves information.

Here, "self-explanatory" means easily understood from the information already given and not needing further explanation.

"Terminological" means Technical words and expressions used in a particular subject.

⇒ IRS is concerned with two basic aspects:

- i) How to store information &
- ii) How to retrieve information.

Design :-

→ An IRS is designed to "analyze process" and store sources of information and retrieve whenever required. (or)

→ An IRS is a set of "rules and procedures" for performing some (or) all of the following operations:

- 4
- * Indexing (or) constructing the representation of documents.
 - * Search formulation / constructing the representation of information needs).
 - * Searching / Matching representations of documents against the representations of needs.
 - * Index language / Generation of rules of representation.

Importance of IRS :-

- Information Retrieval Systems are Very ~~important~~ important to make sense of the data.
- Imagine how hard it would be to find some information on the Internet without Google (or) other search engines.

Aim of IRS :-

- It provide the "best possible" information from a database.

Example of IRS :-

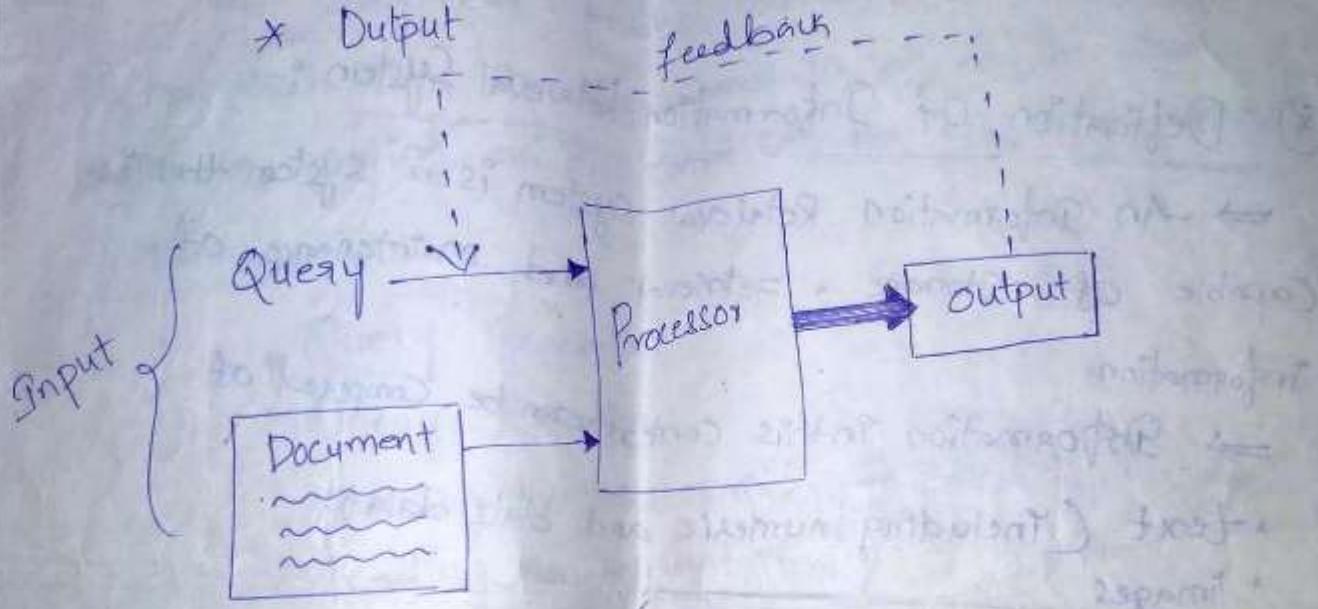
- The most obvious example of an IRS is "Google" and the English language has been extended with the term "Google it" means search for something.

Block Diagram of IRS:-

5

It consists of 3 components. They are

- * Input
- * Processor
- * Output



(fig: IR Block Diagram)

* Input :- It stores the only representation of the "Query" or "document".

Query : It represent the expression of the User information need.

Document : It represent the list of extracted words

consider to be significant.

* Processor :- It involve in performing "actual retrieval".
function and also executing the search strategy in response to a query.

* Output :-

A set of extracted document is called output.

feedback :-

It improving the subsequent run after sample retrieval.

④ Definition of Information Retrieval System :-

→ An Information Retrieval System is a system that is capable of storage, retrieval and maintenance of information.

→ Information in this context can be composed of

→ text (including numeric and date data).

→ images

→ audio

→ video and

→ other multi-media objects.

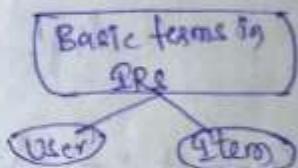
→ The following terms discussed in IRS.

User :-

The term "User" represents an end user of the information system who has minimum knowledge of computers and technical fields in general.

Item :-

The term "Item" is used to represent the "smallest complete unit" that is processed and manipulated by the system.



Objectives of Information Retrieval Systems

7

The general Objective of an IRS is minimizing the overhead of a user by locating needed information.

* Overhead:-

The Over-head can be expressed in terms of "time" i.e how much time the user spends in all of the steps leading to reading an item containing the needed information.

Example:-

- * Query Generation
- * Query Execution
- * Scanning the result of Query to select Item to read,
- * reading non-relevant items.

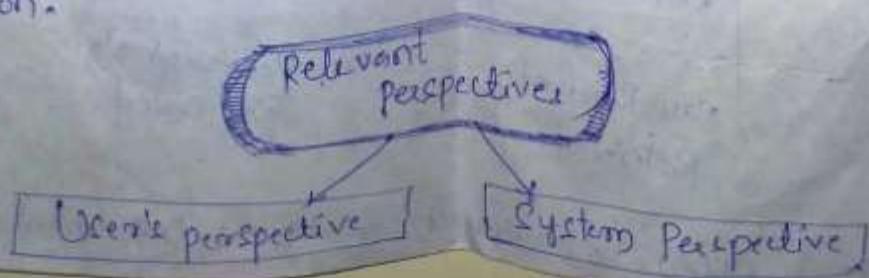
The success of the information system is very subjective, and it depends on two things

i) What information is needed.

ii) The willingness of a user to accept Overhead.

* Relevant Item:-

In Information retrieval Systems the term "relevant" item is used to represent an item containing the needed information.



* User's Perspective :-

From a User's perspective "relevant" and "needed" are synonymous.

* System Perspective :-

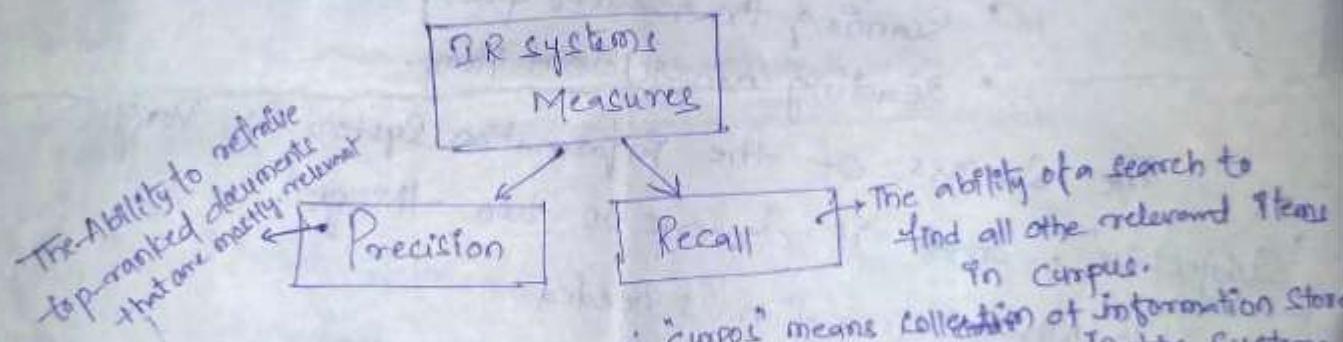
From a System Perspective, Information could be relevant to search statement even though it is not needed / relevant to user.

Eg: User already knew the information.

→ The Objective of IRS is commonly associated with two major measures:

i) Precision

ii) Recall.



Precision

→ When a User decides to search over a topic, the total database is logically divided into 4 segments, as shown below.

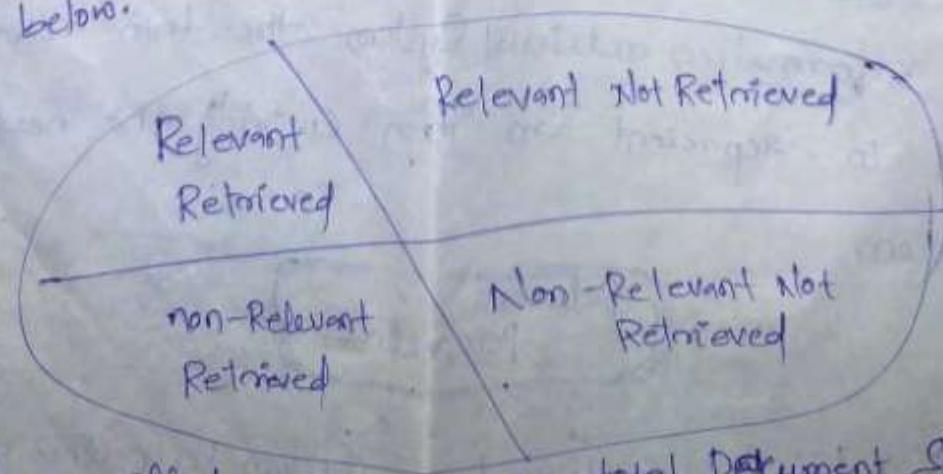


fig: Effects of search on total Document Space.

Let us see, what is "Relevant item" and "non-relevant item".

(Q)

Relevant Item:-

It is used to represent an item containing the needed information. That helps the searcher/user to answering the question.

Non-relevant Item:-

Non-relevant items are those items that do not provide any useful information directly.

⇒ The Precision and Recall formally defined as:

$$* \text{Precision} = \frac{\text{Number - Retrieved-Relevant}}{\text{Number - Total - Retrieved}}$$

$$* \text{Recall} = \frac{\text{Number - Retrieved-Relevant}}{\text{Number - Possible-Relevant}}$$

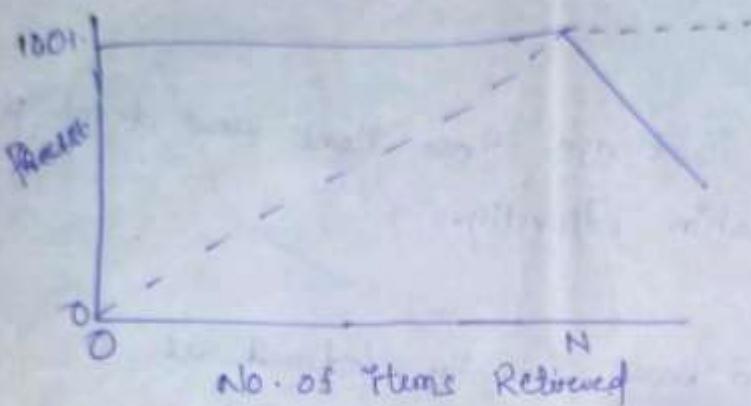
Here,

* "Number - Retrieved - Relevant" is the no. of items retrieved that are relevant to the "User's" needed search.

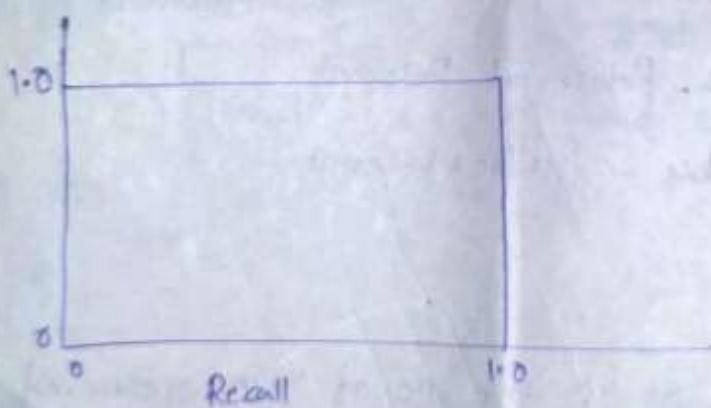
* "Number - Total - Retrieved" is the total no. of items retrieved from the "Query".

* "Number - Possible - Relevant" is the no. of relevant items in the "db" database.

→ In a search of "100%," (or) in 100% search "85%" is precision (i.e. Precision: 85%) and the remaining ~~recall~~ ^{is} "15%". If the user effort is overhead on avoiding the non-relevant items.
 Let us see the "Precision" and "Recall" graph:



(a) Ideal Precision and Recall



(b) Ideal precision/Recall Graph.

→ From the ~~ideal~~ diagram, the basic properties of precision (Solid line) and recall (dashed line) can be observed.

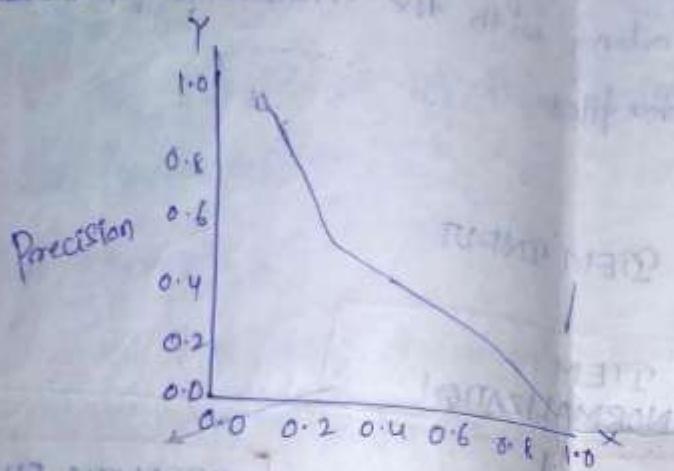
The, "Precision Starts off" at 100% and maintains that value as long as relevant items are retrieved.

"Recall Starts off" close to zero and increase as long as relevant items are retrieved. (until all possible relevant items have been retrieved)

→ Once all "N" relevant items have been retrieved, the only items being retrieved are "non-relevant".

→ Precision is directly affected by ~~mean~~ retrieval of non-relevant items and drops to the number close to zero.

⇒ Recall is not affected by retrieval of non-relevant items and remains at 100%.



Fig(c). Achievable Precision/ Recall graph.

→ fig(b) & fig(c) show the Optimal and currently achievable relationships b/w Precision and Recall.

⇒ To understand the implications on fig(c),

⇒ Assume that there are 100 relevant items in the database and from the graph at Precision of 0.3 (i.e 30%) there is an associated recall of 0.5 (i.e 50%).

This means there would be 50 relevant items in the hitfile ⇒ the recall value.

→ Precision of 30% means user would likely review 167 items to find the 50 relevant items.

* Functional Overview

A total Information Store and Retrieval System is composed of 7 major functional processes:

- ① Item Normalization
- ② Selective Dissemination of Information (i.e. "mail")
- ③ Archival Document Database Search
- ④ Index database Search along with the -Automatic File Build Process that support index files.

Item Normalization:-

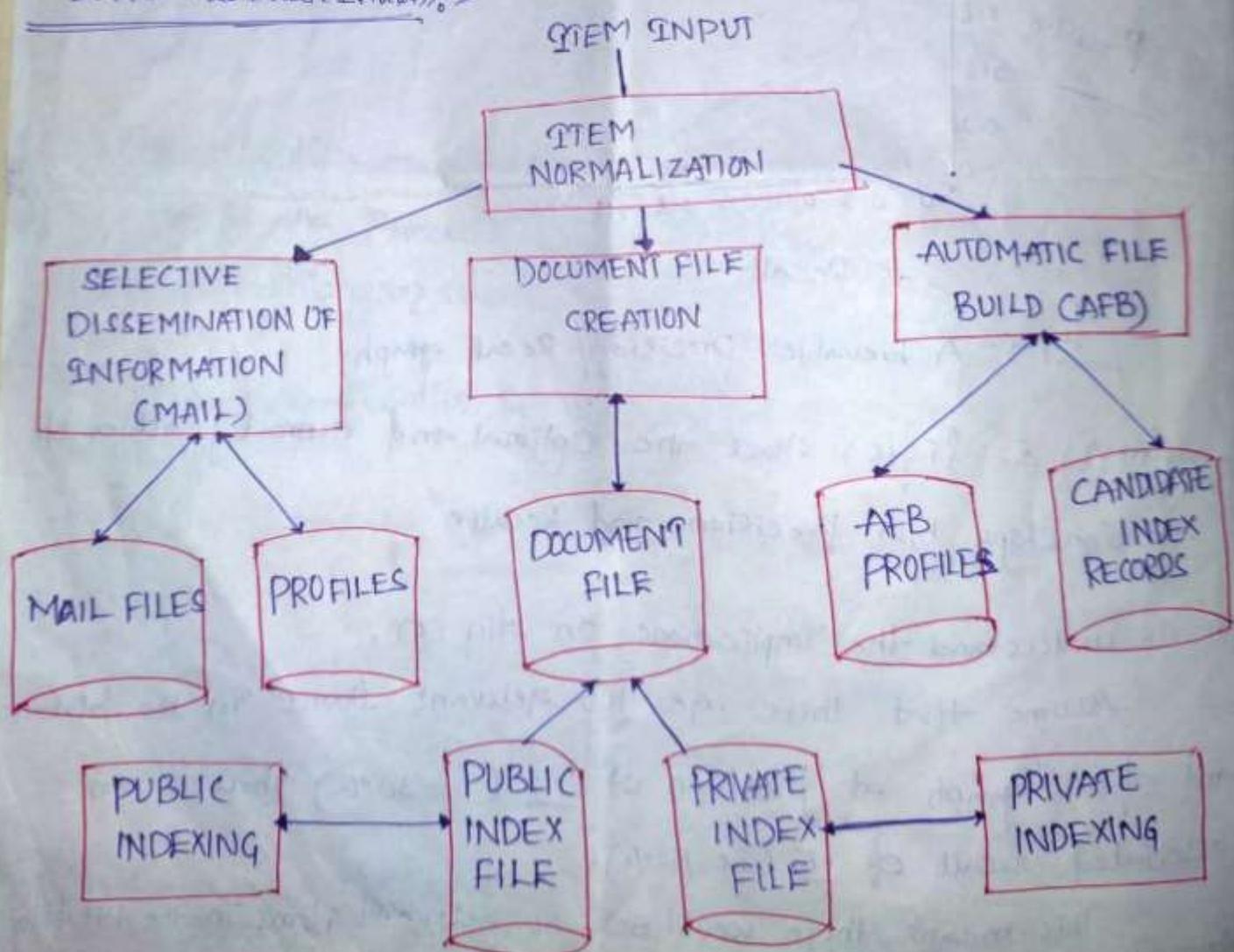


fig: Total Information Retrieval System.

① ITEM NORMALIZATION :-

The First Step in any integrated system is to normalize the incoming items to a standard format.

Item / Item Input :-

The term "Item" is used to represent the smallest "Complete Unit" that is processed and manipulated by the System.

The Definition of item varies how a specific source creates information.

For Example :-

A Complete document, such as a book, a newspaper or magazine could be an item. Sometimes the "chapters" and "articles" inside the magazine also consider as an item.

Normalization :-

It is the process of bringing (or) retrieving something to a Normal Condition (or) state.

⇒ Item Normalization provides "logical reorganization of the item".

The following "Operations" are performed during the item

- Normalization:

- * "Identification" of Processing tokens (Eg: words).

- * "Characterization" of the tokens.

- * "Stemming" of the tokens. (Eg: removing word endings).

Now, Let us See the Text Normalization Process: MARCH 14

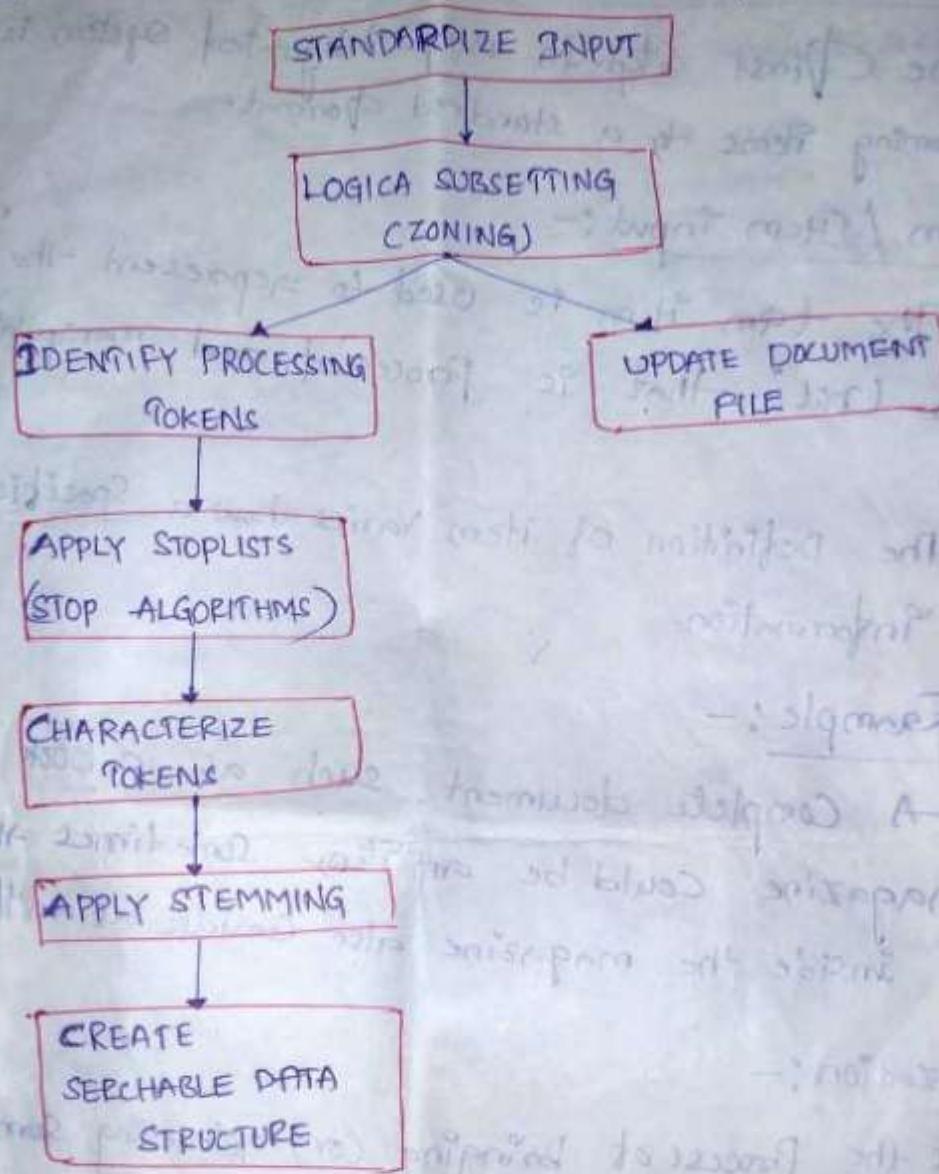


fig: The Text Normalization Process.

Standardize Input:

- Standardizing the input takes the different formats of Input data and perform the translation to and the result of translation is acceptable by the system.
- A system may have a single format for all the items (or) allow multiple formats.

- Potential issues -
Standardization
- for Example :-
- Standardization could be translation of foreign languages into "unicode". Here Unicode is an international standard based upon 16-bit (i.e. 2 bytes) that will able to represent all languages.
 - Every language has a different internal binary encoding for the characters in the language. Ex: UTF-8
 - "One Standard encoding" that covers English, French, Spanish etc is ISO-Latin.
 - The another encoding for other language groups such as Russian, Japanese, Arabic etc is KOI-7 & KOI-8.
 - Multimedia is an extra dimension do the Normalization
 - In order to normalize the textual input, the multimedia input also need to be standardized.

There are lot of options to standardize Normalization.

i.e. Input may be Video (or) Audio (or) Image.

Multimedia Options do the Standardization

Video

Audio

Image

Video :-

→ If the input is "Video" then following "digital standards" are being applied.

→ The digital standards are either MPEG-2, MPEG-1, AVI (or) Real Media.

MPEG → motion Picture Expert Group working group
SMPTE International standards

MPEG Standards are the most universal standards for 16

higher Quality Video where Real media is the most common

standard for lower quality video being used on Internet.

Audio:-

Audio Standards are typically WAV (or) Real media (Real Audio).

* WAV → Waveform Audio file format.

* Real audio is a proprietary audio format developed by Real Networks in 1995.

Image:- Images vary from JPEG to BMP, Bitmap image file

Logical Subsetting (Zoning):-

→ The second step in the text normalization process is to parse the item into logical sub-divisions called Zoning.

→ A typical item is sub-divided into zones and these zones which may overlap and can be hierarchical, such as Title, Author, Abstract, main Text, Conclusion and References.

→ In zoning search is restricted to "specific zones" sometimes.

for ex:-

If the user is interested in discussing Einstein's article.

then the search should not include the Bibliography, instead of that it include references to articles written by Einstein.

→ Once search is completed user is efficiently reviewing the needed information.

→ Zoning is performing the two tasks.

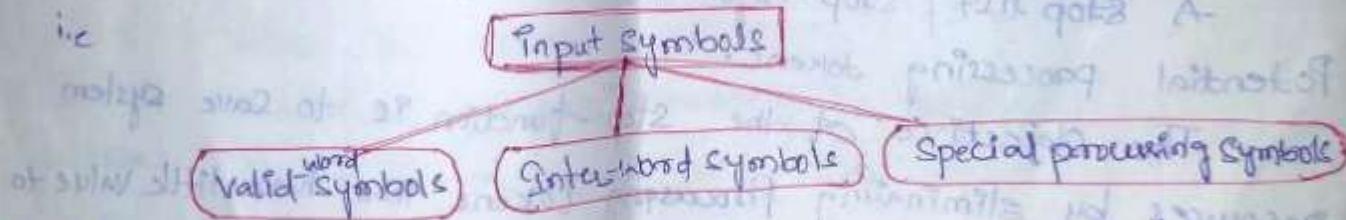
i) Identify processing tokens.

ii) Update Document file.

Identify Processing tokens:-

- Once standardization and zoning has been completed, let us see two terms "Identify" & "Processing token"
- Information (i.e words) that are used in the search process need to be "identified" in the item.
 - The term "processing token" is used because a "word" is not the most efficient unit in the search structure.
 - The first step in the identification of a processing token is determining the "word".
 - System - determine the words by dividing the input symbols into three classes:
 - Valid Word Symbols
 - Inter-Word Symbols
 - Special processing symbols.

i.e



Valid word symbols:-

A word is defined as a contiguous set of word symbols bounded by Interword symbols.

Inter-Word Symbols:-

In many systems, Interword symbols are non-searchable and should be carefully selected.

Examples: Semicolon, blank, an apostrophe

iii) Special processing Symbols:-

There are some symbols, that may required in special processing.

Example: Hyphen (-)

i.e. → Hyphen can be used in many ways,

* If it is used at left to the writer means it gives the Judgment of the Writer Ex: Einstein - eq.

* If it is used at the end of a line, means indicates the continuation of a word.

* If used in other place means it links the independent words.

-Apply stop lists (Stop Algorithms):-

Next Step,

A stop list / Stop Algorithm is applied to the list of Potential processing tokens.

The objective of the stop function is to save system resources by eliminating processing tokens that have little value to the system.

The best technique to eliminate the majority of these words via stop algorithm is first list out them individually.

Example of stop algorithms are:

→ Stop all the numbers greater than "999999" (this was selected to allow dates to be searchable).

* Stop any processing token that has number and character is terminated.

Characterize Tokens :-

19

Next step in finalizing the processing tokens is identification of any specific word characteristics?

The characteristic is used in the system to assist "disambiguation of a particular word".

For example:-

for a word "Plane"

- * the system understand that it could mean "level configuration" as an adjective.
- * "aircraft (or) facet" as a noun.
- * "The act of ~~smoothing~~ or evening" as a verb.

Apply Stemming :-

Once the potential processing token has been identified and characterized, most of the systems apply stemming algorithms to normalize the token.

The decision to perform stemming is tradeoff (or like b/w) between precision of search and recall.

For Example :-

The system must keep singular, plural, past tense, possessive etc.

Create Searchable data structure :-

Once the processing tokens have been finalized, based upon the stemming algorithm, they are used as Update to the searchable data structure.

→ The searchable data structure is the internal representation²⁰
(i.e. not visible to the user) of items that the user query
searches.

② SELECTIVE DISSEMINATION OF INFORMATION:

SDI is a concept that was introduced in information science by Hans Peter Luhn in 1958.

SDI is also called as "mail".

The selective dissemination of information (SDI) process provides the capability to dynamically compare newly received items in the information system against standing statements of user.

The Mail process is composed of 3 things:

i) Search process

ii) User statements of interest (profiles)

iii) User Mail files.

i) Search Process:-

→ After receiving item, it is processed against every user profile.

ii) Profiles :-

→ A profile contains search statements along with list of user mail files.

→ Profiles define all the areas in which a user is interested.

iii) User Mail files :-

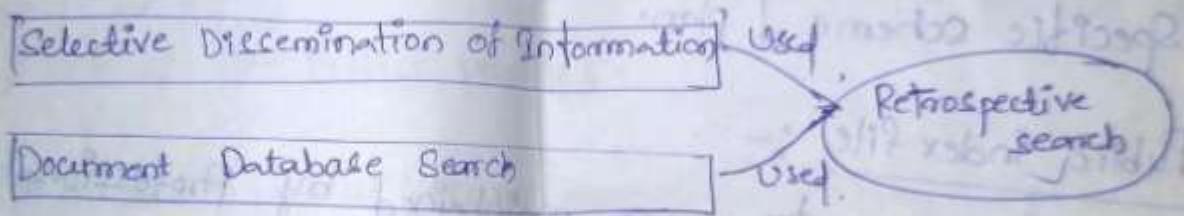
→ When a search statement is satisfied, the item is placed in the "mail files" associated with profile.

→ Items in mail files are typically viewed in the time of receipt order and they are automatically deleted after a specified time period (e.g. after a month etc).

(3) DOCUMENT DATA BASE SEARCH :-

21

- The document database search process is composed of
 - search process
 - User entered Queries
 - the document database.
- The document database, which contains all the items been received, procured & stored by the system.
- The Selective Dissemination of Information and document database search both are used retrospective search. i.e.



- The search can be performed in 2 ways i.e.
 1. Perspective Search
 2. Retrospective Search.

Perspective Search :

Searching in future and it is time valid.

Retrospective Search:-

Searching in past and it is time invalid.

- Document database is very large, it consist of 100's of millions of items.

classified in 2nd

Once the Document Database Search is completed then it creates the document file.

The Document file is classified into 2 types:

i) Public Index file (public indexing)

ii) Private Index file (Private indexing).

Indexing:-

It is originally called as "cataloguing".

→ Indexing refers to Organization of data according to specific scheme / plan.

Public Index file :-

→ Public index files are maintained by professional library services.

→ These files have access right that allow any one to search & retrieve data.

→ Public index files are used to perform the public indexing.

Private Index file :-

→ Private index files typically have very limited access.

→ Every user can have one (or) more no. of private index files that leading to a very large no. of files.

→ Private index files are used to perform the private indexing.

④ Automatic file Build (AFB):-

The automatic file build is also called as "Information Extraction".

AFB is classified in 2 types / performed through

- i) AFB profiles
- ii) Candidate Index Record

AFB profiles:-

The index term extraction process are stored in Automatic file Build profiles.

Candidate Index Record:-

- When an item is processed it results in creation of Candidate Index Records.

④

Relationship to DataBase Management System:-

These are 2 Major categories of Systems available to process items:

1. IRS
2. DBMS

i.e.

Relationship to

DBMS

→ A "confusion" can be arise, when the "Software System" supporting each of these Applications get confused with the data they are manipulating.

① * An Information Retrieval System is "Software" that has the features and functions required to manipulate "Information" items

VS (Versus)

* → DBMS that is Optimized to handle "structured" data.

② * Information is "fuzzy text".

The term fuzzy is used to imply the results from the minimal standards. (or) Controls on the creation of the text items i.e. the minimal standards in terms of Author is.

* The author is trying to present concepts, ideas and abstractions along with Supporting facts.

VS (Versus)

* The structured data is well defined data (facts) typically represented by "tables".

i.e. table is associated with "attributes" and attribute contains semantic description.

for example:

(25)

→ There is no confusion to what values enter in a specific database record.

⇒ we have different attributes like "employee Name" and "employee Salary".

③ * On the Other hand, if two different people generate an abstract for the same item, they can be different.

i.e. One abstract may generally discuss the most important topic in an item.

- Another abstract may specify the details of many topics.

(CRS) Versus

* With structured data a user enters a specific request and the result returned with the desired information.

The results are frequently tabulated and presented in a report format for ease of use.

④

* Search of "information" items in IRS & DBMS.

* A search of "information" items has a high probability of not finding all the items a user is looking for.

so that ^{in search} User enters the additional items of interest, this process is called "iterative search".

→ The information retrieval system gives the capabilities to the user in finding the relevant items such as "relevance feedback".

→ The Result from the Information System search are presented in "relevance ranked Order".

* The Confusion comes, when DBMS software is used to store "Information".

This is easy to implement, but the system lacks the "ranking and relevance" feed back features that are critical to the information system.

→ It is also possible to have structured data used in an Information System, (such as TOPIC).

finally,

from a Practical Stand point, The Integration of DBMS's and Information Retrieval System is Very Important.

* Digital Libraries and Data Warehouses:

Two Other Systems frequently described in the context of Information retrieval are "Digital Libraries and Data Warehouses".

(Data Marts).

There is significant Overlap between these two Systems.

Digital Libraries:

They are also called as "Online Library" (or) "Internet Library" (or) Digital Repositories.

A digital library is a library in which collections are stored in "digital format".

An electronic (or) Digital Library is a type of Information retrieval System (IRS).

It provides the architecture to

- * Model

- * Map

- * Integrate &

- * transform scattered information in digital documents.

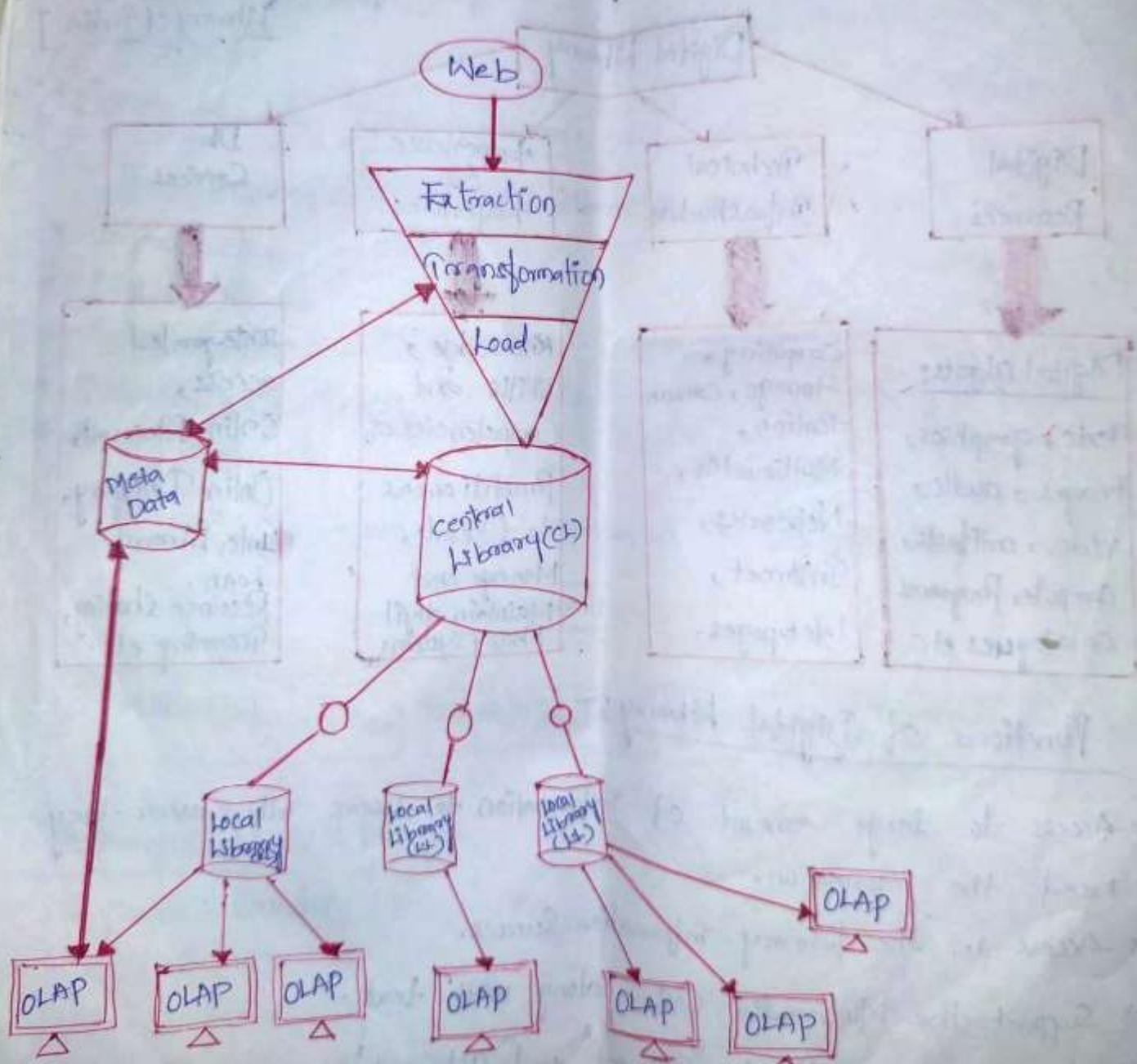


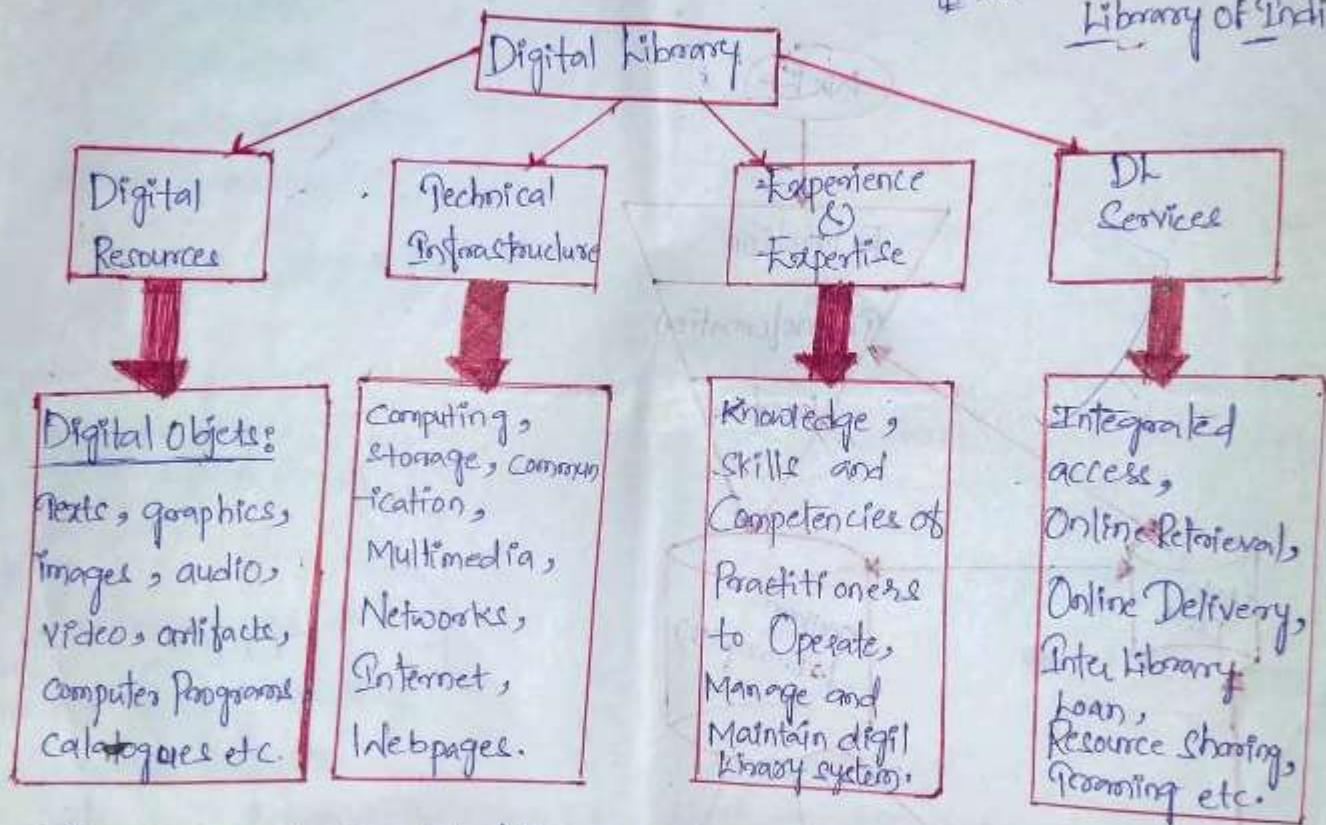
fig: Digital Library

L: OLAP → Online Analytical Processing

- → The term DL's (Digital Library's) was first popularized by the NASA Digital Libraries initiative in 1994.
- → Digital Library is not a single entity. It require technology that link the resources of many collections.

The links between the basic digital libraries and their resources are transparent to users shown in following figure.

[NDL → National Digital Library of India]



Functions of Digital Library :-

- * Access to large amount of information to users whenever they need the information.
- * Access to the Primary Information Sources.
- * Support the Multimedia Content along with text.
- * Network Accessibility in "Internet" and "Intranet".
- * Client - Server Architecture
- * Advanced Search & Retrieval.
- * Integration With Other Digital Libraries.

Purpose of Digital Library :-

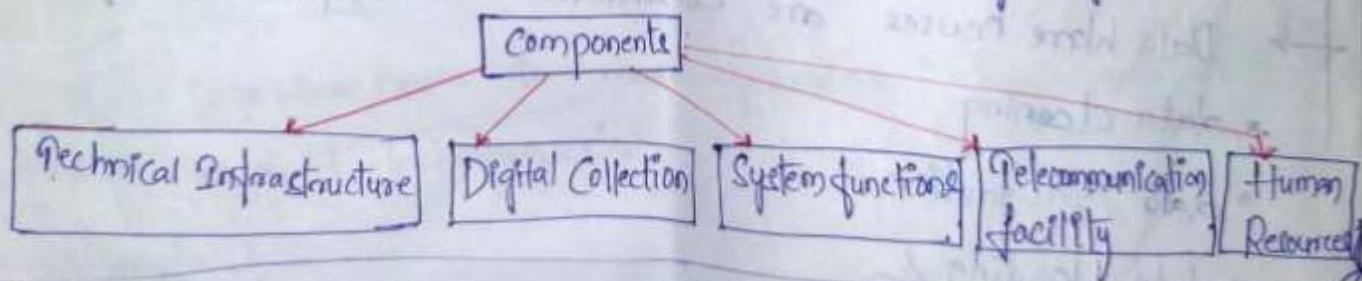
- It expedite the systematic development of procedures to collect, store and organize information in digital form.
- It promote efficient delivery of information economically to all users.
- It strengthens communication and collaboration "between" and "among" educational institutions.
- It takes "leadership role" in the generation and dissemination of knowledge.

Components of Digital Libraries :-

The components of Digital libraries are:

- * Infrastructure
- * Digital Collection
- * System functions
- * Telecommunication facility
- * Human Resources.

You can represent these components in following way.



Advantages of DL's :-

- No physical Boundary
- Multiple Accesses
- Universal Accessibility
- Round-the-clock Availability
- Added Values
- Enhanced IR.

Limitations:

- * Lack of screening (or) Validation.
- * Lack of Preservation of "Best in class".
- * Lack of Preservation of "fixed Copy".
i.e. for the record and duplicating scientific research.

Database house (or) Data Marts:-

→ A data warehouse is a repository of information collected from multiple sources, stored under a "Unified Schema" and that are usually resides at a "Single site", (or)

The process of collecting "Information across organization" is called.

→ A data warehouse refers to a place where data can be stored for useful mining.

where "mining / data mining" refers to process of extracting Useful data from the database.

→ Data warehouse are constructed via a process of

- * data cleaning
- * data Integration
- * data loading

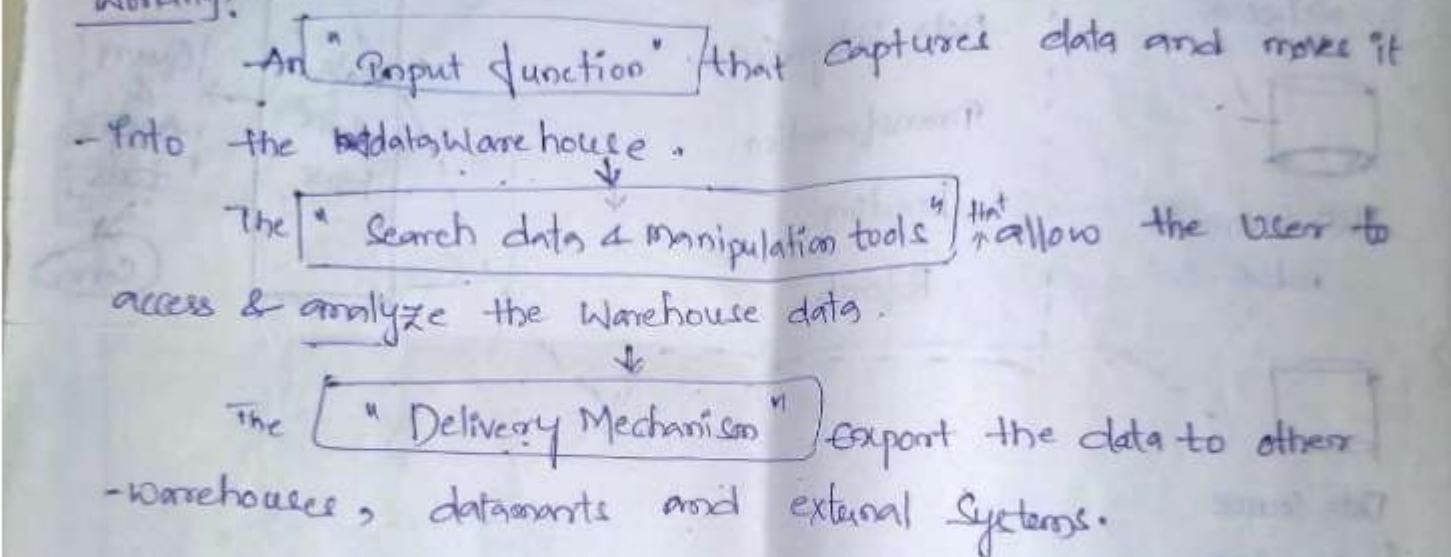
* periodic data refreshing.

→ The "term" "Data Warehouse" came from the Commercial Sector then Academic Sources.

Goal :-

- The goal of the data warehouse is to provide the critical information to decision makers so that they can answer the future queries.
- The data warehouse consists of the "data" carrying an "information directory" that describes the contents and meaning of data being stored.

Working:



⇒ Focus of Data Warehouse:

The Data Warehouse is more focused on "structured data" & "decision support technologies".

For Example:-

The typical functioning for construction and use of data warehouse for all electronic.

32



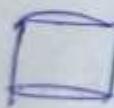
Data Source in Chicago



Data Source in Network



Data Source in Indi



Data Source
in Vancouver

Data cleaning
Data integration
Transformation
Loading
Research

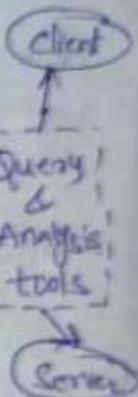
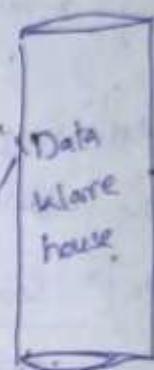
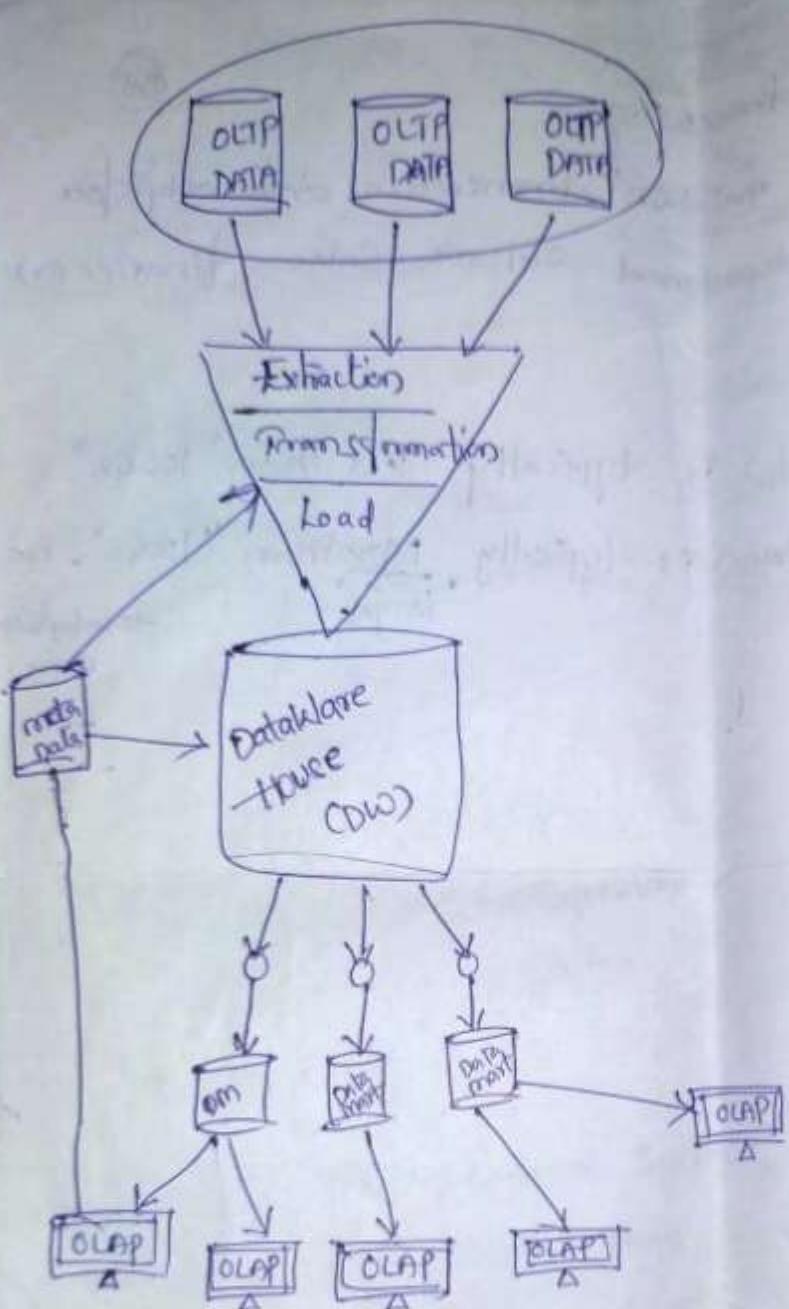


fig (a) Data-Warehouse.

Data mining :-

- The term Data mining is also called "KDD"
Knowledge Discovery in Database
- KDD is the process of ~~converting~~ discovering valuable information from a collection of data. (or)
- The process of converting the raw data into useful information.



DW → Data Warehouse
Dm → Data Mart

fig (b) Data warehouse.

OLAP :-

online analytical processing

→ This Approach is used to answer Multidimensional Data model.

* OLTP :-

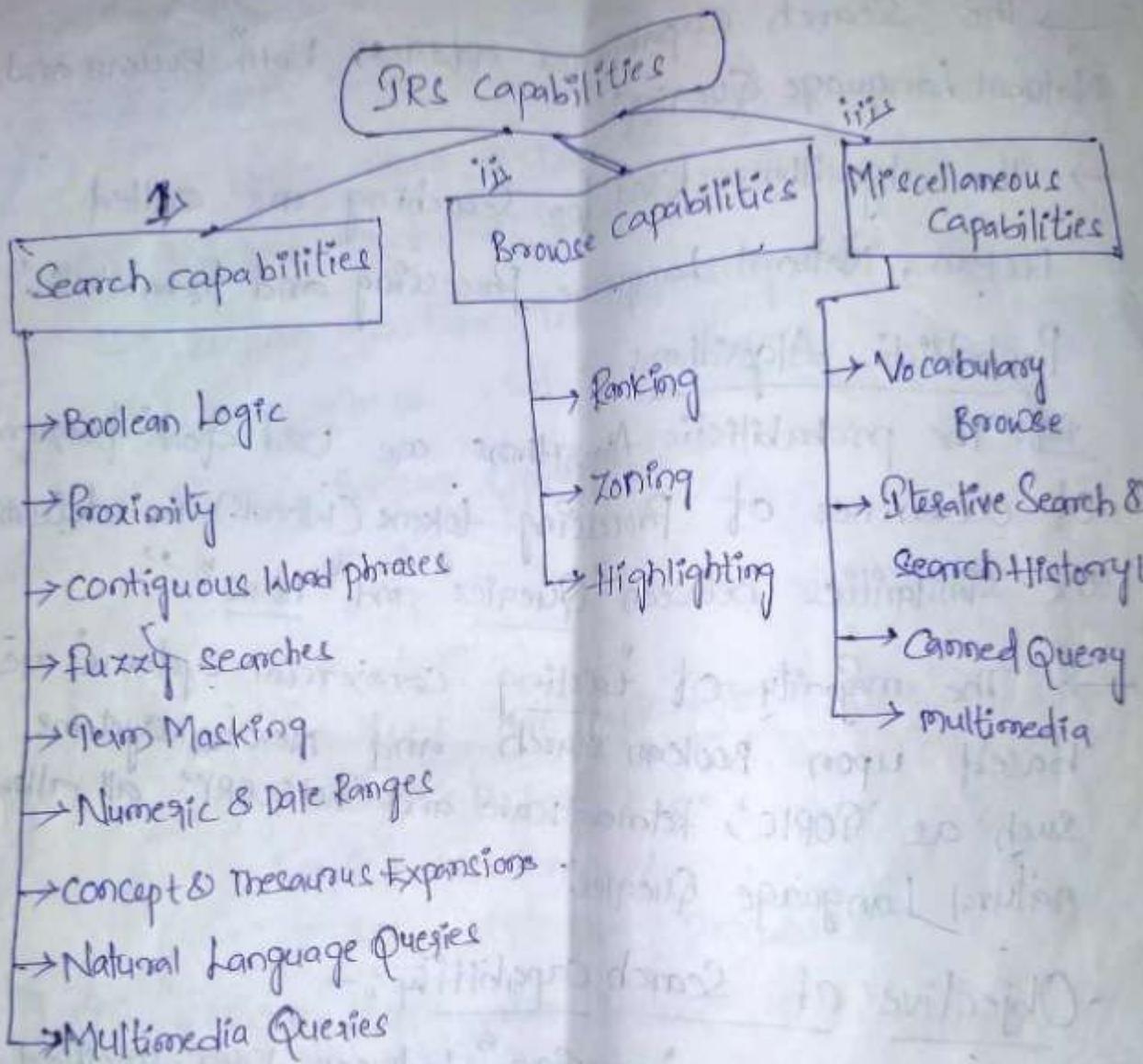
OnLine Transaction Processing

- It "captures and maintains" transaction data only for a specific department such as Sales, finance (or) HR (Human Resources).
- The DM (Datamart) is typically less than "100GB"
- The DW (DataWarehouse) is typically ~~less than~~^{More than} "100GB". i.e larger terabytes more.



IRS Capabilities

(35)



In Information Retrieval System, search and browse capabilities play crucial role to assist the user in locating the relevant items.

1) Search Capabilities :-

(3)

- The Search Capabilities address both Boolean and Natural Language Queries.
 - The algorithms used for searching are called Boolean, Natural Language Processing and Probabilistic.
- Probabilistic Algorithm:
- Def: The probabilistic Algorithms are used for "frequency of occurrence of processing tokens (words) in determining the similarities between "Query" and "Item" "
- The majority of "existing" commercial systems are based upon Boolean search and newer systems such as "TOPIC", "Retrievalware" and "INQUERY" all allow the natural language queries.

Objective of Search Capability:-

- It allows the "mapping" between user specified need and the items in the information database as a result information database will answer the need.
 - One of the concept that has been occasionally in commercial systems are "Weighting" concept
- Weighting:-
- The process of holding the location and ranking of relevant items.

IRS provides the different search capabilities
They are:

(37)

Boolean Logic :-

- Boolean logic allows a user to logically relate multiple concepts together to define what information is needed.
- The Boolean functions in processing tokens are identified anywhere within an item.
- The typical Boolean Operators are:
 - * AND
 - * OR
 - * NOT

These Operations are implemented using

Set intersection

→ Set Union &

→ Set Difference procedures.

→ The Boolean logic also include the "Order of Boolean operations" and "default precedence order for no precedence order".

Order of Boolean Operations :-

Placing the portion of search statement in parentheses are used to specify the Order of Boolean Operations.

default precedence order Operations:-

If the parentheses are not used the system follows a default precedence ordering of operations.

Ex:- typically NOT then AND then OR.

A Special type of Boolean Search i.e called "M of N" logic.

i.e A set of possible search "terms", that are identified and accepted by any item that contains subset of the terms.

* C for example :-

find any item containing any two of the following terms : "AA", "BB", "CC"

This can be expanded into a "Boolean Search" that performs an AND operation between all the combinations of two terms and OR operation to result together.

i.e "AA", "BB", "CC"
AND
AND

i.e $(AA \text{ AND } BB) \text{ AND } (AA \text{ AND } CC) \text{ AND } (BB \text{ AND } CC)$

AND Operation b/w
all combination

$(AA \text{ AND } BB) \text{ OR } (AA \text{ AND } CC) \text{ OR } (BB \text{ AND } CC)$

term₁

term₂

term₃

OR operation
for result.

Example:- Use of Boolean Operations:-

Search statement

* COMPUTER OR PROCESSOR NOT MAINFRAME

* COMPUTER OR (PROCESSOR NOT MAINFRAME)

* COMPUTER AND NOT PROCESSOR OR MAINFRAME

→ select all items discussing "computers" and/or "processors" that do not discuss mainframe.

→ select all items discussing computers and/or the items that discuss processors and do not discuss Mainframes.

→ select all items that discuss computer and not processors OR mainframe in the item.

fig: Use of Boolean Operations.

ii) Proximity :-

Proximity is used to restrict the distance allowed within an item between two search terms.

The typical format for proximity is:

TERM1 within "m" units of TERM2

(or)
TERM1 within m units of TERM2

Here "m" is indicating the "distance operator".

→ The distance Operator m is an integer number and units are in characters, Words, sentences and/or paragraphs.

Proximity Operators:

- * Direction Operator → (i.e before / after)
- * ADJ (Adjacent) Operator → (forward only direction i.e one direction)
- * Distance is set to zero → (within same paragraph/unit)

* Direction Operator:-

Some times, the proximity relationship contains "Direction Operator" indicating the direction (before / after) that the second term must be found within no.of (number of) units specified.

For Example :-

Search Statement

System Operation

- * "United" within five words of "American" → It would "hit on", "United states and American interest", "United Airlines and American Airlines". It would not "hit on", "United states of America and the American dream".

* ADJ (Adjacent) Operator :-

A Special Case of Proximity Operator i.e - the Adjacent ADJ operator, that normally has a distance operator of "one" directly, i.e forward only direction.

For Example :-

SEARCH STATEMENT

SYSTEM OPERATION

- * "Venetian" ADJ "Blind"

→ It would find the items that mention a "Venetian Blind" on a window but not items discussing a "Blind Venetian".

* Distance is set to zero :-

→ In another special case if "Distance is set to zero" means
"with in the same semantic unit".

for Example :-

SEARCH STATEMENT

"Nuclear" within zero paragraphs
of "Clean-up".

SYSTEM OPERATION

It would find the items that
have "nuclear" and "clean-up"
in the same paragraph.

iii) Contiguous Word phrase :-

→ In, A Contiguous Word phrase is two or more words that
are treated as a "single semantic Unit".

for Example :-

"United States of America"

i.e It is "four" words that specify a search term representing
single semantic unit (i.e "country").

→ The single semantic Unit satisfies all the proximity operators
i.e ADJ, Direction Operator & Distance ^{is set to} with Zero Operator.

→ A CKWP is specified in both ways

i.e * Query

* a special Search Operator.

for Example :-

A Query could specify "manufacturing" AND "United States
of America".

(HR)

i.e

which returns any item that contains the word
"manufacturing" and the contiguous words "United States of America"

→ The proximity and Boolean Operators are "Binary Operators"
but Contiguous Word phrases are "N"ary operators
Where "N" is the number of words in the CWP.

→ The contiguous word phrases are called "Literal String"
in "WAIS", and ~~Exact~~ phrases in "Retrieval Ware".

Where WAIS:- WAIS is a standard for

- * indexing
- * storing &

* retrieving the information from source
document that can be located anywhere on the Internet.

→ In WAIS multiple Adjacency (Adj) Operators are used to define a literal string.

Ex:- "United" ADJ "States" ADJ "Of" ADJ "America".

43

IV) Fuzzy search:-

fuzzy search provide the capability to "locate the spellings of words" that are similar to the entered search term.

→ This function is primarily used to compensate for errors in spellings of words.

→ A fuzzy search on the term "Computer" would automatically include the following three words from the information database:

"Computer"

"compiter"

"Computu"

"Computer"

"compute".

→ Maximum utilization of fuzzy search in the system that accepts (OCR) optical character read.

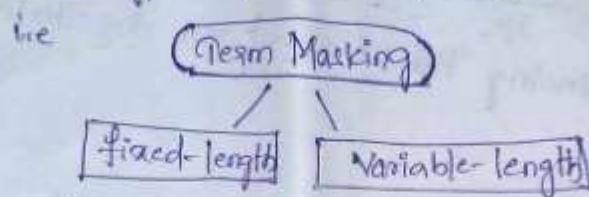
OCR :- In OCR a hardcopy item is scanned in to a binary image. (or)

Process of scanning the hardcopy into binary image is called 'OCR'.

Term Masking :-

(44)

- Term Masking is the ability to expand a query term by masking a portion of the term and that maps to the unmasks portion of the term.
- The value of term masking is much higher in systems that do not perform stemming.
- There are two types of search term masking:



- Some times - the fixed-length and Variable-length.
- Sometimes both "fixed and Variable length" are called "don't care" functions.

* fixed length Masking :- ({ })

- fixed length Masking is a "single position Mask".
- It Mask out any symbol in a "particular position".
- fixed length term Masking is not frequently used.
- It is typically not critical to the system.

for example ? Indicated by " { } ".

SEARCH STATEMENT

- * Multi-national

SYSTEM OPERATION

Matched "multi-national",
"multynational", "multinational"
But it does not match the
"multi national" since it is two
processing tokens.

- * Variable length Masking :- (*) -

→ Variable length Dont-care functions allows masking ~~action~~
of any number of characters within a processing token.

→ The Variable length Masking may be in the front,

at the end, at both ends.

i.e. * Subfix search :-
The front end means, if the masking is performed
at the "Starting position" of the processing token is called

* Subfix Search.

e.g. i.e. [* COMPUTER] → Subfix Search.

* Prefix Search :- If the masking is performed
at the "ending position" of the processing token is called
the "ending position".

Prefix Search.

i.e. [* COMPUTER*] → Prefix Search.

* Imbedded String Search :- If the masking is
performed at "both ends" of the processing token is called
"Imbedded String Search".

i.e. [* COMPUTER*] → Imbedded String Search

→ The Variable length Masking is indicated by "x".

for Example :-

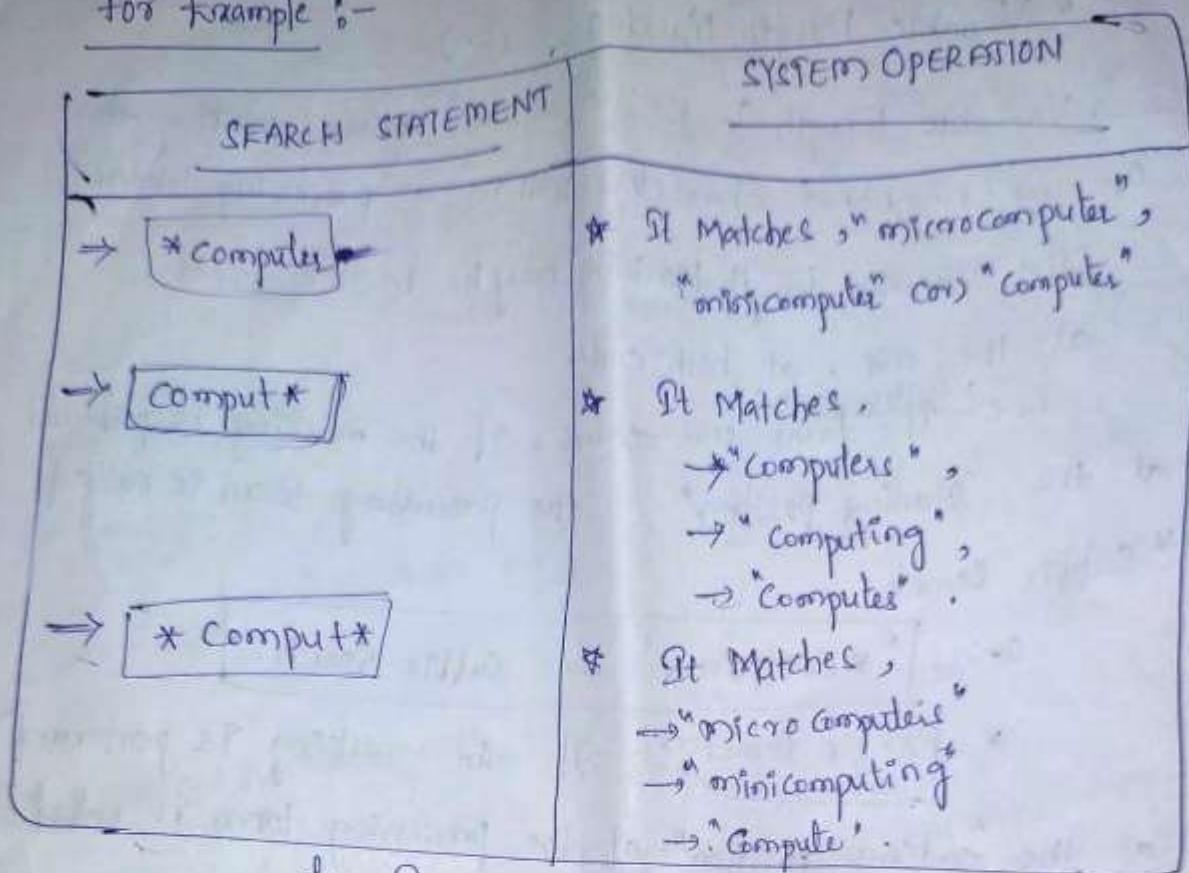


fig. Term masking

VI> Numeric and Date Range :-

- TermMasking is useful "when applied to words, but does not work for finding ranges of numbers (or) Numeric dates.
- To find the numbers larger than "125" using a term "125x" will not find any number.
- As a part of system Normalization process, characterizes words as "numbers (or) dates".

for Example :-

- * for Numbers and Dates :
- ↳ "125 - 425" means finding the numbers ^{range} in between two numbers.

i.e. [4/12/1993 — 5/10/1995] means finding date range
in between given dates. (H)

* for infinite ranges :-

for representing the infinite ranges between the numeric data, it uses special type of operators

i.e. $[> \text{con})$ $\geq (\text{con})$ $< (\text{con})$ $\leq (\text{con})]$

Example : " > 125 " means Larger Number above the 125.

" ≤ 233 " means Smaller (or) equal to 233.

Vii) Concept & thesaurus Expansion :-

* Concept :-

→ A concept ^{class} is a "tree structure" that expands each meaning of a word into potential concept.

→ Concept classes are sometimes implemented as a

"Network structure" that links the word stems.

→ "Concept class representation":

The Concept class representation, assist a user who has minimal knowledge of "concept domain". So that user can easily expand related concepts.

for Example :-

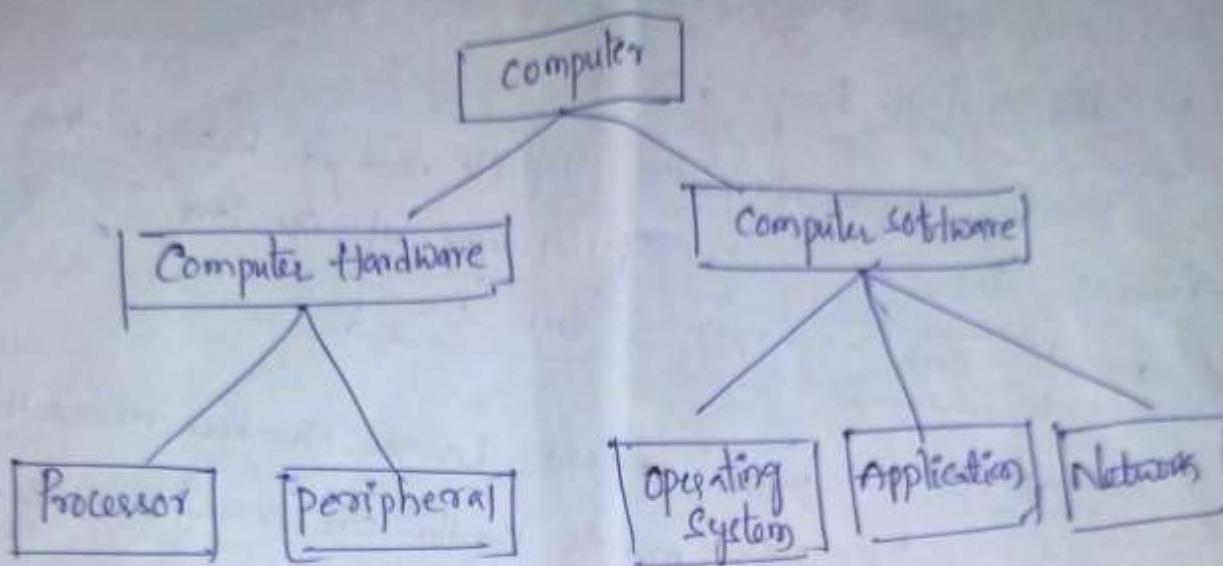


fig:- Hierarchical Concept class structure for "Computer".

Thesaurus Expansion :-

A Thesaurus is typically a one-level (or) two-level expansion of a term to other terms that are similar meaning.

For Example :

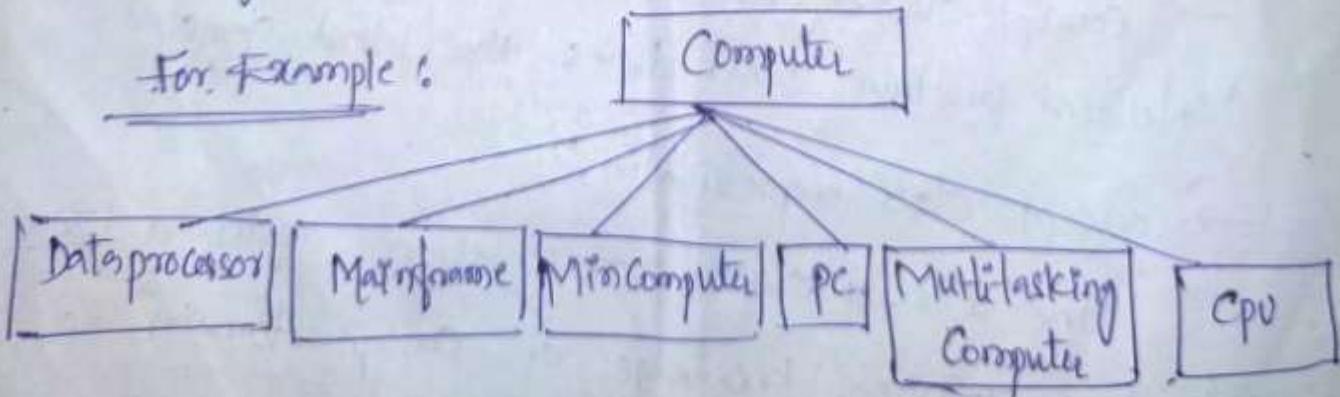


fig: Thesaurus for term "Computer".

viii) Natural Language Queries:-

→ Natural language query consists only of "normal terms" in the user's language, without any special syntax or format.

Natural Language Queries allow a user to enter a prose statement that describes the information, i.e. the user want to find.

The longer prose returned the more accurate results.

The most difficult logic case associated with Natural Language Queries is the ability to specify "negation" in the search statement.

Example :- i.e Example of Natural Language Query:

find for me all the items that discuss "Oil reserves" and current attempts to find new oil reserves. Include any items that discuss the international financial aspects of the oil production process. Do not include items about the oil industry in the United States.

→ The minimized input for above example is

Oil reserves and attempts to find new oil reserves,
international financial aspects of oil production,
not united states oil industry.

⑫

The Boolean query for same search statement is:-

$$("locate" \text{ AND } "new" \text{ AND } "oil reserves") \text{ OR } ("international" \text{ AND } \text{ finance} \text{ AND } "Oil production") \text{ NOT } ("oil industry" \text{ AND } "United States")$$

ix) Multimedia Queries :-

The User Interface becomes far more complex, with the introduction of the availability of multimedia queries.

- The image ^{could be} used to search the images that are part of an item.
- They also could be used to locate specific scene in a video product.

2. Browse Capabilities :-

Browse capabilities provide the user, capability to determine which items are of interest and select those items to be displayed.

There are 2 ways of "displaying" summary of items that are associated with Query.

i.e i) Item line status

ii) Data Visualization.

Item line status :-

from the summary,

The user can select "the specific items" and "Zones within the items" for display.

Data Visualization :-

The system allows easy transitioning between the "Summary display" and "Review of Specific Items".

Browse capabilities are of 3 types.

i) Ranking

ii) Zoning

iii) Highlighting.

⇒ Ranking :-

(F2)

- The introduction of Ranking based upon
 - * Predicted relevance values.
 - * the status summary displayed.
 - * the relevance score associated with items along with
 - brief description.

Relevance score:

- The Relevance Score estimate the item that satisfies the search statement.

- The Relevance Scores are Normalized to a value between 0.0 and 1.0.

→ Here

"The Highest Value of 1.0" means the system is sure that the item is relevant to search statement.

The Lowest Value of 0.0 means "not relevant".

- In addition to ranking based up-on the "characteristics of item" and "the database" → is a technique that contains the item that user might like.
- ⇒ Ranking uses "collaborative filtering" for selecting and Ordering output.

In this case, when user reviewing the items it provide feedback to the system on the relative value.

The "collaborative filtering" has been very successful in sites such as AMAZON.COM, MovieFinder.COM are deciding what products to display to users based on their queries.

ii) Zoning :-

→ The process of dividing the standard input into logical sub-divisions is called zoning / zones.

Zone : The 'title' is frequently called as zone.

It provide the additional information to the user with "relevance weight" to avoid selecting the non-relevant items for review.

→ The Zoning is used for in minimizing what an end user needs to review from a hit item in Passage based Search and retrieval.

iii) Highlighting :-

→ Highlighting indicates how strongly 'the highlighted word' participated in the selection of the item. So that User quickly focus on potentially relevant parts of the text.

→ Highlighting has been useful in Boolean systems to indicate the cause of retrieval.

i.e it provide the direct mapping between the "terms in the search" and "terms in the item".

→ The highlighting may vary by introducing colors and intensities to indicate the importance of a particular word.

3. Miscellaneous Capabilities:-

The additional functions that facilitate the user's ability to input queries, and reduce the time taken to generate the queries.

The Miscellaneous Capabilities are 4 types:

i) Vocabulary Browse

ii) Iterative Search and Search History Log

iii) Canned Query

iv) Multimedia.

i) Vocabulary Browse :-

Vocabulary Browse Provides the capability to display in alphabetical sorted Order Words from the document database.

All "unique words / processing tokens" in the database are kept in Sorted Order along with a number count of unique lines in which the word is found.

The User can enter "a word" (or) "Word fragment" and the system will begin to display the dictionary around the entered text.

i.e. The System indicates what word fragment the User entered, and then it alphabetically displays other words found in the database.

The User can continue scrolling in either direction of reviewing additional terms in the database.

App: It helps the User to determine the Impact of using a fixed (or) Variable length mask on a search term, and potential mis-spellings.

The computer in effect is searching for (iii)
 "Computation" (or) "computative" (or) "compulsory".
 Some one entered the word "computer" they meant
 for potentially "computer".

For Example :-

Term: Comput

TERM	OCCURRENCE
Comput	265
* Computation	1245
* Compute	1
* computes	10,800
* Computer	18
* Computative	29
* computes	

fig: Vocabulary browser list with entered term "comput".

(ii) Iterative Search and search history log:-

"Iterative searching" and "Search history log" summarize previous search activities so that user access the old previous results from the current user session.

(16)

Iterative search :-

- The results of the previous search can be used to create a new query, rather than typing a complete new query. This has some effect as taking the original query, i.e. "result is same as original query" by adding the additional "AND" condition to the search statement.

This process is called "Iterative search".

Search History Log :-

- The search history log is the capability to display "all the previous searches" that were executed during the current session.
- It provides the facility to locate the previous searches as starting points for new searches.

iii) Canned Queries:-

- The capability to "Name" a query, and "Store" it to be retrieved, or "executed" during a later user session. It is called Canned Query / Canned Queries.
- The Canned Queries are also called as "Stored Queries".
- Queries that start with a Canned Query are significantly larger than ad hoc queries.

IV Multimedia :-

(58)

Multimedia introduces new challenges, for potential items, that satisfy the discovered query, and multimedia techniques are used for displaying them.

The "transcribed text" used as navigation technique through the audio.

i.e. they appropriately label this new paradigm
What You See Is (Almost) What You Hear
(WYSIWYH).

→ The transcribed text could be used as an index
onto future retrieval of the audio source.
i.e. The tonal information provided by the original
Speech.

Cataloguing and IndexingCataloguing :-

- Cataloguing is a process that creates "metadata" and represents "information sources" such as books, sound recordings, moving images etc.
- The Cataloguing provides the information such as
 - * the name of the creator,
 - * the title and the subject term that describes the source, through creation of a bibliography record.

Indexing :-

The indexing is one of the most important process in the IR system.

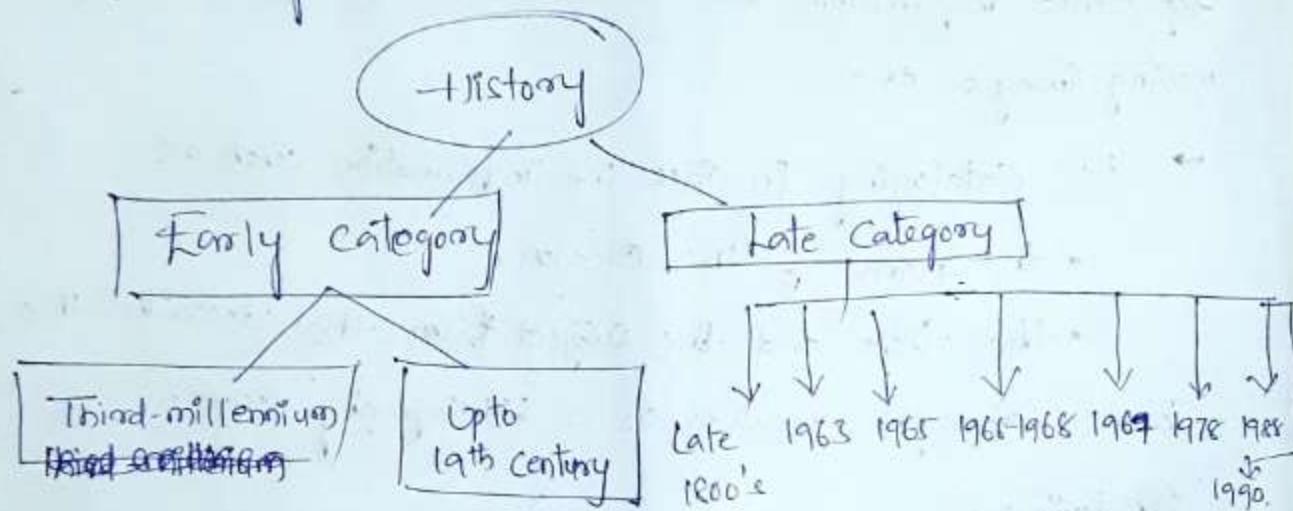
The transformation from the received item to Searchable data structure is called indexing.

Indexing helps to retrieve the information "accurately" and "efficiently."

- ⇒ Indexing Originally called as "cataloguing" and it is the oldest technique for identifying the contents of items.
- Objective: The Objective of the cataloguing is give the access points to a collection, that are expected and most useful to the users information.

* History and Objectives of Indexing:

The basic information required for an item is "what is the item" and "what it is about", has ~~been~~ not change over the centuries.



Early Category :-

As early as,

In third Millennium, in "Babylon"

The Libraries of "Cuneiform tablets" were arranged by the subject (By Hyman - & q).

Up to 19th century :-

Up to 19th century there was little advancement in "Cataloguing".

i.e. The advancement is only change in the methods used to represent the basic information.

Late Category :-

* In the late 1800 :-

The subject indexing became "Hierarchical"
Eg. Dewey Decimal System (DDS)

Dewey Decimal Classification (DDC).
(COS)

DDS/DDC :-

The Dewey Decimal System is the world's
most widely used way to Organize Library Collection.

i.e. The Dewey Decimal System used by the
librarians to arrange "books" via "Subject".

→ DDC allows new books to be added to a library

in their appropriate location based on subject.

It was first published by "Melvil Dewey" in
1876, in United States.

* In 1963 :-

The Library of Congress initiated.

i.e. A study on the computerization of bibliographic
surrogates.

Library of congress (lc) :-

The LC is "research library" that officially serves the "United states Congress".

The Library of Congress "Collections are Universal" not limited by subject (or) format (or) national boundary.

It include "research materials" from all parts of "the world" and in more than 420 languages.

The Library of Congress

→ "Established" in April 24, 1800 and
→ "Located" in Washington.

* In 1965 :- "catalog system - DIALOG"

The earliest commercial cataloging system

"DIALOG", which was developed by "Lockheed Corporation".
for NASA.

* from 1966 - 1968 :-

The Library of Congress ran the MARC I pilot project.
MARC → Machine Readable Cataloging

It standardizes the structure,

- * contents and

- * coding of bibliographic records.

* In 1969 :-

The system (Congress Library System) became operational in 1969.

* In 1978 :-

The DIALOG became commercial, with three government office of indexes to technical publications.

* In 1988 :-

The "DIALOG" was sold to "Knight-Ridder". At that time "DIALOG" contained "320" index databases used by over "91,000" subscribers in 86 countries.

* In 1990 :-

- The significant reduction occur in cost of processing power and memory in modern computer along with full text of item in electronic form, from publishing stage.
- Finally,
- Until recent, Indexing was accomplished by creating a bibliographic citation in a structured file that refers to the original text.

→ The indexing process is typically performed by "Professional Indexers" associated with library organization.

(63)

From the ancient presents, clay tablets, cones and brick fragments "inscribed" using an ancient writing system known as cuneiform

Cuneiform tablets → from the Library of Congress' collections.

- The word "cuneiform" is derived from "Latin".
cuneus for "wedge" and forma meaning "shape".
- Cuneiform is a "Logo-Symbolic" script that was used to write ~~several~~ several languages of the Ancient middle East.
- Cuneiform was developed by the Sumerians, who thrived during the 3rd-millennium B.C.
- Sumerians influenced ^{their} culture and development beyond their original borders in mesopotamia (Present-day Southern Iraq) & earliest civilization.
- Originally Cuneiform signs were "pictograms", later it also became symbolic.

The materials used in Cuneiform - "clay and reeds".

- The "clay and reeds" both are ^{readily} available.
- "Reeds" were used as "writing implements".

- The tip of a reed was impressed into a wet clay surface to the draw strokes of the sign thus it acquiring a "wedge-shape".

i.e. The Sumerians invented this writing system, which involves the use of wedge-shaped reed stylus to make impression in clay.

⇒ The Library of Congress acquired its collection of cuneiform materials in 1929 from Kirkor (an art dealer).

DIALOG :-

DIALOG online search system - (1966 became available).

- DIALOG was the first interactive, online search system addressing 'large database' and ~~and~~ ~~not~~ allowing 'iterative refinement of results'.
- DIALOG was developed at "Lockheed Palo Alto Research Laboratory", and it extended through contracts with NASA.

Why DIALOG became popular:

Its speed, ease of use and wide range of data content attracted the professional users worldwide including

- * scientists
- * attorneys
- * educators &
- * librarians.

- DIALOG preceded major "Internet search tools" by more than two decades.

By 1972,

The commercial introduction extended to all the professions by allowing "the online access" to large collections of digitized materials by a way of "command language" that allowed the researcher to interactively refine results.

DIALOG's ability to allow interactive refinement of results:-

* DIALOG's "major technical innovation is reflected in its name":
i.e It enables a conversion between the "searcher" and the "computer".

* "Search at its best is a conversation":
i.e "An ~~is~~ negative, interactive process where we find we learn".

Other Characteristics of DIALOG Language:

There are 5 important characteristics:

1. The search question is constructed at search time.
2. Dialog is designed for nonspecialists.
i.e It avoids communication barrier.
3. Command language is independent of the particular data it searches.

4. As "an Online-system".
 i.e. It allows continual specification of the search question based on the examination of intermediate results.
5. Control of process lies with the user.
 i.e. Computer merely serve as a data-processing extension of user.

DIALOG Commands! -

The DIALOG system provides the number of commands with which the "Searcher" interact with the "Computer".

A search consist of

- i) Identifying
- ii) Selecting terms and phrases that reflect the user's interest.
- iii) Combining descriptors into search queries expressions
- iv) Examining the retrieved citations and modifying search expressions.

Each of these functions accomplished with a particular command.

The 4 Principal Commands are:

- i) EXPAND
- ii) SELECT
- iii) COMBINE
- iv) DISPLAY

EXPAND :

EXPAND provides list of "synonyms", related terms and similar combination of words defining his/her interest.

SELECT :

SELECT is used to selecting ~~the~~ list of terms (i.e term a, term b, term c ...) ~~(or)~~ or

"a range of terms" i.e (term a - term d) ~~(or)~~
(or)

a list of ranges and terms.

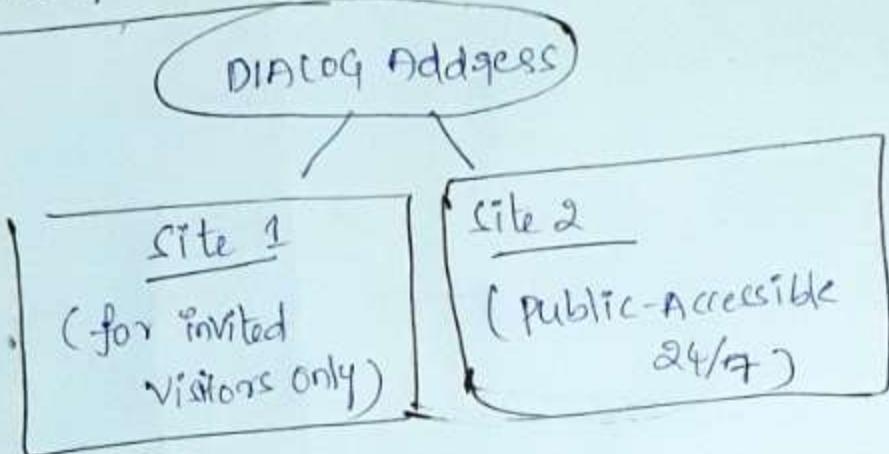
COMBINE :

The list of Boolean Expressions are used to combine the documents.

DISPLAY :

It allows the searcher to successively display the documents contained in the resultant set.

DIALOG Address :-



Site 1 :- (for invited visitors only)

* Lockheed Martin Advanced Technology Center,
formerly Lockheed Palo Alto Research Laboratory
Bldg. 201 (front lobby site)

Site 2 : (publicly-accessible 24/7): mountain view,
Computer History Museum, (on inside face of
front patio brick wall.)

Security / protection of Site 1 & Site 2 :-

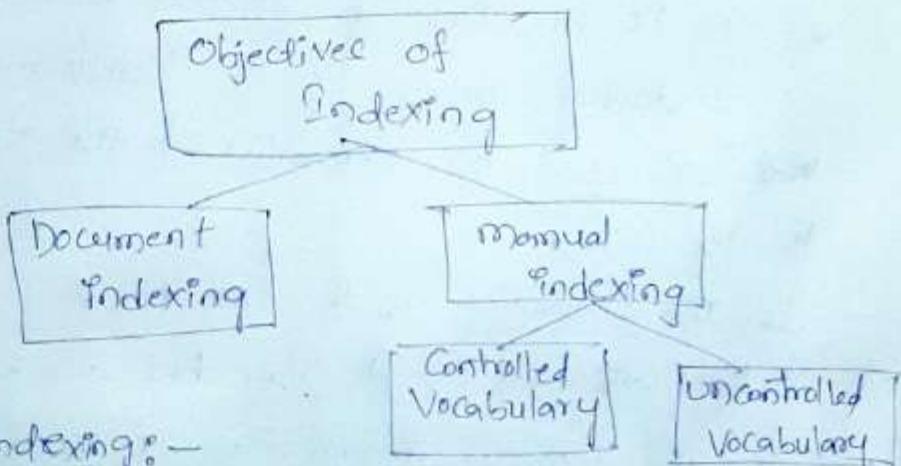
Site 1 : It requiring a specific user security clearance
to access site. (within a secure facility).

Site 2 : freely accessible to the public 24/7.

Objectives :-

The Objectives of "Indexing" change with evolution of Information Retrieval system.

The Indexing Define the "source" and "major concepts" of an item that provide the standard mechanism for "index terms".



Document Indexing:-

The Document file provides the new class of indexing called "Document indexing".

Manual Indexing:-

In Manual Indexing Environment,

* Controlled Vocabulary:-

The Use of Controlled Vocabulary makes the indexing process "slower" but Potentially Simplifies the searchprocess.

* Uncontrolled Vocabulary:-

The Use of Uncontrolled Vocabulary makes the indexing process "faster" but Search process much more difficult.

(70)

The availability of "Item in "electronic form" changes
the objectives of manual indexing.

i.e The source information can be automatically
extracted.

(*)

Indexing Process

→ It is an important process in IR.

→ Indexing process is a transformation of an item
that "extracts" the semantics of the topics discussed
in the item.

→ The Semantics of the item not only refers the
subject discussed in the item but also weighted system.

→ The extracted information i.e used to create
"the processing tokens" and "the searchable data structures".

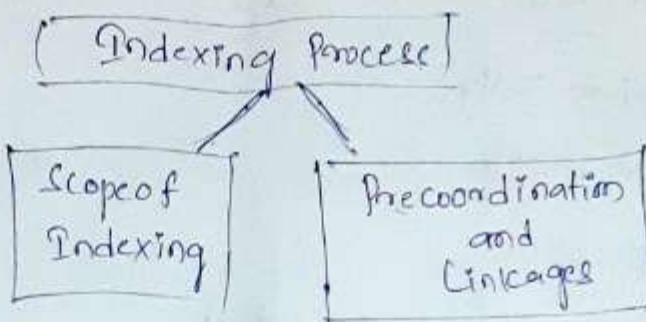
→ The index can be based on the

- full text of the item,

- automatic or manual generation of the item

- natural language representation of the item.

→ The result of this process stored in one of
the data structure called "Inverted data structure".



Scope of Indexing :-

When scope of indexing performed manually,
the following problems are occur with the interaction
of two sources:-

i) The Author

ii) The Indexer.

The Author :-

The vocabulary domain of the author may be different
than that of the Indexer.

The ^{to} Indexer :-

→ The Indexer is not an expert on all the areas
the item is being presented and the result is in
different quality levels of indexing.

→ The Indexer must determine when to stop
the indexing process.

factor effecting the scope of indexing:-

i) exhaustivity

ii) specificity

(f2)

Exhaustivity :-

The extent to which
extension

→ The different concepts in the item are indexed.

Index

Specificity :-

The preciseness or "The exact index terms used
in indexing" called specificity.

Linkages & precoordination :-

Linkages :

The linkages are used to correlate the
related attributes associated with concepts discussed
in the item.

Pre co-ordination :-

The process of creating the term "linkages" at
index creation time is called precoordination.

Post Co-ordination :-

When the index terms are not coordinated at
index time, the coordination occurs at search time
is called Post co-ordination.

The post coordination is implemented by "AND" _{into} terms.

Example :-

The different types of linkages as shown below.

(3)

"The drilling of oil wells in Mexico by CITGO and

The introduction of oil refineries in Peru by the US."

4

3-2

INDEX TERMS

→ Oil, wells, Mexico, CITGO, refineries, } No linking of terms.
Peru, BP, drilling.

→ * (Oil wells, Mexico, drilling, CITGO) } linked (pre-coordination)
* (US, oil refineries, Peru, introduction) }

→ * (CITGO, drill, Oilwells, Mexico) } linked (precoordination)
* (US, introduction, Oil refineries, Peru) } with position indicating
role.

→ * (SUBJECT: CITGO;
ACTION: drilling; drilling → drilling)
OBJECT: Oil wells;
MODIFIER: in Mexico) } linked (pre-coordination)
with modifier
→ radicalizing role: 8
* (SUBJECT: US; modifier
ACTION: introduces; introduction → introduces)
OBJECT: Oil refineries;
MODIFIER: in Peru) } with modifier

fig: Linkage of Index Terms =

* Automatic Indexing:-

Automatic indexing is "automatically" determine the "index terms" assigned to an item.

Automatic indexing is deals with "Simplest case" and "Complex case".

- * Simplest case:-
The simplest case deals with the total document of indexing.

- * Complex case :-
The complex case deals with the "Human Indexer".
The Human Indexer determines the "limited no. of index terms" for major concepts in the item.

Result of Automatic Indexing:-

The automatic indexing results fall into two classes:

- i) Weighted indexing system
- ii) Unweighted indexing system

Weighted Indexing:-

→ Weighted Indexing represents the "term values" placed inside the "index terms".

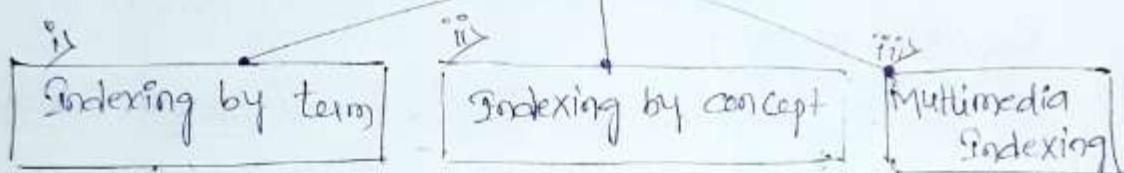
⇒ The Basic function of the weighted indexing associated with "frequency of occurrence of an item".

ii) Unweighted Indexing :-

- The unweighted indexing represents the existence of an index term in a document.
- In unweighted indexing sometimes, "word locations" are kept as a part of searchable datastructure.

Types of Automatic Indexing :-

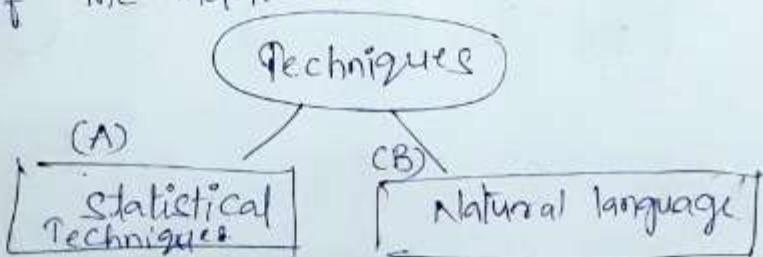
Types of Automatic Indexing



iii) Indexing by term :-

- The "terms" of original item are used as basis of "index process".

There are a major technique used for creation of the index dire:



(A) Statistical Technique :

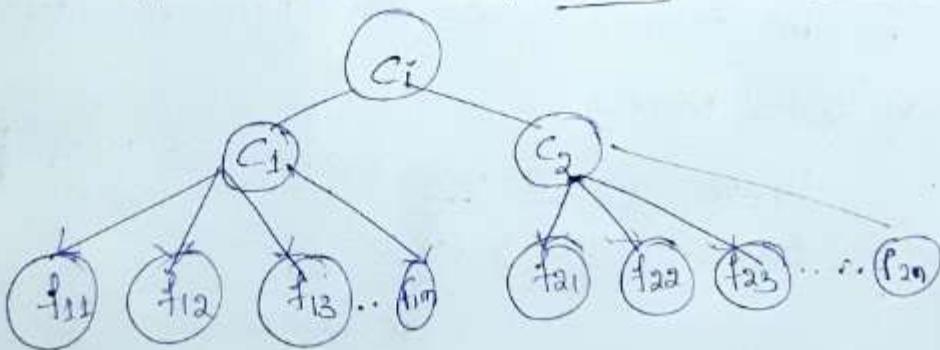
- The statistical techniques can be based on the "vector Models" and "probabilistic models" with a special case being is "Bayesian models".

- All the models (i.e. Vector model, probabilistic model and Bayesian model) classified as "statistical" (f6) because their calculation of weights use the
- statistical information such as "frequency of occurrence of words" and their distributions in the searchable database.

Bayesian Model :-

A Bayesian model is a "directed acyclic graph" in which the node indicate the random variable.

It is the combination of "nodes" and "Arcs".



where $C \rightarrow$ concept, $f \rightarrow$ features.

To calculate the probability of " C " from given " f_{ij} ".

C_i from given f_{ij}

To perform the calculation 2 sets of probabilities are needed:

1. The Prior Probability $\cdot P(C_i)$.

means the item is relevant to concept C .

2. The "Conditional probability"

The conditional probability is indicated by $P(f_{ij}/c_i)$.

where, $P \rightarrow$ Probability

$f_{ij} \rightarrow$ features

$c_i \rightarrow$ Concept

(††)

→ The features F_{ij} where $j=1, m$ are present in an item means the item contains topic c_i .

The automatic indexing task is used to calculate the "Posterior probability":

$$\text{i.e. } P(c_i/f_{i1}, \dots, f_{im}) \quad \text{when } j=1, m$$

The probability of the item that contains concept(c_i) given the presence of features(f_{ij}).

The Bayes inference / Bayesian formula used for automatic indexing is:

$$P(c_i/f_{i1}, \dots, f_{im}) = \frac{P(c_i) P(f_{i1}, \dots, f_{im}|c_i)}{P(f_{i1}, \dots, f_{im})}$$

(or)

$$P(c_i/f_{i1}, \dots, f_{im}) = \frac{P(c_i) P(f_{i1}, \dots, f_{im}|c_i)}{P(f_{i1}, \dots, f_{im})}$$

\Rightarrow The result of search by posterior's, Req
the Bayes rule can be simplified to linear decision rules

$$g(C_1 | f_{11}, \dots, f_{1m}) = \sum_k I(f_{ik}) w(f_{ik}, c_i)$$

(18)

Where,

"g, w" are interpreting the coefficients

i.e. $w \rightarrow$ weight corresponding to each feature / concept pair.
 ex: index term .

$g \rightarrow$ "g" as the sum of the weights of the features.

$I(f_{ik}) \rightarrow$ an indicator variable that equals 1

if f_{ik} is present in the item (otherwise it is equals zero).

$(g/f_{11}, \dots, f_{1m})$ \rightarrow Posterior probability.

(B). Natural language processing :-

\rightarrow The another approach to defining "indexes to items" is called Natural Language processing.

\rightarrow The [DR-LINK] (Document Retrieval through Linguistic knowledge) System processes items at the morphological, lexical, semantic, syntactic and discourse levels.

Each level uses information from the previous level to perform its additional analysis.

iii) Indexing by concept:-

(79)

→ The basis for "concept indexing" is that, there are "many ways to express" the same idea.

→ Indexing by term treat each of these occurrence as a "different index" and then use the other query expansion techniques to expand a query to find the different ways the same thing has been represented.

concept indexing:

Concept indexing determines a "canonical set of concepts" based upon a list set of "terms".

The concept indexing is also called as "Latent Semantic Indexing" because it indexing the latent semantic ~~meaning~~ information in item.

In concept indexing, it does not associated with each "concept" (i.e. words / ^{set of} words that can be used to describe it) but it is a mathematical representation (Ex. a vector / context vector).

Example of Concept Indexing:-

The example of concept indexing is the "Matchplus system" developed by HNC Inc.

→ The words items & queries are represented by high dimensional vectors called "context vectors". i.e. (high dimension at least 300 dimensions - 5000).

The Matchplus system uses "neural networks" to facilitate machine learning concepts of relationships and sensitivity to similarity of Ue.

(80)

The two neural networks are used in concept indexing, age

1. One Neural network : (Stem context vectors)

One Neural network learning algorithm generates "Stem context Vectors" that are sensitive to similarity of Ue.

2. Second Neural network :

It performs query modification based upon user feedback.

Context Vector / Vectors

Context Vector :

Words stems, items and queries are represented by "high dimensional" Vectors called "Context Vectors".

To define context vectors, a set of "n" features are selected on an ad hoc basis.

i.e. for any word stem k, its context vector v^k is an n-dimensional vector with each component j interpreted as follows:

v^k "Positive" if k is strongly associated with feature j.

$v^k \approx 0$ if word 'k' is not associated with feature j.

v^k "negative" if word "k" contradicts feature j.

Once the context vectors for items are determined, they are used to create the index for an item. (21)

iii) Multimedia Indexing :-

- Indexing associated with multimedia differ from the previous discussions of indexing.
- The automated indexing take place in "multiple passes" of the information "versus" direct conversion to the indexing structure.

first pass : The first pass in most cases is a "conversion" from the analog input mode into a "digital structure".

The digital structure algorithms are "extract" the different modalities that will be used to represent the item.

Next pass :-

Indexing video (or) images can be accomplished at the "raw data level".

Ex:- The aggregation of raw pixels.

i.e. The raw pixels, (distinguish) contains the primitive attributes such as "color" and "luminance" and at semantic level meaning full objects are recognized.

Ex: an airplane in Image/Video frame.

Video processing:-

- The system will periodically collect a frame of video input for processing.
- It might compare that frame to the last frame captured to determine differences between the frames.
- If the difference is below the threshold it will discard the frame.
Ex: Virage.

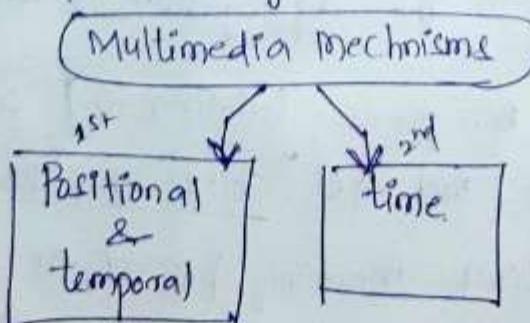
Image processing:-

In Image Processing,
of Objects with in the images.

- Semantic level indexing of "several" pattern recognition
- Ex:
 - MIT's PhotoBook ,
 - IBM's QBIC &
 - Multimedia Database from Informix | Virage .

finally :-(multimedia mechanisms)

The Multimedia items need some mechanism to correlate the different modalities during the search .



first mechanism:

* Positional & temporal mechanism:

This mechanism is used when the modalities are interposed / interleaved in a linear sequential composition.

Example:-

A document that has "images" (or) "audio" inserted can be considered a linear structure.

* Second mechanism: i.e "time" :-

The second mechanism is based upon time because the modalities are executing concurrently.

i.e - Ex:- Television , SMIL

The typical Video source off television is inherently a multimedia source.

It contains video, audio and potentially closed captioning.

SMIL → synchronized multimedia integration language.

It is a mark-up language designed to support multimedia presentations that integrate text with audio, images and video.

* Information Extraction :-

Extraction means → removal,
taking out, drawing out, pulling
out etc.

→ Information Extraction is a task of extracting the 84 ^{int} structured information from unstructured and / or semi-structured machine-readable documents.

→ There are two processes associated with information extraction:

- i) Determination of facts
- ii) Extraction of text

i) Determination of facts:-

Determination of facts go into "structured fields" in a database.

In this case, Only a subset of the important facts in an item may be identified and extracted.

ii) Extraction of text:-

Extraction of text can be used to "summarize an item".

In summarization, all the major concepts in the item should be represented in the summary.

→ The process of extracting facts to go into indexes is called "Automatic file build".

It's "goal" is to process "Incoming items" and "extract index terms" that will go into a structured database.

- An Information Retrieval System's "goal" is to provide an in-depth representation of the contents of an item. (85)
- An Information Extraction System: Only analyzes those portions of a document that potentially contain information relevant to the extraction criteria.

Metric to compare Information Extraction:-

The same previously defined measures of precision and recall are applied with slight modifications to their meaning.

Precision:- Precision refers to how much information was extracted accurately Versus the total information extracted.

Recall:- Recall refers to how much information was extracted from item Versus how much should have been extracted from the item.

i.e. The amount of correct and relevant data extracted Versus the correct and relevant data for the item.

Additional Metrics :-

Overgeneration fallout

The additional metrics used for information extraction are: "Overgeneration" & "fallout".

Ovengeneration :-

Ovengeneration measures the amount of irrelevant/non-relevant information that is extracted.

Fallout :-

Fallout measures how much a system assigns incorrect slot fillers, as the number of potential incorrect slot fillers increases.

These measures are applicable to both "human" and "automated" extraction process.

- ⇒ The best source of analysis of data extraction is from the "Message Understanding Conference proceedings".
- ⇒ The conferences were held in 1991, 1992, 1993 & 1995.
- ⇒ The conferences are sponsored by the Advanced Research Project Agency.

Objective of data extraction :- "Slot"

Objective of data extraction is update structured database with additional facts.

- ⇒ The term "slot" is used to define a particular category of information to extract.
- ⇒ Slots are organized into templates / semantic frames.
- ⇒ Information extraction requires multiple levels of analysis - of text of the item.
i.e. It must understand the words and their context.

→ This focusing is very similar to Natural Language processing.

Data Structure

① Introduction to Data structure :-

81

→ A Data structure is a "storage" that is used to store and organize data.

It is a way of arranging data on a computer so that it can be accessed and updated efficiently.

There are usually "two major data structures" in any information system.

first/One major Data structure:-

→ One major Data structure "Stores and manages" the received items in their "Normalized form".

→ The process of supporting this structure is called the "document Manager".

Other Major Datastructure:-

Other major Data structure contains the "processing tokens", and 'associated data' to support search.

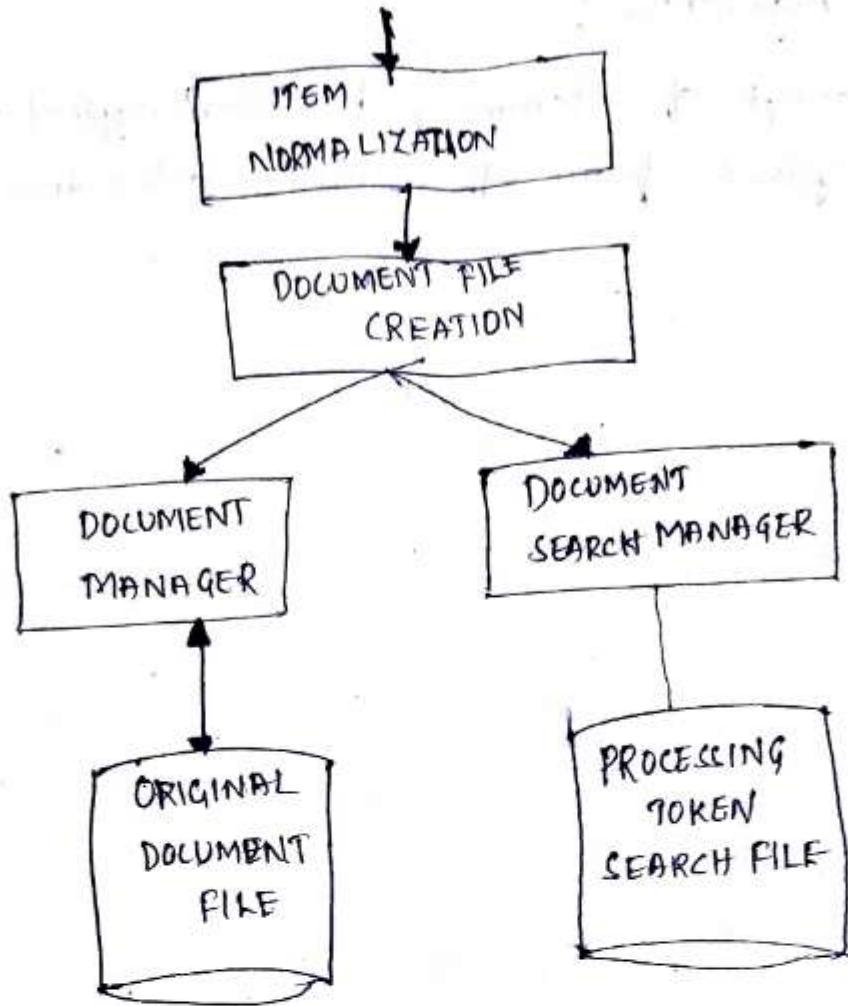


fig: Major Datastructures

- The above figure expands the document file creation function (it was discussed in 1st chapter).
- The results of a search are references to the items that satisfies the search statement, which are passed to the document manager for retrieval.
- To understand this datastructures background, the reader / user should pursue a text on finite automata and language (regular expressions).

* Inverted file structures:-

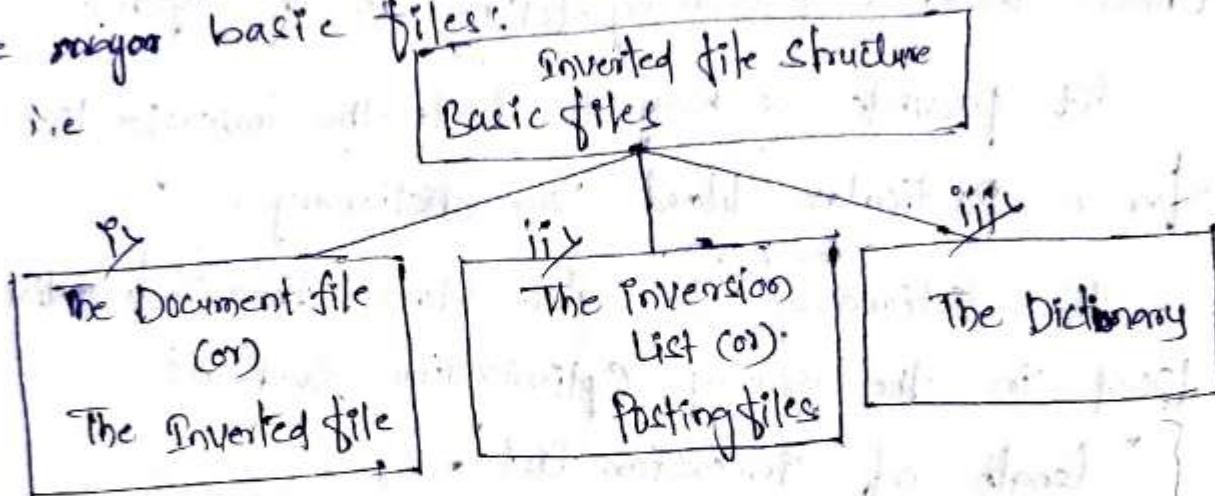
(B9)

The most common data structure used in both Database management and Information Retrieval systems is the inverted file structure.

It is used to minimize the secondary storage access when multiple search terms are applied across the total database.

The Inverted file structures, are composed of three major basic files:

i.e



* Inverted file / Document file:-

The name Inversion "Inverted file" comes from its underlying methodology of storing an inversion of the documents.

from the Inversion of document perspective, for each word, a list of documents for which the word is found.

iii) The Inversion list :-

(a)

Inversion list sometimes called as "posting list".

In this each document in a system give a

"Unique numerical Identifier".

The identifiers are stored in the inversion list.

iv) The Dictionary:-

The Dictionary is typically a stored list of all unique words (processing tokens) in the system.

It provides a way to locate the inversion list for a particular word via dictionary.

The dictionary can also store other information used in the query optimization such as

"length of inversion list".

DOCUMENTS

Doc #1, Computer, bit, byte
Doc #2, memory, byte
Doc #3, Computer, bit, memory
Doc #4, byte, Computer

DICTIONARY

bit (2)
byte (3)
computer (3)
memory (2)

INVERSION LISTS

bit - 1, 2, 3

byte - 1, 2, 4

computer - 1, 3, 4

memory - 2, 3

fig : Inverted file structure

from the figure/diagram :-

(9)

- The Document contains information of words (such as bit, computer, bytes) i.e. how many words are there in a particular documents.
- The Dictionary contains how many number of times the particular word is appear in the document. It is written inside parenthesis () .
- The Inversion List Contains in which document the particular word is appearing.

finally,

The inverted file structure uses some silent features

i.e.

1. It increases precision and recall.
2. It also uses browse capabilities such as ranking and Zoring.
3. It uses the NLP (Natural Language processing)
4. It is mainly used to store the concepts and relationships.
5. Rather than using a dictionary to point the inversion list B-Trees can be used.

In B-Tree Inversion list may appear at
"leaf level / leaf node".

(92)

→ Bxix A B-tree of Order "m" is defined as:

- * A "root node" with between "2 and $2m$ keys".
- * All other Internal nodes have between "m" and " $2m$ " keys.
- * All keys are kept in Order from smaller to larger.
- * All leaves are at the same level (or) differ by at most one level.

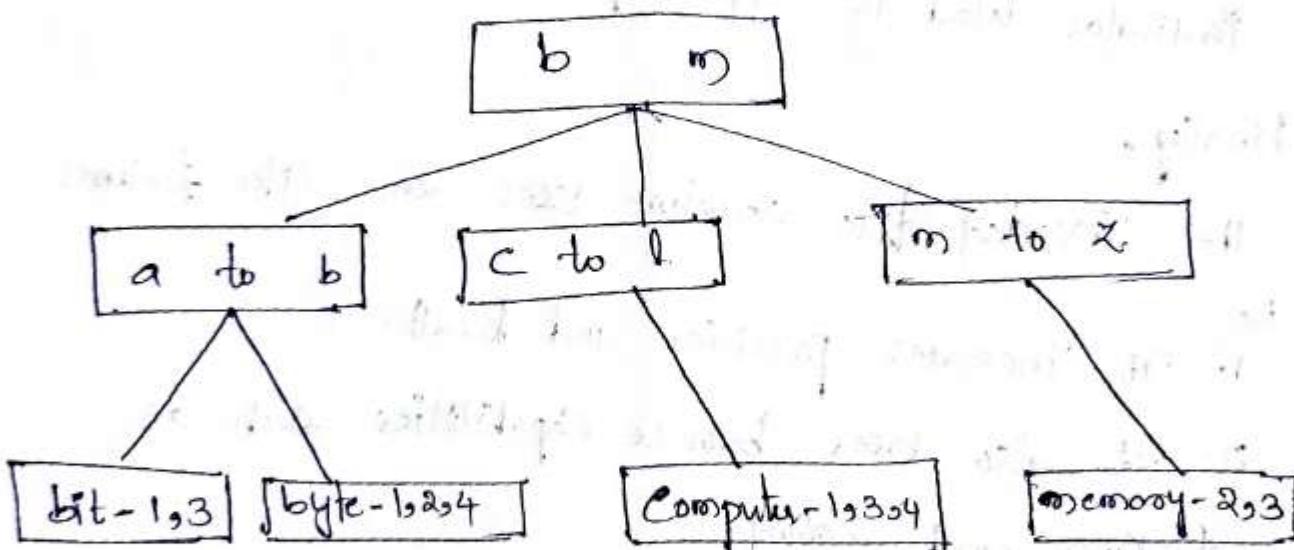


Fig: B-Tree Inversion list.

*

N-Gram Data Structures

(93)

- A Variant of the 'Searchable data structure' is N-Gram Data Structure, that breaks processing tokens into smaller string units and uses the token fragments for search.
- N-grams can be viewed as a special technique for conflation (stemming). It is used as a unique data structure in information systems.
- N-grams are 'fixed-length' consecutive series of "n" characters.

Unlike stemming,

i.e. The stemming determines the stem of a word that represents the semantic meaning of the word but n-grams do not care about semantics.

Instead of that, they are algorithmically based upon a fixed number of characters.

→ The Searchable data structures is transformed into Overlapping n-grams, which are then used to create the Searchable data base.

Example : Bigramme, Trigramme and pentagramme for the word phrase "sea colony".

(94)

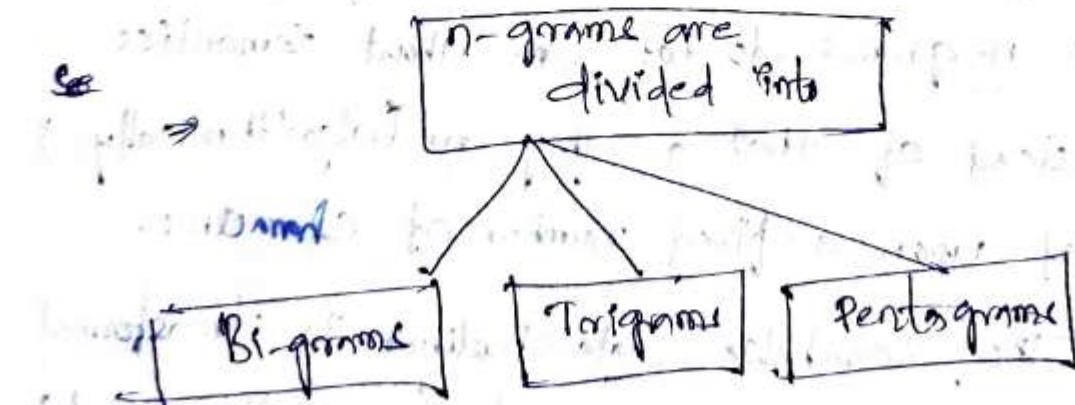
→ for n-grammes with $n > 2$, some systems allow interword symbols to be part of the n-gram.

The symbol "#" be used to represent the interword symbol.

i.e. # → interword symbol

⇒ The # may be represented any set of symbols such as "blank, period, semicolon, colon" etc.

⇒ It is possible that same "n-gram" can be created "multiple times" from a single word.



For example :-

(Q5)

Bigrams, trigrams and pentagrams for the word phrase "Sea colony".

* Bigrams for word "Sea colony"

Word Bigram
Sea → se ea }
colony → co ol lo on ny } (Bigrams
 no interword symbols)
∴ [se ea co ol lo on ny]

* Trigrams for word "Sea colony"

Word Trigram
Sea → sea }
colony → col olo lon ony } (Trigrams
 no interword symbols)
∴ [sea col olo lon ony]

* Pentagrams for word "Sea colony"

Word Pentagram
Sea → sea (no change) } (Pentagrams
Colony → colo olon lony } no interword
 symbols)

∴ [sea colo olon lony]

- * Trigramme with Fosterblad symbol "#" for word "sea colony"
- 96
- Place the Trigramme inside bigrams and place # intervord symbol at both ends.
- i.e. $\# \underset{\substack{\text{trigrams} \\ \text{bigram}}}{\text{se sea ea}}$ # $\underset{\substack{\text{bigram} \\ \text{intervord symbol}}}{\text{co col olo lon ony ny}}$ #
- $\underset{\substack{\text{intervord symbol}}}{\text{co col olo lon ony ny}}$ #
- $\therefore \boxed{\# \underset{\substack{\text{trigrams} \\ \text{bigram}}}{\text{se sea ea}}} \# \underset{\substack{\text{trigrams} \\ \text{intervord symbol}}}{\text{co col olo lon ony ny}} \#$

- * Pentagramme with intervord symbol "#" for word "sea colony"
- place the Pentagrams inside trigrams and place # intervord symbol above trigrams at both ends.
- Sea \rightarrow $\underset{\substack{\text{4-grams}}}{\text{colo}}$ $\underset{\substack{\text{5-grams}}}{\text{olon}}$ $\underset{\substack{\text{4-grams}}}{\text{olony}}$ \rightarrow 4-grams
- colony \rightarrow $\underset{\substack{\text{5-grams}}}{\text{colon}}$ $\underset{\substack{\text{4-grams}}}{\text{olony}}$ \rightarrow 5-grams
- $\# \underset{\substack{\text{4-grams}}}{\text{colo}} \underset{\substack{\text{5-grams}}}{\text{colon}} \underset{\substack{\text{5-grams}}}{\text{olony}} \underset{\substack{\text{4-grams}}}{\text{olony}} \#$
- $\therefore \boxed{\# \underset{\substack{\text{4-grams}}}{\text{sea}} \# \underset{\substack{\text{5-grams}}}{\text{colo}} \underset{\substack{\text{5-grams}}}{\text{colon}} \underset{\substack{\text{4-grams}}}{\text{olony}} \underset{\substack{\text{4-grams}}}{\text{olony}} \#}$

Example 2

Read the details & find Bigramme, Trigramme & Pentagramme.

(q7)

Use of N-gramme

① The first use of n-gramme was at the time of World War-II, it was used by "crypto grapher".

The Crypto grapher namely "Fletch Pratt" states / defines the "bigram and trigram tables" so that any cryptographer can dismember simple substitution cipher taking partition/division.

Another Crypto grapher namely "Adamsen" describes the use of Bigramme, as a method for conflating terms (or) stems.

i.e. The conflating terms do not follow normal definition of stemming, because the "stemming" has "semantic meaning" whereas "n-gram" has "no semantic meaning".

② Another Major use of N-gramme (particularly trigrams) is in "Spelling error detection and correction".
i.e.

i.e There are 4 categories of errors.

(98)

Error Category

- * Single character Insertion
- * Single character Deletion
- * Single character Substitution
- * Transposition of two adjacent characters

- Example → Computer

~~compuer~~ Computer → Insertion at single character i.e u

comptee → Single character deletion i.e u

comptee → substitution i.e u → t

comptee → computer i.e ut

fig: categories of spelling errors.

③ N-Gram patterns also used for "Identifying the language of the item".

④ Trigrams have been used for "text compression" to manipulate the length of index terms.

① Advantages & Disadvantages of N-grams:-

Advantages :-

i) longer n-grams, ignoring the word boundaries.

ii) ^{Theorem} A finite limit on the ~~no. of~~ no. of searchable tokens.

$$\text{Max seq}_n = (\lambda)^n$$

Where, $\text{Max seq}_n \Rightarrow$ Maximum no. of Unique "N-grams" that can be generated.

- i.e MaxSeq \rightarrow calculated as a function of 'n'. (9)
 n \rightarrow indicates the length of N-gramme.
 λ \rightarrow indicates the no. of processible symbols from the alphabet.

Disadvantages :-

- i) The disadvantage of N-gram is "increased size of inversion list" that store the linkage datastructure.
- ii) N-gramme are poor representation in "concepts and their relationships".

* PAT Data Structure

The name PAT is the short form of ^{PATRICIA} Patricia trees.

PATRICIA stands for

- P \rightarrow Practical
- A \rightarrow Algorithm
- T \rightarrow To
- R \rightarrow Retrieve
- I \rightarrow Information
- C \rightarrow Coded
- T \rightarrow In
- A \rightarrow Alphanumeric.

i.e Practical Algorithm To Retrieve Information Coded In Alphanumeric.

→ The PAT Data structures were described by (100)

frakes - 92

Ronni - 83

Knuth - 73 &

morrison - 68 has "Patricia trees".

Concept.

→ The Original Concepts of PAT tree data structure were described as "Patricia trees", and they are used for searching text and images.

These applications are in "Genetic database".

The PAT datastructure represents the information in arrays.

i.e. *

PAT trees &

* PAT arrays.

PAT trees and PAT arrays are addressing the different View of continuous "text" \rightarrow Input".

The input stream transformed into a searchable data structure consisting of "substrings".

Creation of PAT trees:-

In creation of PAT trees each position in the Input strings is the "Anchor Point" for a sub-string.

The substring start at "Anchor point" and include all new text upto end of the input.

→ All substrings are 'Unique'.

(101)

Substring :-

- * A substring can start at "any point in text" and uniquely indexed by the "starting location" and length.
- * If all the strings are to the end of the input, Only the starting location is needed.
Because the length is different from the location and the total length of the item.
- * It is possible to have a substring go beyond the length of the input stream by adding the additional null characters.

These substrings are called listing
listing means semi-infinite string).

Example : Some possible listing for an input text given below.

TEXT :

String 1

String 2

String 3

String 4

String 5

String 6

String 7

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30
Economics for Warsaw is Complex. → with space count
Economics for Warsaw is Complex. → without space count
conomics for Warsaw is complex.
omics for Warsaw is complex.
for Warsaw is complex. (with out space count)
is complex. (with space count)
ex. (with space count).

fig: Examples of listing.

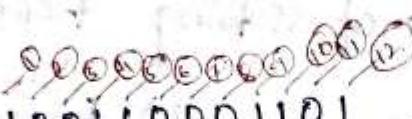
In PAT trees,

(102)

* The key values are stored at the "leaf nodes"
in the PAT tree.

for a text input of size "n" there are "n"
leaf nodes and "n-1" at most higher level nodes.

Example Of the strings Used in generating a
PAT.

INPUT:  (∴ Binary representation
of HOME)

Sistring 1 1001....

H = 100

Sistring 2 001100....
② place

O = 110

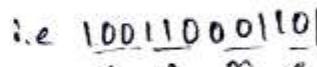
Sistring 3 01100....
③ replace

M = 001

Sistring 4 11....

E = 101

Sistring 5 1000....

i.e. 

Sistring 6 000....

Sistring 7 001101

Sistring 8 01101

fig: Sistings for input "100110001101".

→ the word home produces the input

100110001101

\Rightarrow No. of Sistring = 8 (node's)

$$\text{i.e. } \boxed{n=8}$$

\Rightarrow We get $\boxed{n-1} = 7$ levels

$$8-1 = 7$$

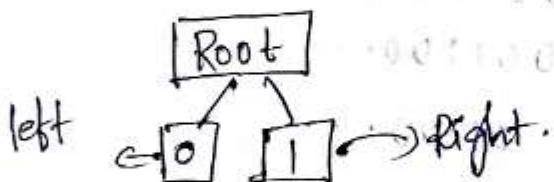
i.e. for ' n ' no. of node, we get the ' $n-1$ ' levels.

Rules for PAT Tree:-

The individual bits of sistring decides the branching

- \rightarrow The pattern with Zero (0-bit) branching 'left'.
- \rightarrow The pattern with one (1-bit) branching 'Right'.

i.e.



Note

Take any query | sistring to reach external node.

Query = 001100
 ↑
 bit
 first

We first inspect bit-1 (if it is '0', branching will be done at left)
then bit-2 (if it is '0' branching will be done at left)

bit-3 (if it is '1' branching will be done at Right)

bit-4 (if it is '1' branching will be done at right)

bit-5 (if it is '0' branching will be done at left)

Once we reach desired node we have to make one final comparison with one of sistring.

i.e Example: INPUT : 100110001101

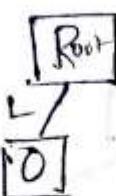
nodes: $n=8$, $n-1=7$ (levels)

(104)

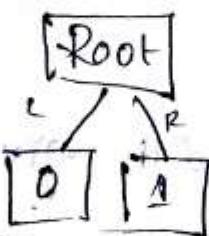
ii) Initial tree contains 'Root - R'.



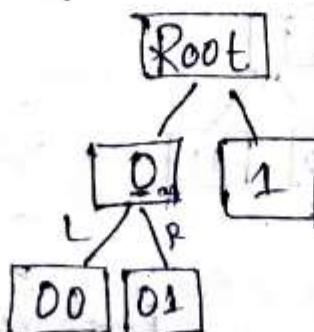
iii) add '0-bit' as left branching



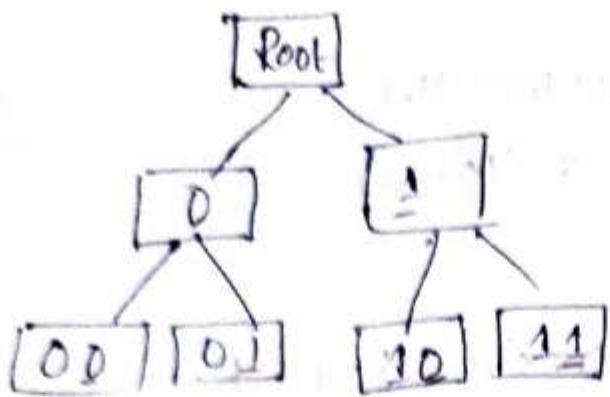
iv) add '1-bit' as right branching:



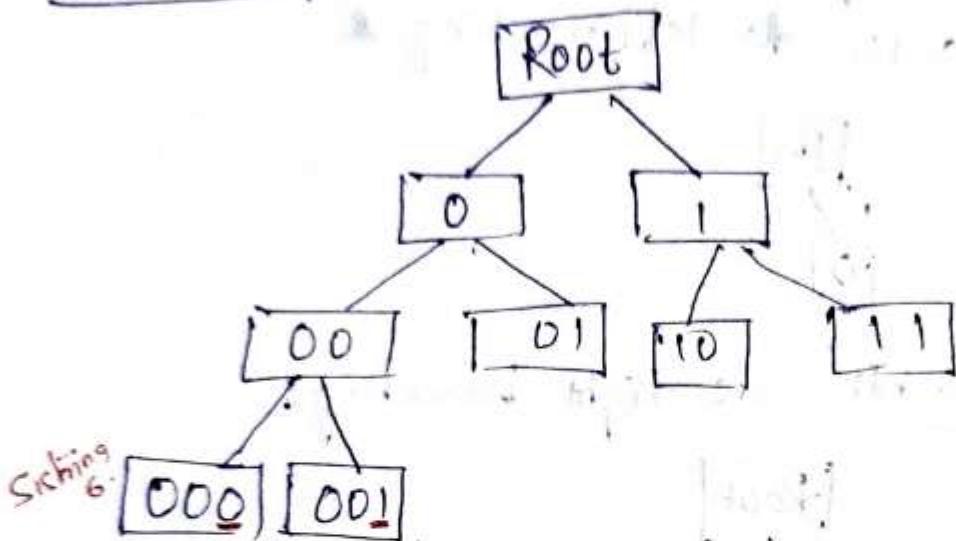
v) for desired position '0' (Zero) compare with two bits i.e again zero & one.



vi) similarly for 1 (Right most bit). Compare with two bits i.e again zero & one.

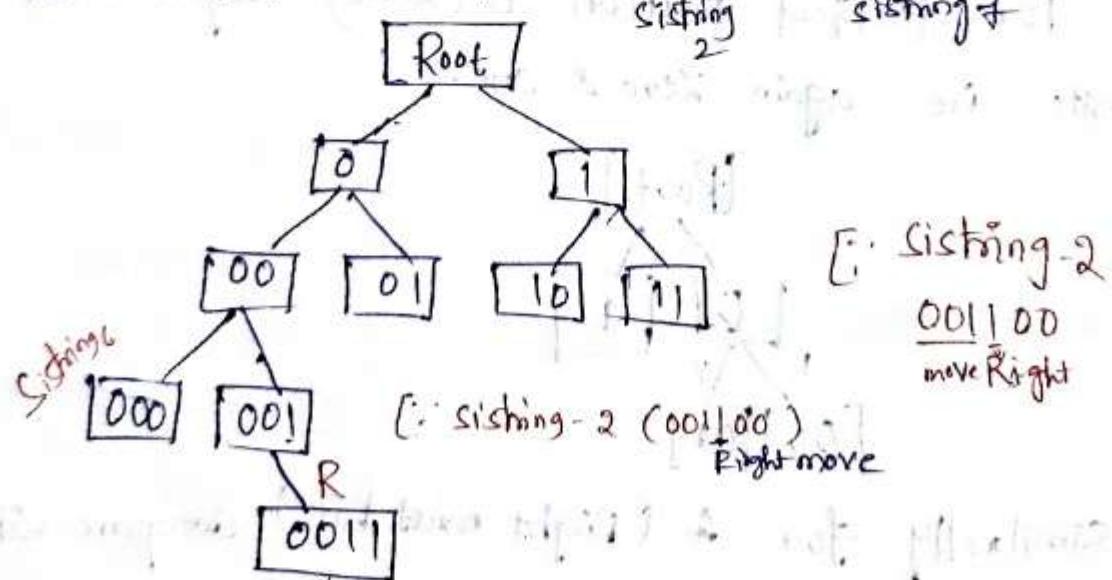


from Given strings we are going to construct PAT tree

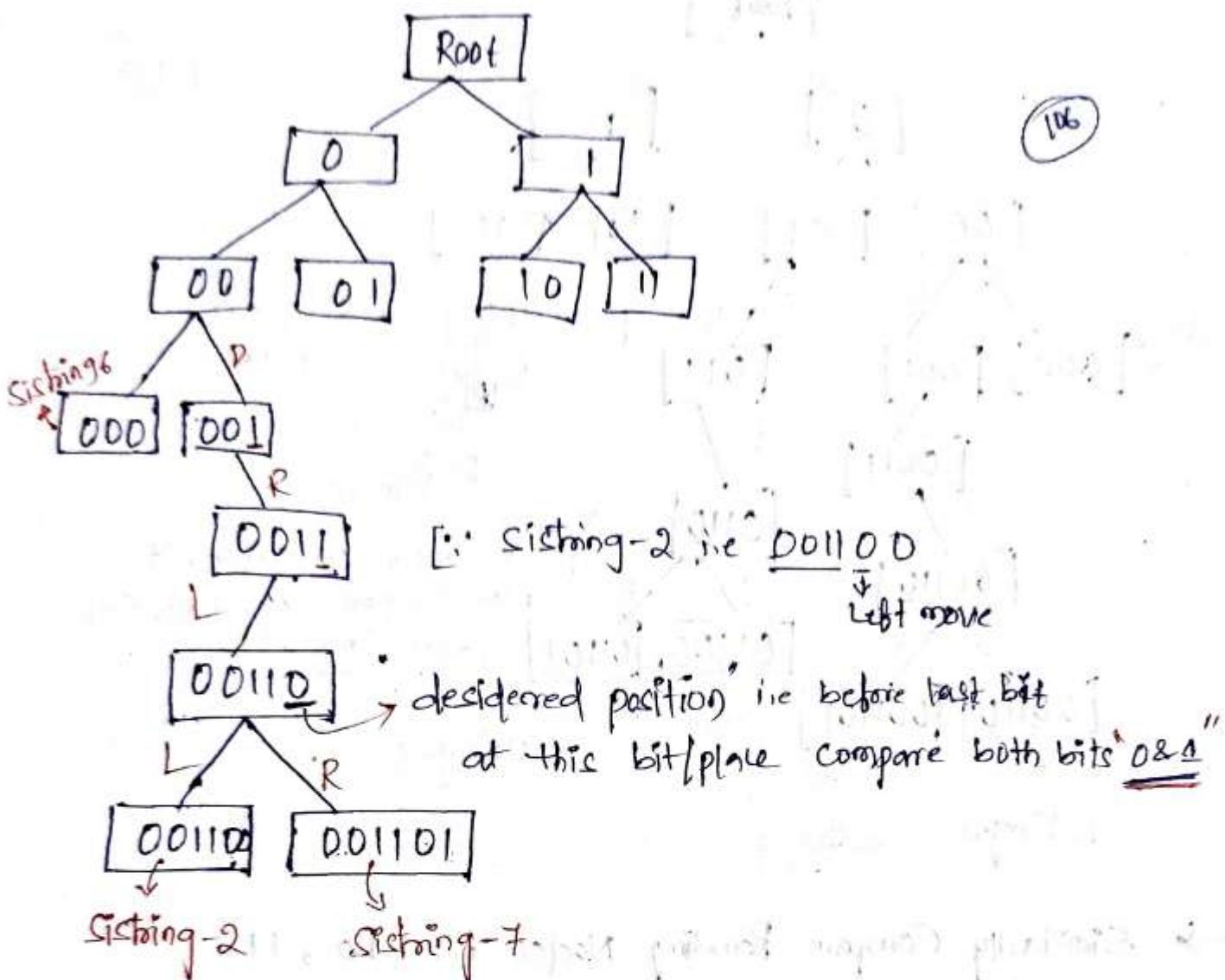


→ Now compare 001 node with given strings.

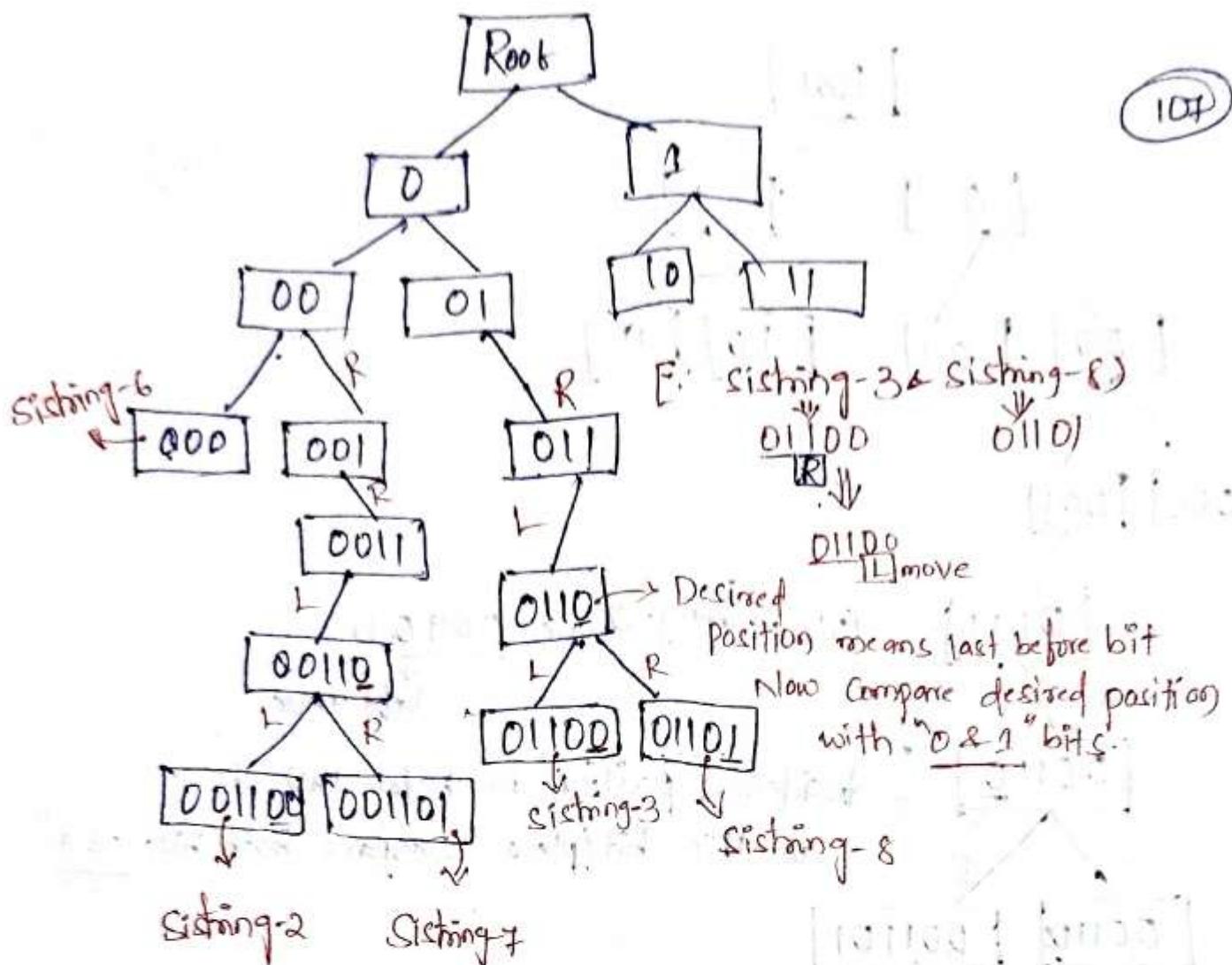
In given input strings : 001100, 001101
i.e. slicing 2 slicing 1



Next compare running bit in slicing 2. i.e. 001100
Left (zero) move



- Similarly Compare Remaining Nodes. Now comp.
i.e (01, 10, 11)
- Now compare 01 node :-
- check is there any sifting start with '01'
at sifting-3 (i.e 0100) & sifting-8 (i.e 001101)
- Start with 01.
Now construct \uparrow ^{Remaining} PAT Tree.



→ Similarly Compose Remaining Nodes i.e 10, 11.

Now compare "10" Node.

- * Check if there any 'sifting' start with 10.
- * At sifting-1 (1001) and sifting-5 (1000)
- are start with 10.

Now construct Remaining PAT tree.

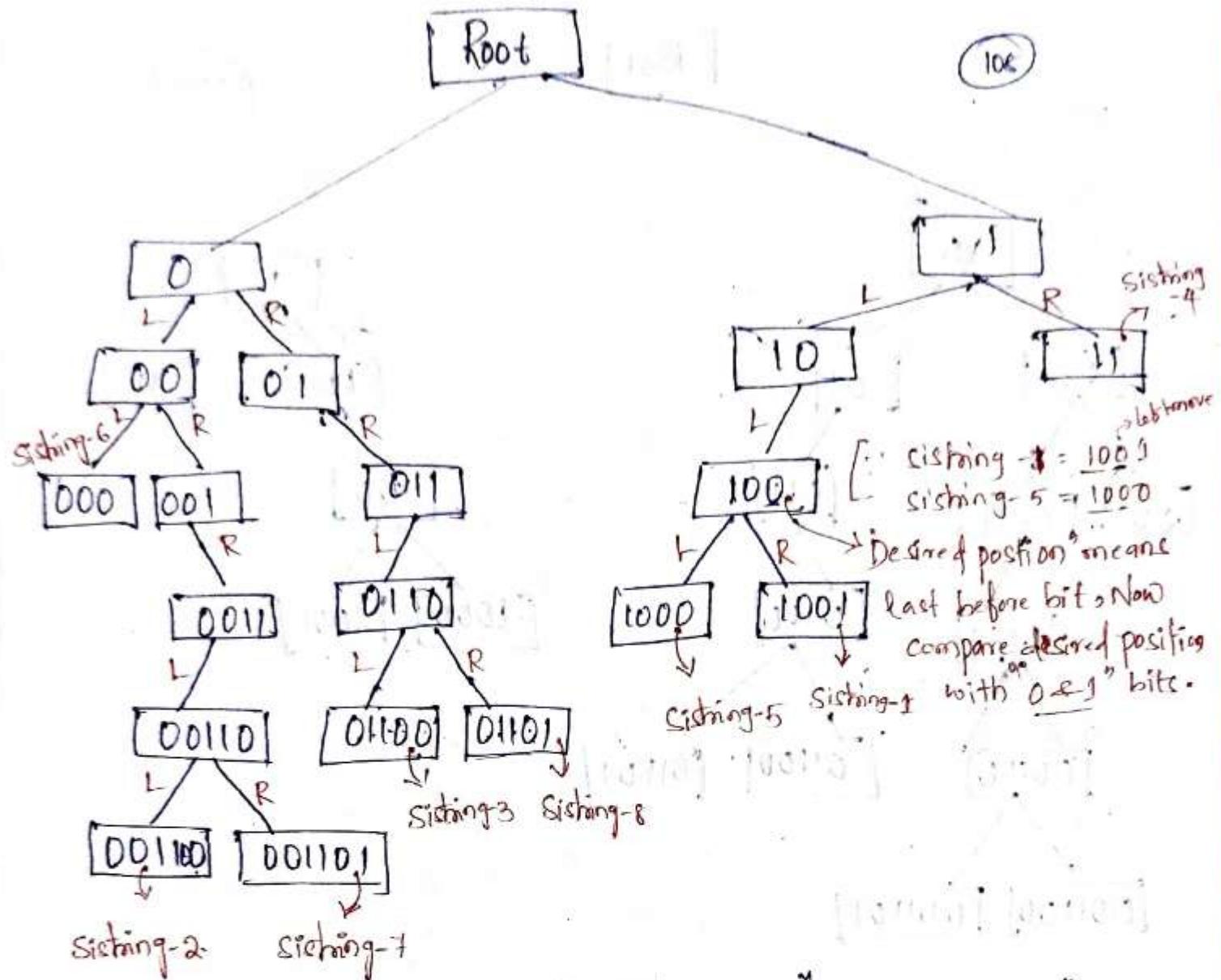
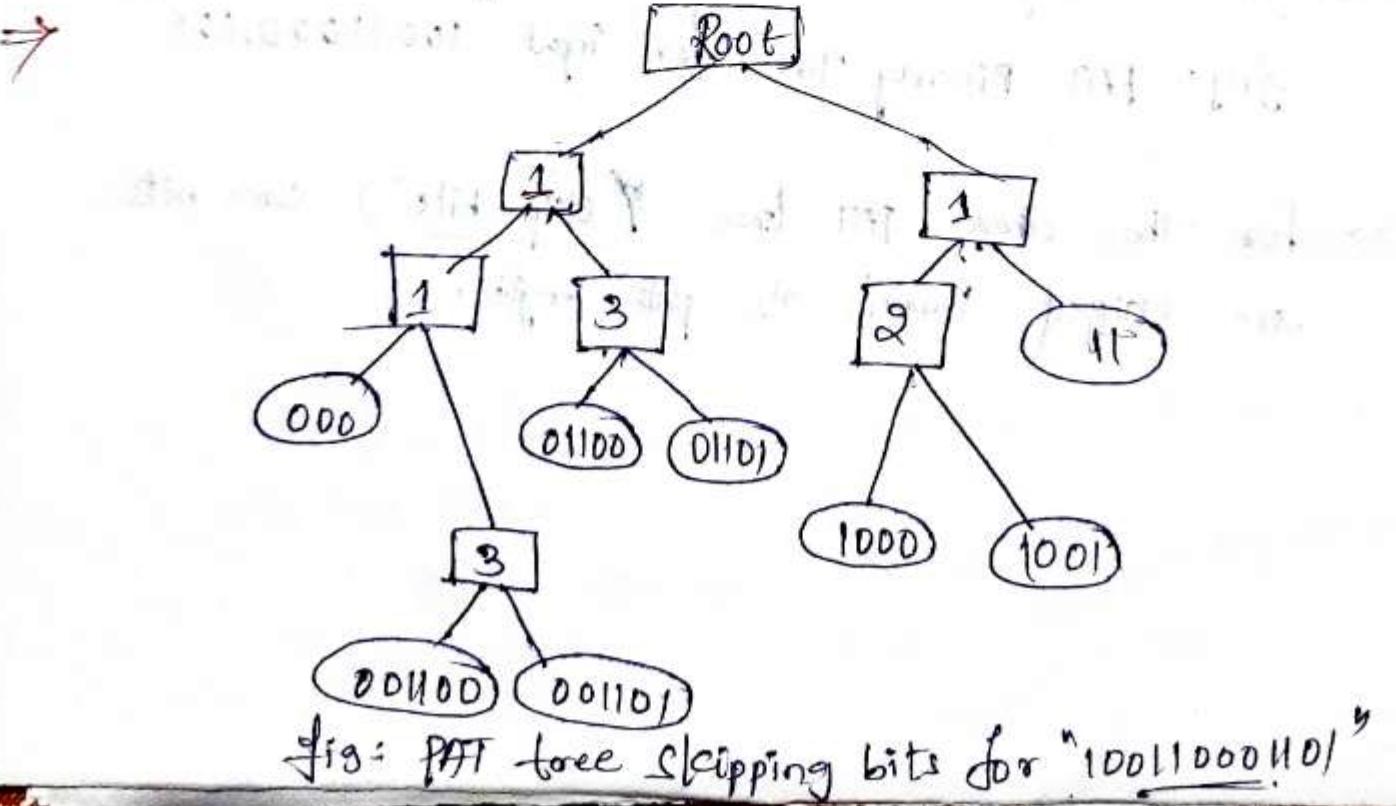
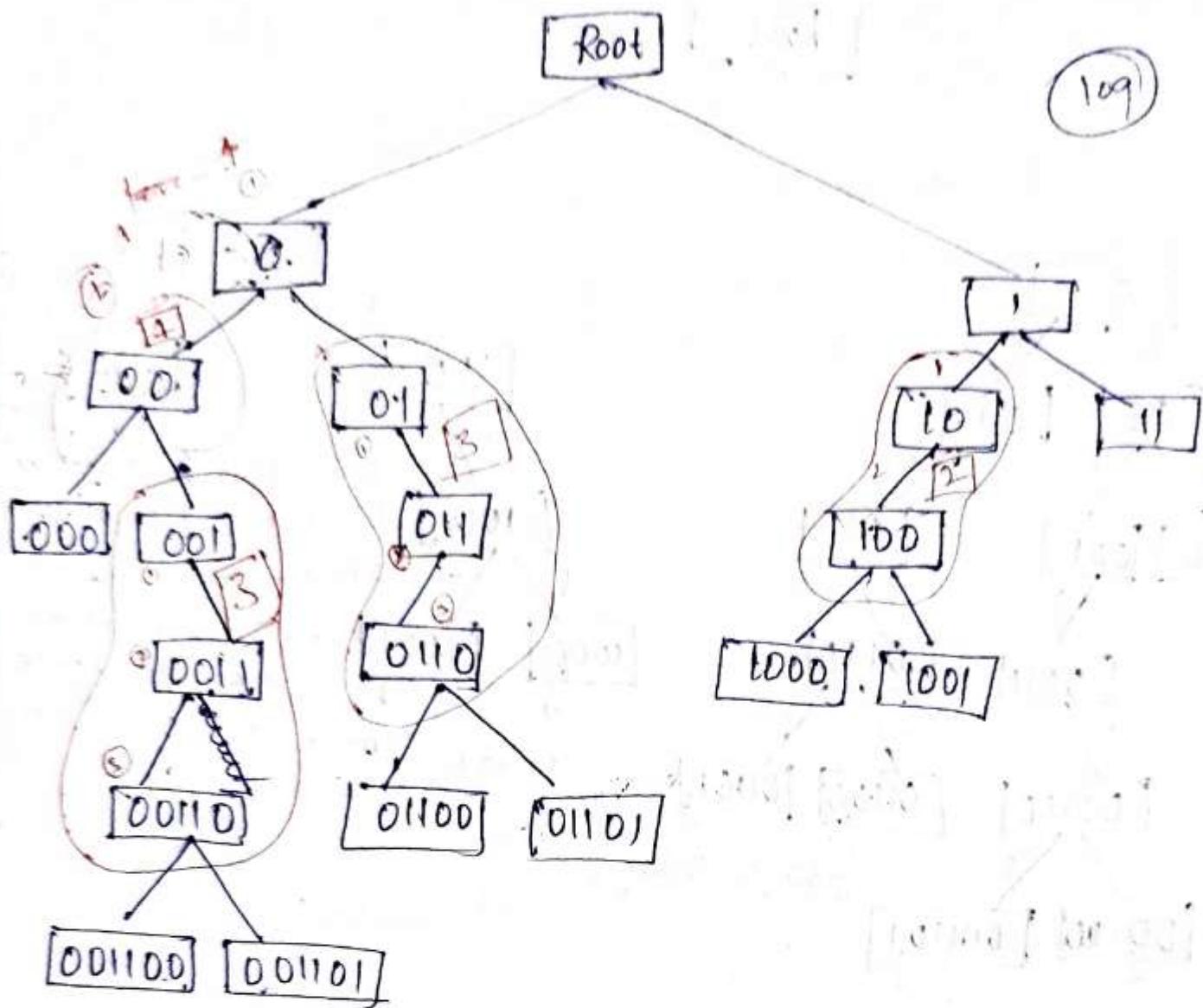


fig: PAT Binary Tree. for input "100110001101"

⇒ for the above PAT tree (skip-bit) some bits are skipped through its path length.



* Signature file structure :-

A signature is a compressed version of a database.
 All the signatures that represent the documents are kept in a file called "Signature file".

Goal :

The Goal of a signature file structure is to provide a 'fast test' to eliminate the majority of items that are 'not related' to a query.

The items that satisfies the test can either be evaluated by another search algorithm to eliminate the additional false delivered to the user to review.

The text of the items is represented in a highly Compressed form that facilitate the 'fast test'.

Because, the file structure is highly compressed and 'unordered', so it requires significantly less space than inverted file structures, and

New items can be concatenated to the end of the structure. 'Versus' significant inversion list update.

In this, the items are deleted from information database marked as deleted.

(11)

Signature file Search is a "linear scan" of (Compressed) Version of items that producing the response time linear with respect to file size.

Process creation of signature file:

The signature file typically uses the "Superimposed Coding" to create the signature of the document.

The coding is based up on the words in the item. The words are mapped into Word Signatures.

Word Signatures—

A word signature is a "fixed-length code" with a fixed number of bits set to 1.

The bit positions that are set to "1" are determined via a "hash function" of the word.

The longest ^{to long} signatures with "1" partitioned into "blocks" of size.

for example — The Block size is set at five words, the code length is 16-bits, and number of bits that are allowed to be "1" for each word is five.

i.e here * Block size is indicated by (D)

* Code length / fixed length indicated by (F)

* Bits per Word is indicated by (m)

→ distinct non common words

(42)

method - I

TEXT : Computer science graduate students study (assume block size is five words).

WORDsignature

Computer

0001 0110 0000 0110

Science

1001 0000 1110 0000

Graduate

1000 0101 0100 0010

Students

0000 0111 1000 0100

Study

0000 0110 0110 0100

Block signature :

1001 0111 1110 0110

-fig: superimposed Coding.

i.e
Here (1) Block size = "5" five words

five words are : Computer, Science, Graduate, Students, Study.

(F) code-length = "16-bits"

$$16 = \frac{4}{2} \Rightarrow 16 \text{ bits}$$

so bit length is 4,(m) no. of bits allowed to 1 is 5 (no. of bits per word)

- * The search of the signature file requires $O(N)$ search time.

for example: II - method

(115)

(e) TEXT : DataBase Management System (Assume block size is 4)
code length / fixed-length of bits = 20 bits

<u>WORD</u>	<u>signature</u>				
Data	0000	0000	0000	0010	0000
Base	0000	0001	0000	0000	0000
management	0000	1000	0000	0000	0000
System	0000	0000	0000	0000	1000
---	---	---	1000	---	---
Block Signature	0000	1001	0000	0010	1000

Here

D → Distinct non-common words $D=4$

They are: Data, Base, management, system

f → Signature size / fixed-length bits

$$f = 20 \text{ bits}$$

m → (How many bits set to 1) no. of bits per word.

$$m=4$$

The signatures ~~of file~~ have been used in following environments:

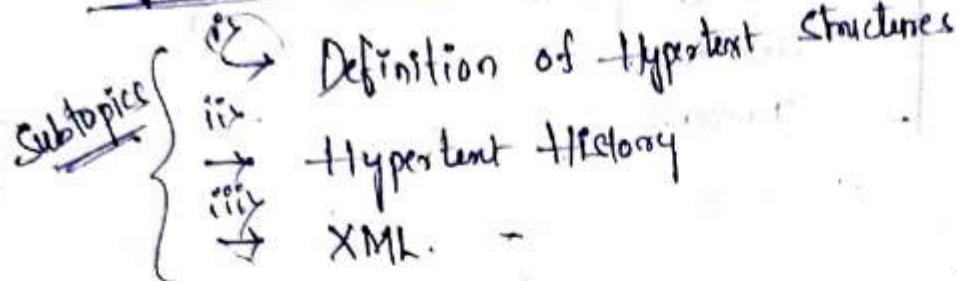
i) PC-based media size databases.

ii) WORM device Write Once Read many (Optical disks).

iii) Parallel machines.

iv) distributed environments.

⑤ Hypertext and XML Data Structures



Introduction of Basic points :

The advent / arrival of the Internet and its exponential growth and wide acceptance has introduced new mechanisms for representing information.

This structure is called "Hypertext" and it is different from traditional information storage data structure in format and use.

The Hypertext is stored in Hypertext Markup Language (HTML) and extensible Markup language (XML).

Definition of Hypertext :-

The Hypertext data structure is used extensively in the Internet environment and it requires an electronic media storage for the item.

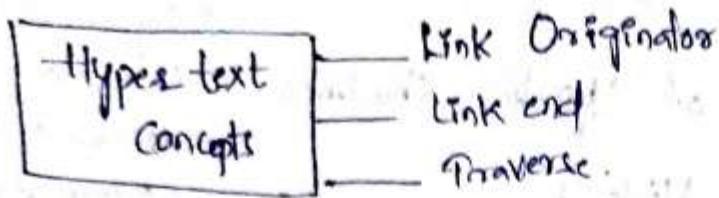
→ The HTML defines the internal structure for information exchange across the World Wide Web on the Internet.

→ The HTML document is composed of text of the form along with HTML tags that describes how to display the document.

→ Tags are expressed b/w < and > symbols.
i.e. <body> </body> <head> </head>

Hyper-text concepts :

(115)



Link Originator :-

- * It is the starting point to the link.
- * It is surrounded by the "symbols" do indicate that it is a hyper-text link.
- * The Link-Originator is an anchor.

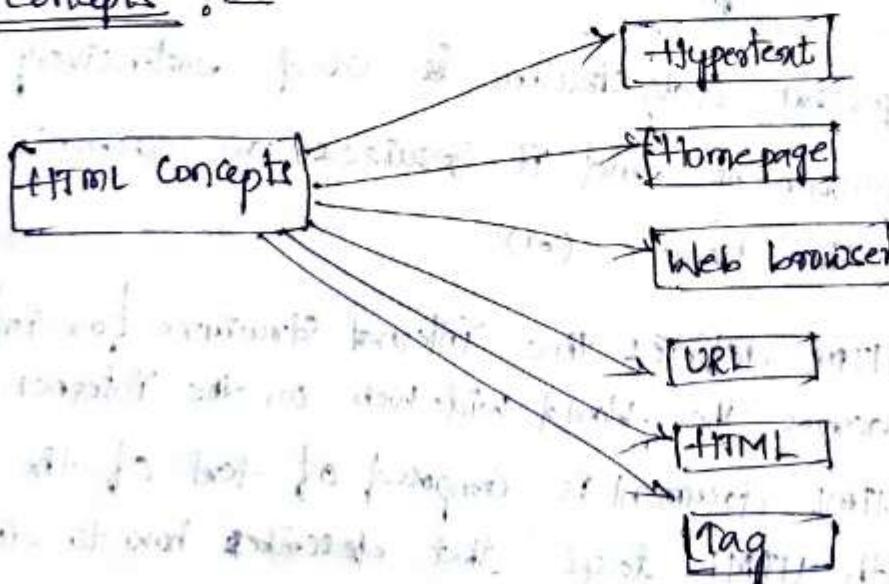
Link end :-

The other side of the Link Originator is link end.

Traverse :-

The act of travelling from a link Originator to its associated link end.

HTML concepts :-



Hyper text :-

Text ^{with} links the other text.

Homepage :-

The Root of the Hyper text structure.

Web browser :-

A tool to retrieve a document using "URL" from the "Web server".

URL :- Uniform Resource Locator.

The address of Web page is called "URL".

HTML :- Hyper text Markup language.

^{If the} HTML language is used to design the webpage.

Tags :-

HTML is made with the "tags" and they are enclosed in "angular brackets" i.e. < >.

iii) XML :-

The extensible markup language is starting to become a standard datastructure on the web.

Its first recommendation (1.0) was issued on February 10th 1998.

Hyperlinks for XML are being defined in the "Xlink" (XML linking language) and Xpoint (XML pointer language) specifications.

XML Objective is extending HTML with Semantic Information.

The logical data structure within XML defined by a
(DTD) Data Type Description

XML was designed to store and transport the data.

→ It is designed to be both ~~anthuman~~ and machine
 readable.

The ~~late~~ W3C (World Wide Web Consortium) is
 developing HTML as a suite of XML tags.

The following is a simple example of XML tagging:

<company> Widgets </company>

<city> Troy </city>

<state> NY </state>

<product> Widgets </product>

The W3C is also developing a Resource Description
 Format (RDF) for representing Properties

④ Hidden Markov Model :- (HMM)

(118)

In Hidden Markov Model the data is hidden to you / unknown to you.

Q : Why it is called Hidden Markov model ?

A : The reason is - to constructing an interface model based on the assumptions of Markov process.

The Markov process assumption is simply that

"future is independent of the past given the present".

* Andrei Andreyevich Markov (1856-1922)

- Andrei Andreyevich markov was a Russian Mathematician. He is best known for his work on stochastic processes.

- A primary subject of his research later became known as "Markov chains" and "Markov process".

- Hidden Markov models (HMM) have been applied for the last 20 years to solving problems in "Speech recognition", and lesser extent in the areas
 - Locating named entities (Bikel-97),
 - Optical character recognition (Barzilay-98)
 - topic identification (Kubala-97).

- More recently HMM's have been applied to information retrieval search with good results. (119)
- One of the first comprehensive and practical descriptions of Hidden Markov Models was written by Dr. Lawrence Rabiner (Rabiner - 89).

The easiest way to understand HMM is by an example.

Example :

Let's take the example of three State Markov Model of the "Stock Market".

The states will be one of the following that is observed at the closing of the market.

i.e. State 1 (S_1): Market Decreased.

State 2 (S_2): Market did not change.

State 3 (S_3): Market Increased in Value.

The movement between "states" can be defined by a "State transition matrix" with "state transitions" (this assumes you can go from any state to any other state):

$$A = \{a_{ij}\} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

$$A = \{a_{ij}\} = \begin{bmatrix} 0.5 & 0.3 & 0.4 \\ 0.1 & 0.6 & 0.3 \\ 0.6 & 0.1 & 0.5 \end{bmatrix}$$
120

Given that the Market fall (fall) on one day (state-1), the matrix suggests that, the probability of the market not changing the next day is "0.1".

then it is allow questions such as

The probability that the Market will increase for the next 4 days then fall.

This would be equivalent to the sequence of

$$\text{SEQ} = \{S_3, S_3, S_3, S_3, S_1\}.$$

In Order to Simplify this model,
Let's assume instead of the current state being dependent upon all the previous states, let's assume it is only dependent upon last state. (This order is called discrete, first order Markov chain).

Example of a sequence of states: $S_1, S_2, S_3, S_4, S_5, S_6, S_7$

→ It is calculated by the formula:

$$P(SFQ) = P[S_3, S_3, S_3, S_3, S_1]$$

$$\Rightarrow P[S_3] * P[S_3/S_3] * P[S_3/S_3] * P[S_3/S_3] * P[S_1/S_3]$$

initial state

previous state (similarity all)

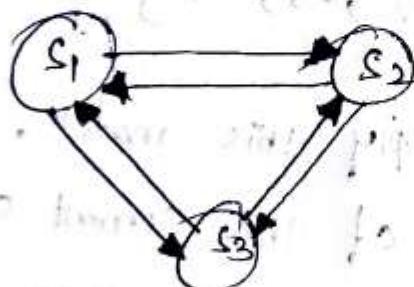
$$\Rightarrow S_3(\text{init}) * a_{3,3} * a_{3,3} * a_{3,3} * a_{1,3}$$

$$\Rightarrow (1.0) * (0.5) * (0.5) * (0.5) * (0.4)$$

$$\Rightarrow 0.05$$

→ In this equation we also assume the probability of the initial state S_3 is $\boxed{S_3(\text{init}) = 1.0}$

The following "graph" depicts the Model.



from Graph:

- * The Directed lines indicate the state transition Probabilities ϵa_{ij} .

- * There also an "implicit loop" from every state back to itself.

- * In example every state corresponded to Observable event (i.e. change in the Market).

Formal Definition of HMM :-

(122)

A more formal definition of a discrete Hidden Markov Model is summarized by Mittendorf and Schubert (Mittendorf-94) as consisting of the following.

1. $S = \{s_0, \dots, s_{n-1}\}$ as a "finite set of states"
Where "s₀" is always indicates the "initial state".
2. $V = \{v_0, \dots, v_{m-1}\}$ is a "finite set of output symbols".
3. $A = S \times S$ is a transition Probability Matrix

Where a_{ij} represents the probability of transitioning from State i to State j.

Such that $\sum_{j=0}^{n-1} a_{ij} = 1$ for all (\forall) $i = 0, \dots, n-1$.

* Every value in the Matrix is a Positive Value between 0 and 1.

* for the case;

Where every state can be reached from every other state, every value in the matrix will be non-zero.

4. $B = S \times V$ is an Output Probability Matrix, 123

Where element $b_{j,k}$ is a function determining the probability and $\sum_{k=0}^{m-1} b_{j,k} = 1$ for all $j = 0, \dots, n-1$.

5. The Initial State distribution.

From the given HMM Definition; it can be used for both "Sequence of output" and their Probabilities.

The complete specification of HMM requires

* Specification of states,

* Output Symbols

* Three Probability Measures for

① → The state transitions,

② → Output probability functions

③ → The initial states.

The distributions are frequently called A, B and π .
The following notation is used to define the model:

$$\lambda = (A, B, \pi)$$

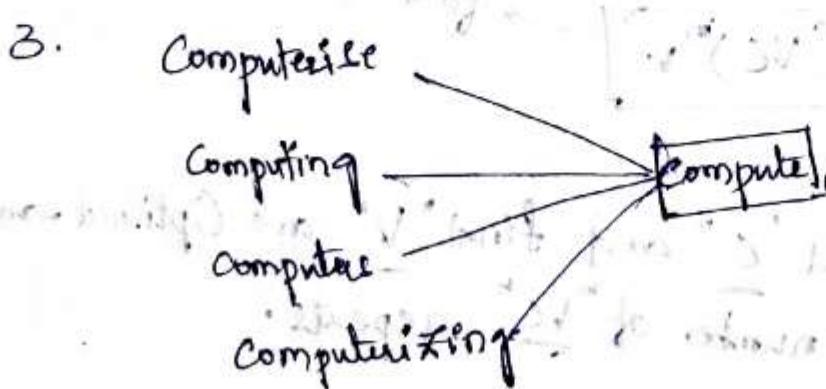
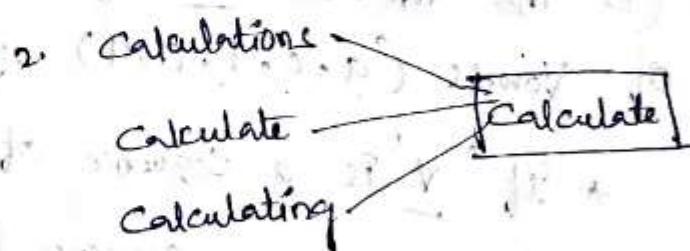
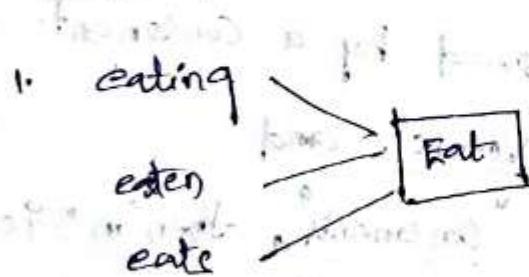
(*) Stemming - Algorithm :- 124

- Stemming means cutting (or) trimming the unnecessary data.
- Stemming is a technique used for base/ forms abstract the base form of affixes:
- Affixes indicate the ending terms.
- Stemming Programmes are commonly referred as Stemming Algorithms and Stemmers.

The main goal is to improve the performance of the system.

"Stemming is just like cutting the branches of tree".

Examples :-



Subtopics of stemming algorithms are:

(125)

- i) Porter stemming Algorithm
- ii) Dictionary look-up stemmers
- iii) Successor Stemmers
- iv) conclusions

ii) Porter Stemming Algorithm

* It is introduced in 1980 to perform "actions" based on "conditions".

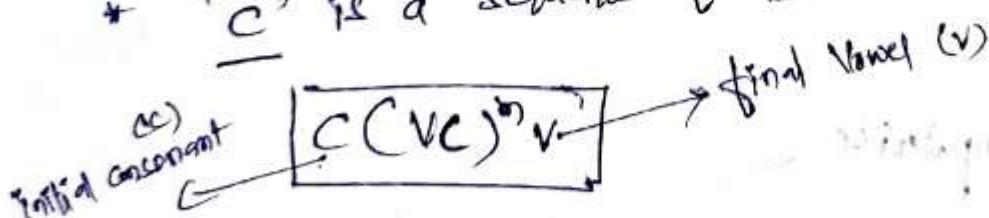
- * It is a rule based "Algorithm", based on conditions.
- * It is a rule based "stemming/porter" based on the ^{set of} stem, suffix, prefix and associate conditions.

Conditions :-

① The measure "m" of a stem is a function of sequence of Vowels (a, e, i, o, u) followed by a consonant.

+ If "V" is a sequence of "Vowels" and

+ "C" is a sequence of "consonants", then "m" is:



Where

initial "C" and final "V" are optional and "m" is the number of "VC" repeat.

Examples :-

① FREE, Why

find for given word no. of consonants and no. of vowels
 Consonants → Vowels
~~FREE~~ Why
 ↓↓↓ m=0

Hence $m=0$ Because there is no "VC" combination.

② FREES, Whose

FREES
 ↓↓↓↓↓
 C C V V S

Whose
 ↓↓↓↓↓
 C C V C V

Hence $m=1$ Because there is 1 "VC" combination.

③ Prologue, Compute, Private

for above words find no. of vowels & consonants and later find the no. of (VC) combination.

Prologue
 ↓↓↓↓↓↓
 C C V C V V

Compute
 ↓↓↓↓↓↓
 C V C C V C V

Private
 ↓↓↓ W A H
 C C V C V C V

② *X* :-

This statement is indicating "stem ending with letter X".

Example :-

Index, tax, Xerox, prefix, suffix etc.

③ *V* :-

Stem contains a Vowel.

Example :-

Pen → vowel.
 ↓↓
 CVC

④ *d* :-

Stem ending with double consonant.

Example :-

sky, dry, strong
 ↓↓↓
 CCC CCC
 ↓ double
 consonants

⑤ *o* :-

It indicates the stem ends with Consonant - Vowel - consonant sequence (C-V-C) where final consonant

is not W, X, Or Y.

Example :-

~~Pen~~, ~~Pen~~

⑥

Pen
 ↓↓↓
 CVC → C-V-C sequence
 last consonant is
 is not w, x, y.

Rules :

(128)

The Rules are divided into steps to define the Order of applying the rules.

The following are some examples of rules:

PORTER STEMMER STEPS

Step ① \Rightarrow A, B, C

(A)

i) SSES \rightarrow SS

Eg: processes \rightarrow process, stresses \rightarrow stress

ii) IES \rightarrow I

Eg: ties \rightarrow t

iii) SS \rightarrow S

Eg: class \rightarrow class

iv) S \rightarrow null

Eg: cats \rightarrow c \Rightarrow i.e. cat

(B)

B) (m>1) EED \rightarrow EE

If the stem is ending with "EED" is replaced by "EE".

Example:-

* Agreed

→ for the above word find the no. of (VC) combination so that we get "m" value. If the "m" value is ($>$) greater than 1 then it satisfies the above condition.

i.e. Agreed

$\downarrow \downarrow \downarrow \downarrow \downarrow$
V C C V V C
① ②

$m=2$: for above word 'm' value is '2'.

So that it satisfies the condition.

$m=2$

~~(m>1)~~ EED \rightarrow EE

\rightarrow (2>1) Agreed \rightarrow Agree

ii) * V * ED \rightarrow ε

If the stem containing vowel, followed by ED is replaced by ε (empty).

Example: Died \rightarrow Di null

~~Died~~ \rightarrow Di ϵ \Rightarrow Di

iii) (* V *) ING \rightarrow ε

If the stem ending with ING is replaced by ε.

Example: Dancing \rightarrow Dance ϵ null Working \rightarrow work

* Dancing \rightarrow Dance ϵ \Rightarrow Dance

Making \rightarrow Mak ϵ null \Rightarrow watching \rightarrow walk

* Making \rightarrow Make ϵ \Rightarrow Mak Searching \rightarrow search

(C)

136

i) AT → ATE

Ex: * Confatedconfated → confat_e^{null} [∴ from 1(b) condition
* V* ED → e]i.e. confated → confat → confate [∴ from 1(c) condition
AT → ATE]

ii) BL → BLF

Example : 1: Troubling → another condition i.e. 1(c)

for word Troubling → 1(b) condition goes.

i.e. Troubling → Troubl. [∴ 1(b) condition
* V* ing → e]

Remaining words: Troubl

Troubl → Trouble [∴ BL → BLF, 1(c) condition]Example : 2 "Bubbled"→ Bubbled find "BL" in word i.e. BubbledBubbled → Bubblee [∴ * V* ED → e]Bubbled → Bubbl → Bubble [∴ BL → BLF]

iii) $(\text{Ad}! (\text{*L}(\text{or}) \text{*S}(\text{or}) \text{*Z})) \rightarrow \text{single character.}$ (31)

The stem is ending with "double consonant" and that consonant is "not L(or) S(or) ~~Z~~?"

Example :-

① HOPPING

↓↓↓↓bblb↓
CVCCVCC

HOPPING → HOPPE [∴ *V*ing → e]

HOPP. → HOP ^{single consonant.}

→ single character

∴ Hopping → Hopp → Hop

② Panned

Panned → Pann [∴ *ed → e]

Remaining stem: Pann

↑↑↑↑
Pann → pan

- Ending with double consonant is replaced
- by single character.

Pann → Pan

∴ Panned → Pann → pan

③ swimming

$\rightarrow \text{swimming} \rightarrow \text{swimm} [\because * \text{VXing} \rightarrow \epsilon]$

Remaining stem 'swimm' :

Swimm \rightarrow swim

$\downarrow \downarrow \downarrow \downarrow \downarrow$
C C V C C
double
consonants

replaced by
single consonant

$\therefore \boxed{\text{swimming} \rightarrow \text{swimm} \rightarrow \text{swim}}$

④ falling

$\text{falling} \rightarrow \text{fall} [\because * \text{VXing} \rightarrow \epsilon]$

Remaining stem : fall :

$\text{Fall} \rightarrow \text{Fa} \cancel{\text{l}} \times \text{(condition not verified)}$

$\therefore \boxed{\text{falling} \rightarrow \text{fall} \rightarrow \text{fa} \cancel{\text{l}} \times}$

IV) $(m=1 \& *0) \rightarrow F$ stem ending with double consonants and final consonant is not VXY

The measure 'm' (Vc combination) is 1 and the stem is ending with VXY (CVC) combination is replaced by F.

Example : ~~file~~ m=2 (initial)

(183)

① ~~Re~~ filing → fil, E: *vating → e

fil → file [cvc sequence followed by E.
↓
CVC
↓
m=1]

∴ VC combination m=1, cvc followed by E)

Filing → fil → file

Condition Verified.

② fail → fail (not Verified)

cvc (VC combinations m=1)

m=1 but cvc sequence is not there

so condition failed.

v) (*V*)Y → I

Our aim is to take out "Y" from the Word
and the remaining stem containing at least, one vowel.

If the condition is satisfied "Y" is replaced
with "I".

Example : Happy → HappI

Step 2 :- Derivation Morphology - I

i) (m>0) MIONAL \rightarrow MIKE

Example :

Relational \rightarrow Relate.

∴ Relational
 $\downarrow\downarrow\downarrow\downarrow\downarrow\downarrow\downarrow$
 CVCVCVCV
 ① ② ③ ④
 $\therefore m=4$

m>0 i.e. 4>0]

ii) (m>0) IZATION \rightarrow IZE

Example :

Generalization \rightarrow Generalize

Generalization
 $\downarrow\downarrow\downarrow\downarrow\downarrow\downarrow$
 CV CV CV CV CV CV
 ① ② ③ ④ ⑤ ⑥

$\therefore m=6 \Rightarrow 6>0$ & IZATION \rightarrow IZE

iii) (m>0) BILITI \rightarrow BLE

Example : Sensibiliti \rightarrow sensible

$\downarrow\downarrow\downarrow\downarrow\downarrow\downarrow$
 CV CCV CV CV CV
 1 2 3 4

Give combinations = 4
 $m=4$

$\therefore m=4$ i.e. 4>0 BILITI \rightarrow BLE

22(rboop)

Step 3 :- Derivation morphology - II

(125)

i) $(m > 0) \text{ICATE} \rightarrow \text{IC}$

Ex 3: duplicate

Example: Tripligate \rightarrow Tripli
 $\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$
 $\text{CVCVCV} \quad \text{VCV}$
① ② ③

$m=3 \in \text{ICATE} \rightarrow \text{IC}$

i.e. $(3 > 0) \text{ICATE} \rightarrow \text{IC}$

ii) $(m > 0) \text{FUL} \rightarrow \epsilon$

Example: HOPEFUL \rightarrow HOPE

H O P E F U L
 $\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$
 $\text{CVCV} \quad \text{VC}$
① ② ③ $m=3$

$\therefore 3 > 0 \text{ FUL} \rightarrow \epsilon$

i.e. $\text{hopeful} \rightarrow \text{hope}$

iii) $(m > 0) \text{NESS} \rightarrow \epsilon$

Example: Goodness \rightarrow good
 $\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$
 CVCVCC

$m=2$ → VC combinations 2

$2 > 0 \text{ Ness} \rightarrow \epsilon \rightarrow \text{Goodness} \rightarrow \text{Good}$

Step 4

Derivation Morphology - III

(136)

i) ($m > 0$) ANCE $\rightarrow \epsilon$

Example : Allowance \rightarrow Allow

-Allowaⁿce
↓d↓↓↓ ↓↓↓
vccvc vccv
① ② ③

$m=3$

\therefore | ($3 > 0$) Allowance \rightarrow Allow.

ii) ($m > 0$) ENT $\rightarrow \epsilon$

Example : Dependent \rightarrow Depend

↓↓↓↓↓↓↓
cvccvccvcc
④ ⑤ ⑥

$m=3$

\therefore | ($3 > 0$) Dependent \rightarrow Depend

iii) ($m > 0$) IVE $\rightarrow \epsilon$

Example : Elective \rightarrow elect

Elective
↓↓↓↓↓↓↓
vc vccvcv
① ② ③

$m=3$

\therefore | ($3 > 0$) Elective \rightarrow elect

Step 6 $\rightarrow A \& B \rightarrow -$

(A)

i) $(m > 1) F \rightarrow \epsilon$ Bye \rightarrow By.

Example: Probate \rightarrow Probat

ii) $(m > 1 \& ! * 0) NECC \rightarrow \epsilon$

condition $m > 1$ and stem not ending with CVC combination.

Example: Goodness \rightarrow Good

$\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$
 CVC CVC
 $\frac{1}{2}$

M=2 Goodness $\xrightarrow{\text{Good}}$

(B) $((m > 1) \& * d \& * l)) \rightarrow$ Single letter

Example: controll \rightarrow control

$\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$
 CVCCCVCC
 $\frac{1}{2}$

Replaced by
single letter L

check the condition m (i.e m=2) and the stem

is ending with double consonant and 'l' is replaced
by single letter.

Example 2: Roll \rightarrow Roll

m=1, Condition not Verified.

The two other techniques in stemming are

- * i) Dictionary lookup stemmers
- ii) Successor stemmers
- iii) Conclusion

ii) Dictionary lookup stemmers:

An alternative way to determine a "stem" is a

"Dictionary lookup mechanism".

In this approach, simple stemming rules are may be applied.

The original term (or) stemmed version of the term is looked up in a dictionary and replaced by the stem that is the best representation of it.

This technique has been implemented in the "INQUERY and Retrieval ware systems".

The INQUERY system uses a stemming technique called "K-stem".

K-stem: K-stem is a morphological analyzer that conflates word variants to a root form.

Eg:- "Memorial" and "Memorize" reduced to "memory".

But memorial and memorize are not synonyms and they have different meanings.

K-Stem requires a word to be in the dictionary before it
reduces one form word form to "Other".

Some endings are removed, even if the root form is not found in the dictionary.

Ex: "need", "ly"

* K-stem uses the following six major "data files to control and limit the stemming process:

- ① Dictionary of Words (Lexicon).
- ② Supplemental list of words for the dictionary.
- ③ Exception list for those words that should retain an "e" at the end (Ex: "suite" → "suite" but "suited" → ~~"suit"~~)
- ④ Direct - Conflation: → allows definition of direct conflation via word pairs that override the stemming algorithm.

⑤ Country - Nationality:

It indicates conflations between nationalities and countries. ("British" → ^{maps} "Britain").

⑥ Proper Nouns:

A list of Proper Nouns that should not be stemmed.

Retrieval Index System:-

This system lies in the thesaurus / semantic network support the data structure that contains over 4,00,000 words.

iii) Successor stemmer :-

(140)

Successor stemmers are based on the "length of prefixes" that Optimally expand the stem. It is also known as "successor variety".

It is used to determine the "word", word is divided into "segments".

for Example:-

→ The Successor Variety for the first three letters (i.e word segment) of a five letter word is no. of words that have the same first three letters but fourth letter is different.

Graphical Representation:-

→ A Graphical representation of Successor Variety is shown

"Symbol tree".

→ The "symbol tree" for the tame

bag, barn, bring, both, box and bottle.

→ The successor Variety for any prefix of word is the no. of childrens that are associated with the node in the "symbol tree". representing that prefix.

for example:- The Successor Variety for the first letter "b" is "three".

The Successor Variety for the prefix "bg" is "two".

B	3
<u>Baa</u>	2
→ <u>Baq</u> (first successor varying 1 and and successor varying 2)	0
<u>Bar</u>	1
→ <u>Barn</u>	0
<u>Ba</u>	1
<u>Bai</u>	1
<u>Brin</u>	1
→ <u>Bring</u>	0
<u>Bo</u> (: Bo followed by x and Bo followed by t)	2
→ <u>Box</u> Bo followed by x	0
<u>Bot</u> Bo followed by t	2
<u>Both</u>	0
<u>Bott</u>	1
<u>Bottl</u>	1
→ <u>Bottle</u>	0

fig: Successor Vari

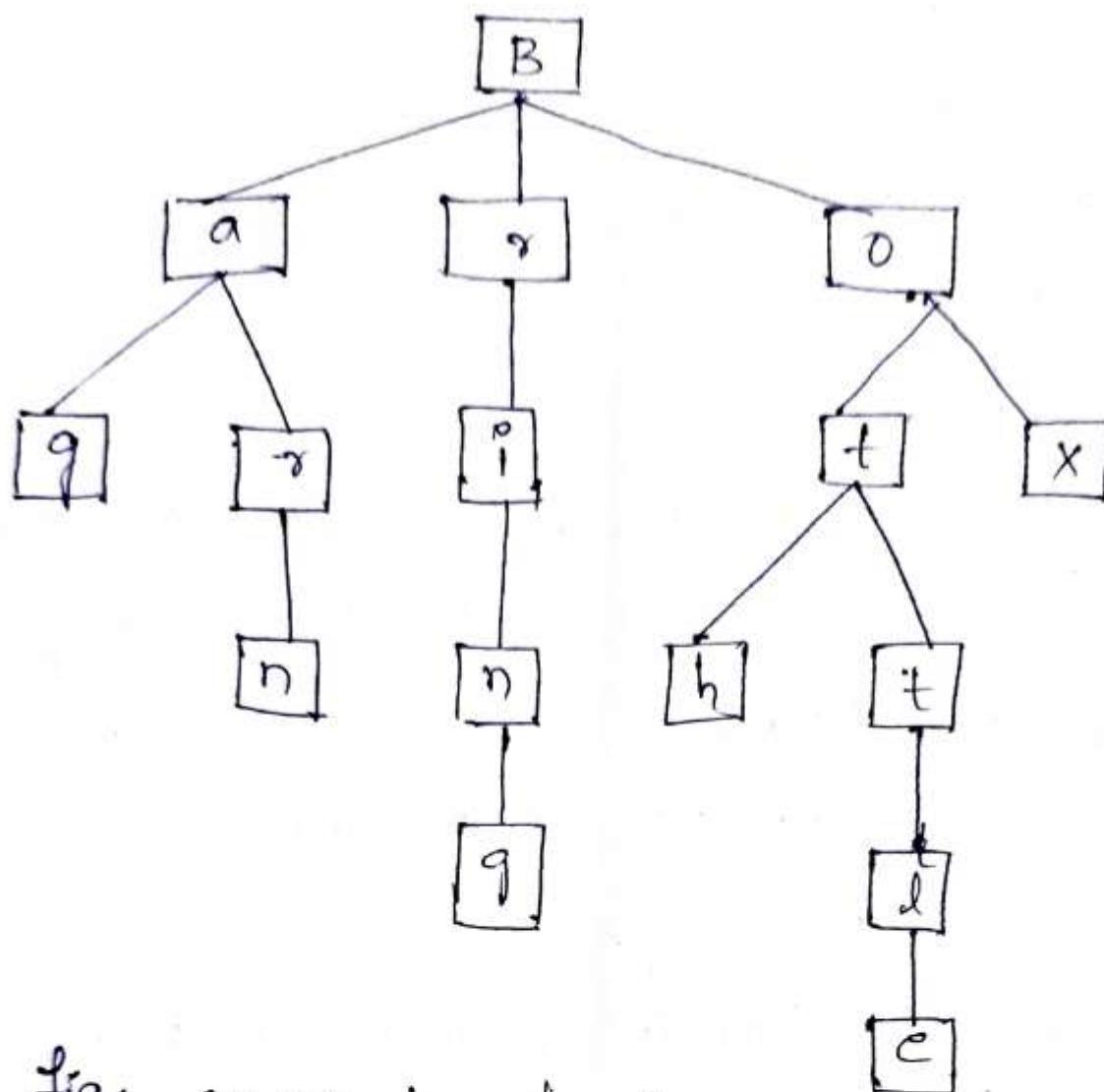


fig: Symbol tree for terms bag, barn, bring, box, bottle and both.

Successor Variety Methods :-

The Successor Varieties of a word are used to segment a word by applying at one of the following 4-methods:

1. Cutoff Method
2. Peak and plateau
3. Complete Word Method
4. Entropy Method.

① cutoff Method :-

cutoff value is selected to define "stem length" and the values are varies from word-to word/ each possible set of words.

② peak and plateau :-

"a segment break" is made after a character whose successor variety exceeds that character and immediately preceding it.

③ Complete Word Method :-

Break on boundaries of complete words.

④ Entropy Method :-

Used to Distribution of Successor Variety letters.

(Let $|D_k|$ be no. of words begin with "k" length sequence of letters "a".)

iv) Conclusions

frakes summarized the study of various stemming.

frakes comes to following conclusions:-

if stemming can effect retrieval (recall) and where effects were identified they were positive.

- iii) Stemming is as effective as manual collation.
- iii) Stemming is dependent upon the nature of the vocabulary.

Extra Information:

→ To quantify the impact of stemmers, "parce" has defined a stemming performance measure, called

"ERROR RATE RELATIVE TO TRUNCATION (ERRT)".

ERRT :- ERRT can be used to compare stemming algorithms.

→ This approach depends upon "concept groups".

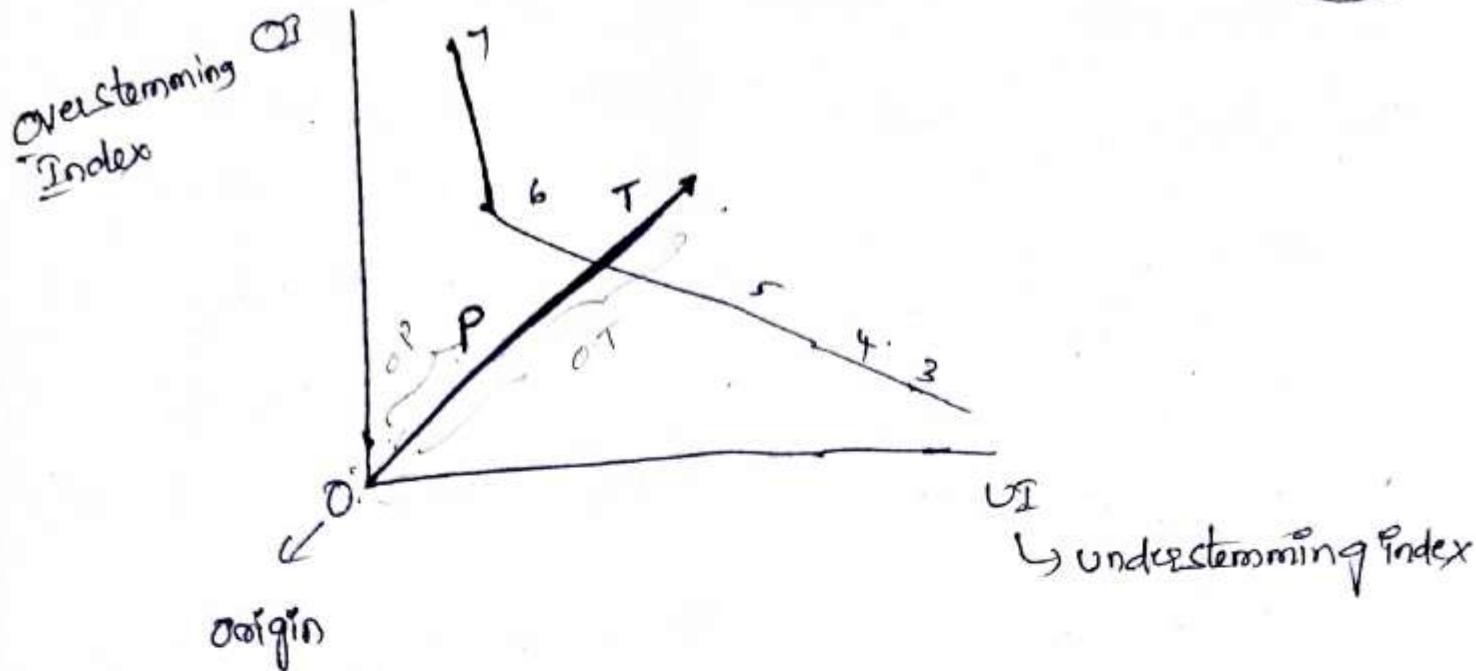
→ After applying a stemmer that is not perfect, "concept groups may still contain multiple stems rather than one".

→ This introduces an error reflected in the understemming index (UI).

This also possible that some stem is found in multiple groups, is ^{error state} reflected in the overstemming index (OI).

⇒ The UI & OI values can be reflected calculated based on truncated word lengths.

The Perfect case is where $UI = OI = 0$ (zero).



⇒ FRRT is calculated as the distance from the Origin to the (UI, OI) coordinate of the stemmer being evaluated (O_P) versus distance from Origin to the pure truncation (OT)