

Unit-II

Cataloguing and Indexing

History and objectives of indexing:-

- Indexing: The transformation from received items to searchable data structure is called indexing.
- Indexing (originally called cataloguing) is the oldest technique for identifying the contents of items to assist in their retrieval.
- The objective of cataloguing is to give access points to a collection that are expected and most useful to the users information.
- The basic information required on item is, what is the item and what it is about, has not changed over the centuries.

As early as the Hyman-89:-

* Third-millennium (in Babylon):-

The Libraries of cuneiform tablets were arranged by Subject. (Hyman-89).

* Cuneiform tablets indicates the logo-symbolic script, used in Ancient.

* Upto 19th century there was little advancement in Cataloguing.

i.e. only changes in the methods used to represent the basic information. [Norris-69].

* In the Late 1800's :- Subject Indexing became hierarchical.

exit Dewey Decimal System (DDS)

* In 1963 :-

The Library of Congress initiated a study on the computerization of bibliographic surrogates.

* From 1966 - 1968 :-

The Library of Congress ran its "MARC I Pilot project".
MARC → (Machine Readable Cataloging).

Here,

MARC standardizes

* the structure

* Contents &

* Coding of bibliographic records.

* In 1969 :-

The system became Operational. [Avramis]

* In 1965 :- (DIALOG)

The earliest commercial cataloguing System is DIALOG, which was developed by "Lockhead Corporation" for NASA.

* In 1978 :-

It became (i.e., DIALOG) commercial with three government files of indexes to technical publications.

* In 1988 :-

→ DIALOG sold to KPA Knight-Ridder, and it contained over 320 index databases used by over 91000 subscribers in 86 countries. [Harper-8].

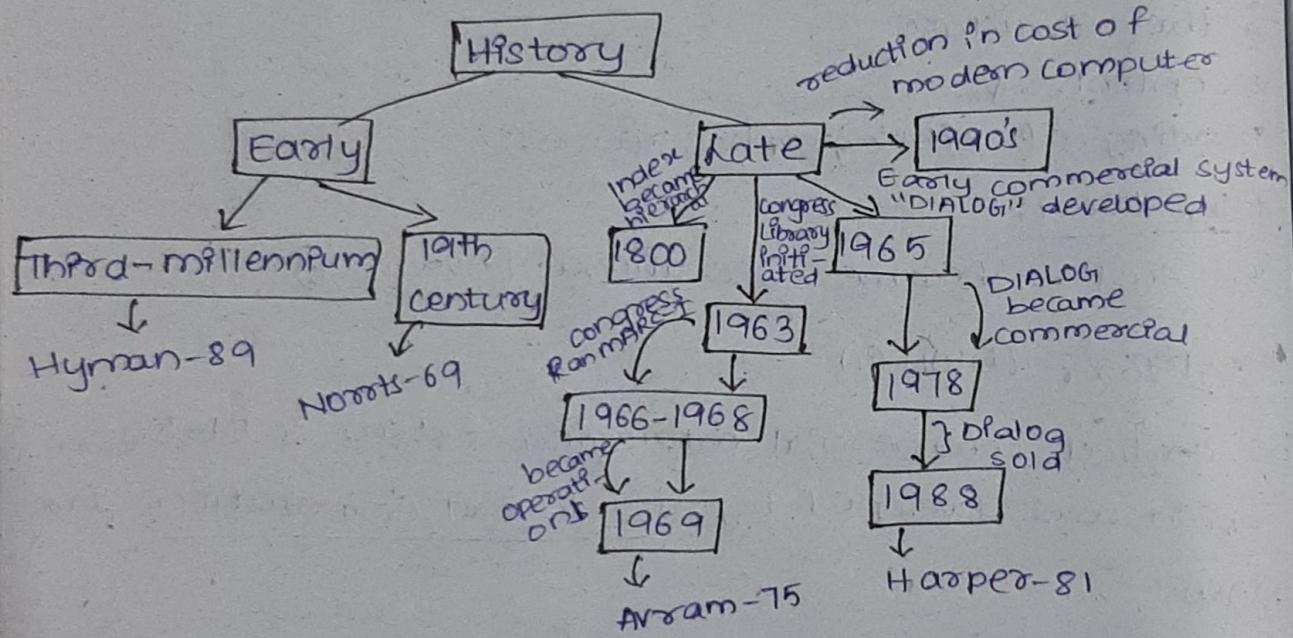
→ Indexing (Cataloguing) until recent:-

The indexing was accomplished by creating a bibliographic citation in a structured file that references the original text.

- Citation means source information.
- These files contains citation information, "about item", "Keywording", the subjects of "item" and "length of free text field used for abstract/summary".

* In 1990's :-

A significant reduction in cost of processing power and memory in modern computer.



→ The indexing process is typically performed by professional indexers associated with library organizations.

Objectives:

The objectives of indexing have changed with the evolution of IRS.

→ Availability of the full text of the items in searchable form alters the objectives historically used in determining guidelines for manual indexing.

→ In manual indexing environment, The use of controlled vocabulary makes the indexing "process slower", but it potentially simplifies the "search process".

- The availability of items in "electronic forms" changes the objectives of manual indexing.
 - The words used in an item do not always reflect the value of the concepts being presented.
- Indexes < public file index
 private file index
- The public file indexer needs to consider the information needs of all users of library system.
 - Items overlap b/w full items indexing, public & private indexing of files.
 - Users may use "public index files" as part of search criteria to increase recall.
 - They users can constrain these search by private index files.
 - The primary objective of representing the concepts within an item to facilitate user's finding relevant information.
 - The other objectives of indexing :-
ranking, item clustering.

* Indexing process :-

- Indexing is an important process in Information Retrieval (IR) systems.
- The indexing process is typically performed by professional indexers associated with library organizations.
- The indexing process depends on 3 thing/files, i.e.,
 - Document file
 - Public Index file
 - Private Index file.

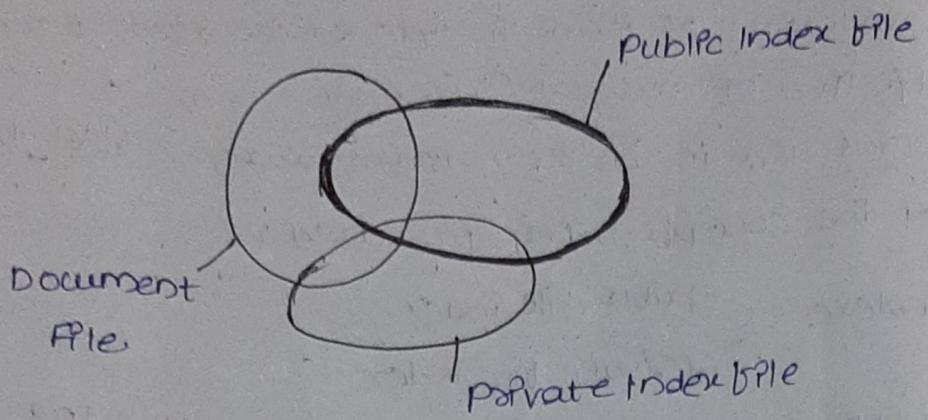


fig:- Indexing process.

- ⇒ The above diagram showing the items overlap blur full item & indexing, public file indexing and private file indexing.
 - ⇒ Linking index terms is needed when there are multiple independent concepts found with in an item.
- Document file:-
- (The full text Searchable data structure for items in)
- "The Document file provides a new class of indexing called total document indexing".

Public Index files:-

Users may use public index files as part of their search criteria to increase the recall.

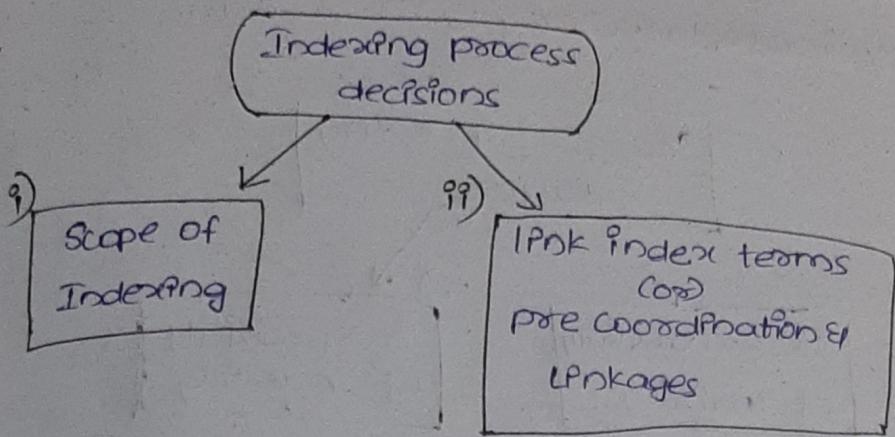
Private Index files:-

Users may want to constraint the search by their private index file to increase the precision of search.

Precision.

- ⇒ Indexing process takes the 2 types of decisions.

i.e

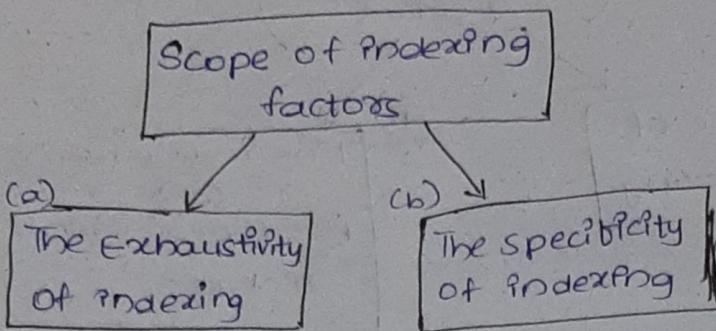


Manual indexing:- The process of reliably & consistently determining the bibliographic terms that represents the concepts in an item is extremely difficult.

i) Scope of Indexing

- The first decision is the scope of the indexing to define what level of detail the subject index will contain.
- This is based upon usage scenarios of the end users.
- When indexing is performed manually, the problems arise from two sources:
 - i.e "the author" and "the indexer".
- The vocabulary domain of the author may be different than that of indexer.
- So the results is in different quality levels of indexing.
- The indexer is not an expert on all areas & the indexer must determine when to stop the indexing process.
- There are two factors involved in deciding on what

level to index the concepts in an item.



(a)

The Exhaustivity:-

Exhaustivity of indexing is the extent to which different concepts in the item are indexed.

for ex:-

If two sentences of a 10-page item on microprocessors discusses on-board caches. Should this concept be indexed.

(b) The specificity:-

Specificity relates to the preciseness of the index terms used in indexing.

for ex:-

whether the term "processor" (as)

the term "microcomputer" (as)

the term "pentium"

Should be used in the index of an item is based upon the specificity decision.

⇒ High Exhaustivity and Specificity indexes almost every concept in the item using as many detailed terms as needed.

⇒ "Low Exhaustivity" has an adverse effect on both precision & recall.

- "Low specificity" has an adverse effect on precision but no effect to a potential increase in recall.
- Another decision on indexing Ps, what portion of an item should be indexed.
- The simplest case to limit the indexing Ps to
 - * The title(oz)
 - * The title & Abstract zones
- This case leads to loss of both precision & recall.

i) Pre-coordination and linkages (oz) Link index terms:

The another decision is the need to link index terms together in a "single index" for a particular concept. Linkages are used to correlate related attributes associated with concepts discussed in an item.

→ The process of creating "terms linkages" at index creation time is called precoordination.

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

→ Post coordination is implemented by "AND"ing index terms.

i.e. which only finds indexes that have all of search terms.

Forexit

The capability to link the source of a problem, the problem & who is affected by the problem may be desired.

*Linkage of Index Terms:

INDEX TERMS

Oil, wells, mexico, CITGO, refineries,
Peru, BP, drilling

Methodology

No linking of terms.

(Oil, wells, mexico, drilling, CITGO)

Linked (pre coordination)

(U.S., Oil refineries, Peru, introduction)

(CITGO, drill, oil wells, mexico)

Linked (pre coordination)

(U.S., introduction, oil refineries,
Peru)

With position indicating
role.

(SUBJECT : CITGO;

ACTION : drilling;

OBJECT : oil, wells;

MODIFIER : in mexico)

SUBJECT : U.S.;

ACTION : introduces

OBJECT : oil refineries;

MODIFIER : in Peru)

Fig: Linkage of index terms.

The above fig shows the different types of linkages.
 → It assumes that an item discusses the drilling of oil wells in 'mexico' by CITGO & the introduction of oil refineries in 'Peru' by the U.S.

→ When the linked capability is added, the system does not erroneously relate Peru & Mexico since they are not in the same set of linked items.

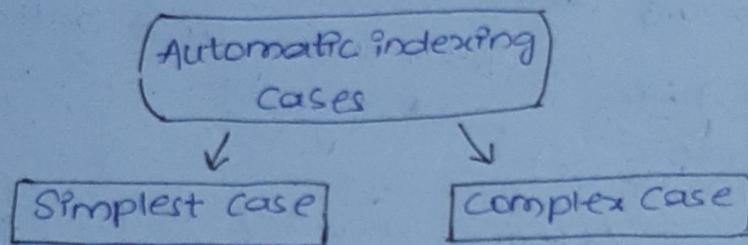
⇒ Positional role?

The U.S was introducing OPI defines in Peru, Bolivia and Argentina then the positional role technique would require three entries, only the difference is affected country position.

Here all 3 countries would be listed with three "MODIFIERS".

* Automatic Indexing

→ Automatic indexing is capability to automatically determine the index terms to be assigned to an item.



Simplest Case:

It deals with "total document indexing".
(or)

All words in the document are used as possible index terms.

Complex Case:

In this process the objective is emulate a "human indexer" and determine a limited no. of "index terms" for major concepts in the item.

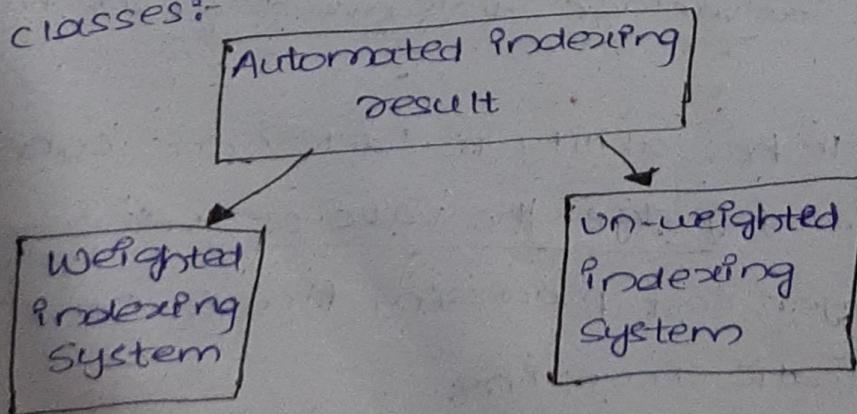
Here we have advantage & disadvantage of human indexing.

Advantage: It has ability to "determine concept abstract and judge the value of a concept."

Disadvantage: Cost, processing time & consistency over the automatic indexing.

Human Indexing Vs Automatic Indexing:-

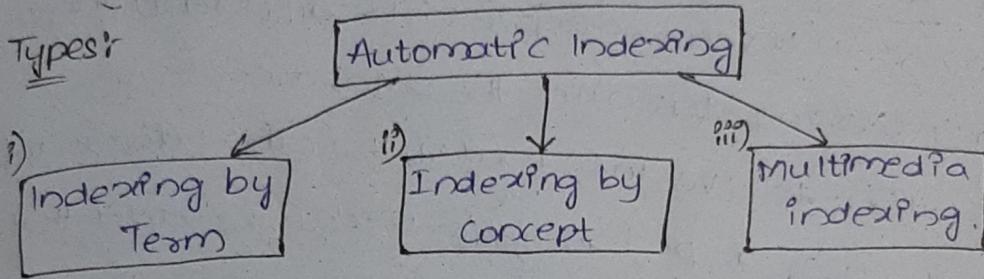
- The processing time of an item by a human indexer varies significantly based upon the "Indexer's knowledge" of the concepts being indexed with the exhaustivity & specificity guidelines & the amount & accuracy of processing via Automatic File Build.
- Usually it takes at least "five minutes" per item.
i.e. Human iteration with computer.
e.g.: typing speeds, cursor positioning, correcting spelling errors, taking breaks b/w activities.
for short items it takes 5min time pattern.
e.g.: (300-500) words)
- Automatic indexing requires only a "few seconds" less of computer time based upon the size of the processor and the complexity of the algorithms to generate the index.
- If indexing is being performed automatically by an algorithm, there is consistency in the index term selection process.
- Human indexers typically generate different indexing for the same document.
- ⇒ Indexes resulting from automated indexing fall into two classes:-



Weighted Indexing System: Weight is based on the function associated with the frequency of occurrence of item. In a weighted indexing system an item is made to place or value on the index terms that represents concept in the document.

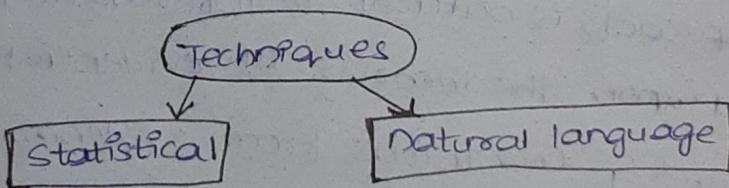
Unweighted indexing System: In an un-weighted indexing system, the existence of index term in a document & some times its "word & location(s)" are kept as part of the searchable data structure.

Types:



i) Indexing by Term:

The terms of the "original item" are used as a basis of the index process. These are 2 major techniques used for creation of the index.



* Statistical Techniques:

→ These techniques can be based upon "vectors" and "probabilistic models" with a special case being "Bayesian models".

→ They are classified as "statistical" because their calculation of weights use statistical information such as frequency of "occurrence of words" & their distributions in the "searchable DS"?

⇒ Often weighted systems are discussed as vectorized information systems.

This association comes from the "SMART Systems" at Cornell University created by Dr. Gerald Salton. (i.e. Salton → 73, Salton → 83).

→ The system emphasizes weights as a foundation for "information detection" & stores these weights in a "vector form".

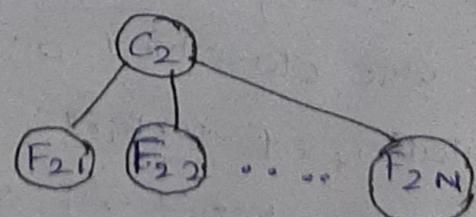
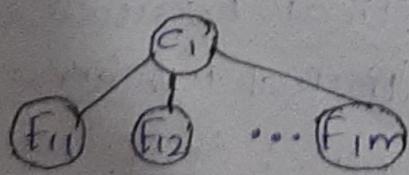
Each "vector" represents a document and each "position" in a vector represents different unique word (processing token) in the database.

→ In addition to a vector model, the other dominant approach uses a probabilistic model. The model that has been most successful in this area is the Bayesian Approach.

Bayesian Approach

This approach is natural to information System & is based upon the theories of evidential reasoning (i.e. drawing conclusions from evidence).

A Bayesian net is a directed acyclic graph in which each "node" represents a random variable and the "arc(s)" b/w "nodes" represents a probabilistic dependence b/w the node and its parents.



HERE

Nodes C_1, C_2 represents

"the item contains concept C_i ".

If nodes represent "The item has feature F_j "
eg: words.

The goal is to calculate the probability of ' C_i ' given F_j . To perform that calculation two sets of probabilities are needed:

1. prior probability $P(C_i)$ \Rightarrow that item relevant to concept C_i .

2. conditional probability $P(F_j | C_i)$ \Rightarrow that the features F_j where $j=1, n$ are in an item contains topic C_i .

CP.

* Natural language processing / Techniques:

\rightarrow It process the items at the

* morphological,

* Lexical,

* Semantic,

* Syntactic E1.

* discourse levels.

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

\rightarrow The discourse level is abstracting information beyond the sentence.

ii) Indexing by Concept:

The basis for concept indexing is that, there are many ways to "express the same idea" E1 increase the retrieval performance comes from single representation.

→ Hence indexing by term treats each of these occurrences as a "different index" & then uses the same (or) other query expansion techniques to expand the query.

→ Concept indexing determines a canonical set of concepts based upon a test set of terms & uses them as basis for indexing all items.

iii) Multimedia Indexing:

Indexing associated with multimedia differs from the previous discussions of indexing.

Indexing video (or) images can be accomplished at the raw data level, (e.g.: aggregation of raw pixels) & the feature level distinguish primitive attributes such as "color" & "luminance" & at the semantic level meaningful objects are recognized.

Ex:- An example is "processing of video".

→ The system will periodically collect a frame of video I/P for processing.

It might compare that present frame to the last frame captured, to determine the difference b/w the frames.

→ If the difference is below threshold it will discard the frame.

⇒ processing of image:

To process an image, semantic level indexing requires for the pattern recognition of objects with in the images.

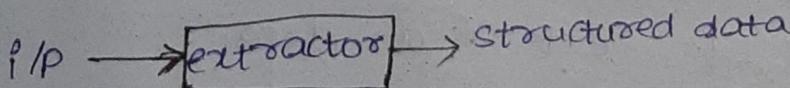
exit photobook (pentland - 94)
IBM's Qbic.

Audio:

- ⇒ If you consider an analog audio I/P & the system will convert the audio to digital format.
- ⇒ In multimedia indexing two mechanisms are used
 - "positional" and "temporal".
- ⇒ The first mechanism is used to represent "a linear sequential composition".
- ⇒ The second mechanism is based upon 'time', because the modalities are executing concurrently.
modalities means mode in which something exists.
- ⇒ The typical video source off television is inherently a multimedia source.
It contains video, audio & potentially closed captioning.

Information Extraction

IE is the task of automatically extracting structured information from unstructured (or) semi-structured documents & other sources and store them in a database.



*Information Extraction?

There are two processes associated with the information extraction.

i.e., 1. determination of facts to go into structure fields in a database.

2. Extraction of text, that can be used to summarize an item.

⇒ In first case, only a subset of the important facts in an item may be "identified and extracted" whereas in summarization all of 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 (AFB). Its goal is to process incoming fact items & extract "index terms" that will go into a structural database.

⇒ An IRS goal is to provide an indepth representation of the total content of an item.

⇒ An Information Extraction system, only analyzes those portions of a document that potentially contain information relevant to the extraction criteria.

⇒ The term "slot" is used to define a particular category of information to be extracted.

Here

→ Slots are organized into "templates" (or) "Semantic frames".

→ Information Extraction requires multiple levels of analysis of text of an item.

→ It must understand the words & their context.

→ The processing is very similar to the Natural language processing ; that described under indexing.

Metric to Compose Information Extraction

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

Recall:-

Recall Refers to how much information was extracted from an item versus how much should have been extracted from the item.

precision:-

precision refers to how information was extracted accurately "versus" total information extracted.

Additional metrics:-

The additional metrics use are over generation & fallout.

Over Generation:-

Over Generation measures the amount of irrelevant information that is extracted.

Fallout:-

Fallout measures how much a system assigns "incorrect slot fillers" as no. of potential incorrect slot fillers increases.

⇒ Another related information technology is document summarization.

The goal of document summarization is to extract a summary of an item maintaining the most important ideas that significantly reducing the size.

⇒ The term "slot" is

Data Structure

* Introduction to Data structures

Data structures

- A data structure is a specialized format for organizing, processing, retrieving & storing data.
- There are usually "two" major data structure in any information system.

* One data structure stores & manages the received items in their Normalized form.

This process is called the "document manager".

* The other major data structure contains the processing tokens & associated data to support search.

This process is called the "document search manager".

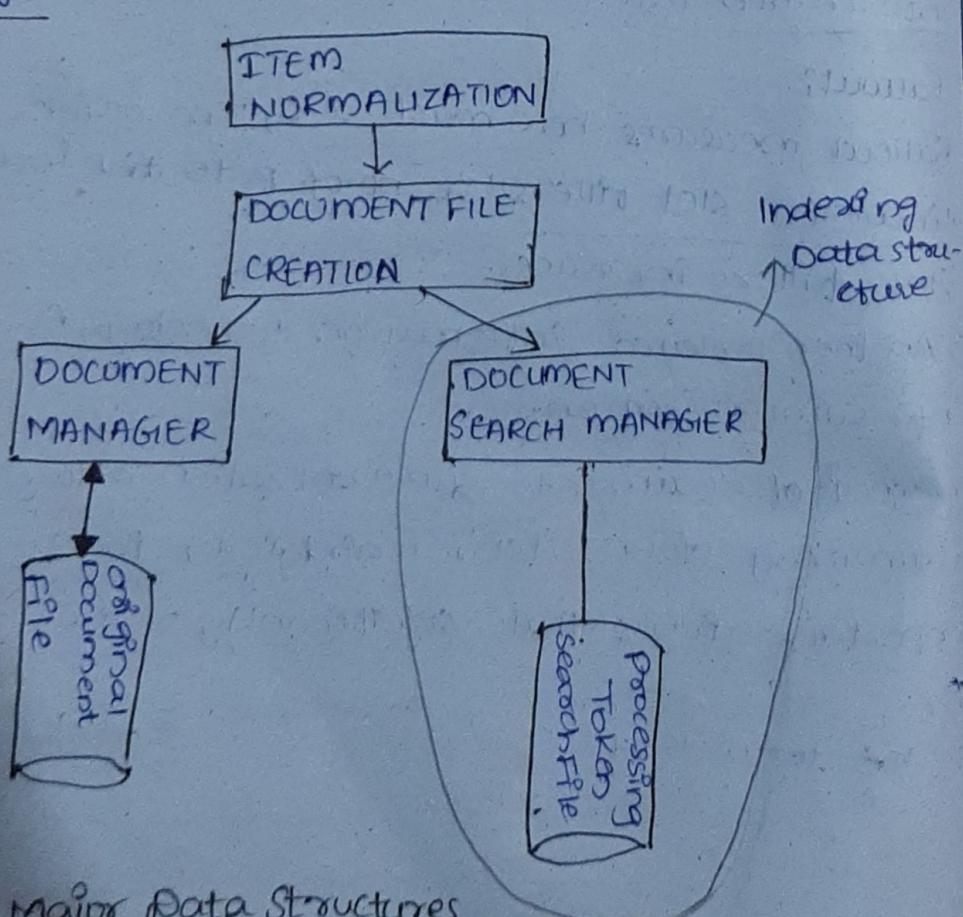


Fig:- Major Data Structures

→ The above big expand the Document file creation.
i.e document manager &
document search manager.

The result of search are references to the items that
satisfy the Search statement which are passed to
the document manager for retrieval.

→ The most common data structure encounters in
both "database" and information systems is the
"Inverted file system". → (4th chapter)

Inverted File System:

It minimizes secondary storage access when multi-
ple search terms are applied across the total
database.

→ A various of the Searchable data structure is the
"N-gram structure".

The N-gram structure that breaks processing
tokens into smaller string units & uses the token
fragments for search.

→ A "special data structure" that is becoming common
place because of its use on the internet is

"Hypertext"

Various Data structures are:-

- * Stemming algorithms
- * Inverted structures
- * N-Gram data structures
- * PAT data structures.

*Signature data structures

*Hyper text data structures

Stemming Algorithms:-

Stemming means cutting and trimming.

i.e. un-necessary things will be trim/stem.

→ The concept of stemming is introduced in the 1960's.

→ The main goal of stemming is to improve performance as it require less system resources by reducing the no. of unique words that the system contain.

→ Stemming algorithms are used to improve the efficiency of the information system & to improve recall.

Stemming:-

→ Stemming is the process of producing morphological variants of a root/base word.

→ Stemming programs are commonly referred as stemming algorithms (or) stemmers.

i.e. reduce the word forms to proper class.

→ The concept of stemming has been applied to the information systems from 1960's.

→ The original goal of stemming was to improve performance by reduce the no. of resources that system contain.

Introduction to the stemming process

Stemming algorithms are used to improve the efficiency of the information system. The another major use of stemming is to improve the recall.

→ Stemming is a technique used to extract the base form of the words by removing affixes from them.

It is just like cutting down the branches of a tree to its stems.

for ex: ① The stem of words

eating, eats, eaten is eat

② The stem "Comput" could associated with

Computable, Computability, computation, Computational, computed, computing, Computer, Computerise, Computerize to one compressed word

i.e. Comput

③ Stemming / Stem word "Calculate"

Calculate, calculates, calculation, calculations, calculating.

→ The most common stemming algorithm removes suffixes & prefixes sometimes recursively to derive the final stem.

→ The other Techniques such as

"Table / Dictionary look-up" in stemmers and

"Successor Stemming" / nstemmers provide attending

Alternatives to the stemming algorithms

Porter Stemming Algorithm:-

It is one of the most popular stemming method proposed in 1980.

It is based upon set of conditions of the stem, suffix and prefix & associated actions given the condition.

Some examples of stem conditions are:

1) The measure 'm' of a stem is a function of sequences of vowels [a,e,i,o,u,y] followed by a consonant.

e.g. $\underbrace{a,e,i,o,u}_{v} \underbrace{y}_{c}$

If 'v' is a sequence of vowels and 'c' is a sequence of consonants.

Then 'm' is : $\boxed{c(vc)^m v}$

Where, the initial 'c' and final 'v' are optional 'm' is the no. of 'vc' repeats.

Measure

$$m=0$$

$$m=1$$

$$m=2$$

Example

bree

(This word doesn't contain any combination of vc).

Fees

TROUBLE

compute

$\begin{matrix} \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ c & v & c & c & v \end{matrix}$

2) *(X)

Stem ends with letter X.

3) *V*

Stem contains a vowel.

∴ Based on 'm' value the Porter Stemmer will work.

4) *

Stem ends in "double Consonant"

5) *

Stem ends with "Consonant-vowel-Consonant" sequence
(CVC)

⇒ 2nd 'C' Should not be w,x,y.

Where the final consonant is not

w,x,(or) y

ex:- WIL, HOP

CVC CVC

Steps in Porter stemmer

The porter stemmer; step 1(a)

Conditions:-

① SSES → SS

② IES → I

③ SS → SS

④ S → E

① SSES → SS:-

If we have a word / term that is ending with "SSES"
is simply replaced by "SS".

Ex:- caresses

i.e. caresses after applying the porter stemmer the word became caress

i.e. caresses → caress

② IES → I :-

If we have a word/term that is ending with "IES", it is replaced by "I".

ex:- Pon_{ies} → Pon_i
 |
 ties → ti
 |
 I

③ SS → SS :-

If we have a word/terms that is ending with "SS" it remains same i.e., "SS".

ex:- Care_{ss} → Care_{ss}
 |
 ss

④ S → E :- (i.e. null);

If the stem is ending with single 's' then it will be replaced by null 'E'. ~~with~~

ex:- cat_s → cat_E
 |
 E

Step 1(b)

Condition's :-

① (m > 1) EED → EE

② (*v*)ED → E

③ (*v*)ING → E

Step 1(b) (stemming with conditions)

① (m > 1) EED → EE (condition 1)

If the measurement of stem is > 1 (i.e. combination of 'vc' is appearing > 1) & it is ending with "EED" will be replaced by "EE".

ex:- "Agreed"

Check 'vc' combination i.e., Agreed

$\therefore m=2$ & our condition is $m>1$

$2>1$

i.e. $2>1$ is satisfied.

\therefore EED is replaced by EE.

i.e. Agreed \rightarrow Agree

\therefore condn verified : Agreed \rightarrow Agree.

ex:- "Feed"

check 'vc' combination i.e., Feed

$\therefore m=1$

our condition is $m>1$ but we have only '1' VC combination

trn $\boxed{1>1}$ x not satisfied.

\therefore So the word is remains same.

i.e. Feed \rightarrow Feed

\therefore condition not verified : Feed \rightarrow Feed.

② $(\ast v \ast) ED \rightarrow e$ (condition 2)

If our stem containing vowel followed by "ED" then
it will be replaced by 'e'.

ex:- "plastered"

The above stem containing vowels(v) (i.e. plastered)
'a' and 'e') followed by "ED" is replaced by null, so
the word will become "plaster".

i.e. plastered \rightarrow plaster

condition verified : plastered \rightarrow plaster.

ex:- "bled"

In the above word check 'vowels', but in this
word non-of them is vowel so it is remain

Same.

i.e. bled there is no vowel

i.e. bled → bled

Condition not verified : bled → bled.

③ (+v*)ING → e (Condition 3):-

If our stem containing vowel followed by "ING" then it will be replaced by 'e'.

e.g. "Motoring"

→ Find out "ING" is given word Motoring after finding out "ING" the remaining stem contains vowel 'o'.

i.e. Motoring. so the "ING" is replaced by null.

i.e. Motoring → Motor

Condition verified : Motoring → Motor.

e.g. SING

→ Find out "ING" in given word i.e. SING and the remaining stem contains only 's' and it is "not vowel" i.e. SING.

This rule will not be applied here and the word will be remain same.

i.e. SING → SING

Condition not verified : SING → SING

Step 1(c) (clean up)

Conditions

① AT → ATE

② BL → BLE

③ (*d & !(*L & *S or *Z)) → single letter.

⑨ $(m=1 \& * \times 0) \rightarrow E$

These rules are ran if second (or) third conditions
in "1b" apply.

① AT \rightarrow ATE

Here our word ps contlat \rightarrow contlate
check the word "contlated" with 2nd condition in
1b i.e. $(*v*)ED \rightarrow E$
so the word became contlat(ed) \rightarrow contlat
the remaining stem contlat contains 'v' vowel and
"ED" is replaced by 'E' i.e. contlated \Rightarrow contlat
then convert the word "contlat" into proper word,

i.e. contlate

$\therefore \underline{\text{contlat}} \rightarrow \underline{\text{contlate}}$ [$\because \text{AT} \rightarrow \text{ATE}$]

② BL \rightarrow BLE :-

or Troubling \rightarrow TroublE

Here take a word ps "Troubling".

i.e. E Troubling \rightarrow TroublE

check the word "Troubling" with "3rd" condition in

1b.

i.e. $(*v*)\text{ING} \rightarrow E$

so the word became Troubling \rightarrow TroublE

findout the 'ING' in given word and then check
remain stem i.e. "Troubl" contain vowel(v) and then
replace "ING" by "E"

i.e. Troubling \rightarrow TroublE

then convert "Troubl" into proper word

i.e. Trouble.

$\therefore \underline{\text{Troubl}} \rightarrow \underline{\text{Trouble}}$ [$\therefore \text{BL} \rightarrow \text{BLE}$]

③ (*d & !(*L or *S or *Z)) \rightarrow single letter.
The stem is ending with double consonant (*d) and that consonant is not *L (or) S (or) Z then we convert it into single letter.

For Example (i) : "Hopping."

Here the word hopping contains "ING" so verify this hopping with 3rd condition step 1b (step)

Hopping

\rightarrow Find out "ing" and the remaining word 'HOPP' then check "v" vowel in remaining stem.

i.e. $\begin{array}{c} \text{Y} \\ \text{I} \\ \text{E} \end{array}$ HOPP so in hopp vowel is there and "ING" will be replaced by E.

HOPP^{ing} \rightarrow HOPP

Here the word "Hopp" stem is ending with double consonants $\begin{array}{c} \text{C} \\ \text{V} \\ \text{C} \end{array}$ and that consonant is not L (or) S (or) Z then we convert it into single letter.

i.e. $\text{PP} \rightarrow \text{P}$

HOPP \rightarrow single letter
 \downarrow
(HOP)

condition verified : HOPP^{ing}) \rightarrow hop

ex:- tanned

i.e. check V & C \rightarrow in above word tanned Here the stem is ending with "Ed" & i check the remaining stem contains "v" vowel followed by "ED". So it will be replaced by null "e".

tanned → tann
E

Here the stem is ending with double consonants i.e "nn" and then convert double consonant into single letter /consonant.

i.e tann → tan
 ↑

∴ Condition verified : tanned → tan.

ex² falling → ball

fall it is ending with double consonant i.e 'll'. But our condition is not "L" (or) "S" (or) "Z" so it will not converted into single letter, so the word is remain same.

∴ Condition not verified : fall(ing) → ball.

④ ($m=1$ & *0) → E

Here the measure 'm' is 1 (i.e vc combination) and the stem is ending with CVC combination is replaced by single letter.

i.e. CVC → E
 ↓↓ [not w^wly]

ex² Filing → file

Fil^wng
CVC m=1

Fil^wng → File

the word is ending with "ing" so the remaining stem is in CVC combination is replaced by 'E'.

i.e fil is in CVC combination.

i.e Fil → File.

∴ Condition Verified : bplng → file.

ex:- Fail → Fail.

Fail
[CVC]

The word 'Fail' is not ending with 'CVC' combination. So the condition is not validated and the word will be remain same.

∴ Condition not verified : fail → fail.

Step 1c and 2c

The poster stemmer : Step 1c and 2

Step 3c :-

Condition:

Y Elimination

$(\ast V \ast) Y \rightarrow I$

Our aim is to take out 'Y' from word and the remaining stem contain at least one vowel if this condition is satisfies then 'y' is replaced by 'I'.

For ex:-

Happy → Happi

Happy

→ find out 'Y' in given word
→ check is vowel is present in remaining word or not if 'v' is there then it will be replaced by 'i'.

f i f f f
Happy

so 'y' is replaced by 'i'.

∴ Condition verified : Happy → Happi

ex:- SKY → sky

SKY
S C U I

If we take out 'y' from stem the remaining stem is 'SK' which is not containing vowel. So the rule will not be applied here and the word is remain same. ∴ Condition not verified. SKY → SKY.

Step 2: Derivational Morphology - I.

conditions:

* $(m > 0)$ ATIONAL → ATE

ex:-

Relational → Relate

first check the measure 'm' from given word.

How many VC combinations are there. 'm'.

Relational
↓↓↓↓
V C V C V C V C m=4

our condition is $(m > 0)$

Here $m=4$ i.e. $4 > 0$ condition satisfied the "ATIONAL" is replaced by "ATE"

* $(m > 0)$ IZATION → IZE

ex:- Generalization → Generalize

check 'm' value and verify condition then after replace "IZATION" with "IZE".

i.e. Generalization
↓↓↓↓↓↓↓↓↓↓

∴ $m=5$ our condition is $m > 0$
 $5 > 0$ satisfied.

* $(m > 0)$ BILITI → BLE

ex:- Sensib^{pli}t^o → sensible

Check 'm' value and verify condition then after replace
"BILITI" with "BLE"

i.e. Sensibliti^p
 $\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$
C V C C V C V C V

$\therefore m=4$ our condition is $m>0$

i.e. $4>0$ satisfied.

The poster stemmer

= =

Step 3 and 4^p
= =

Step 3 :- Derivational Morphology - II

Conditions :-

(*) ($m>0$) ICATE \rightarrow IC

ex:- triplicate \rightarrow triplc

i.e. check 'm' value and if it is satisfies condition

then "ICATE" is replaced by "IC".

triplicat
 $\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$
C V C C V C V C

$m=3$ our condition $m>0$ i.e. $3>0$ satisfies.

(*) ($m>0$) FUL \rightarrow E

ex:- hopeful \rightarrow hope
 $\downarrow \downarrow$ E

check 'm' value and verify it with condition

hopeful
 $\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$ m=3 >0 ✓ satisfied

FUL is replaced by null i.e. E

∴ hopeful \rightarrow hope
 $\downarrow \downarrow$ E

⊗ (m>0) NESS → E

ex: goodness → good
 ↓
 E

i.e. check 'm' Condition for given word after verify the condition 'ness' is replaced by 'E' null.

goodness
↓↓↓↓↓↓
C V V C C V C C
m=2 >0.

Step 4:- Derivational Morphology - III

Conditions :-

⊗ (m>0) ANCE → E

ex: allowan_{ce} → allow
 ↑↑↑↑↑↑
 E
m=2 >0 ✓

⊗ (m>0) ENT → E

ex: dependen_t → depend
 ↑↑↑↑↑↑↑↑↑↑
 E
m=3 >0 ✓

⊗ (m>0) IVE → E

ex: effecti_{ve} → effect
 ↑↑↑↑↑↑↑↑
 E
m=3 >0 ✓

The porter Stemmer:-

Step 5) Clean up :-

Step 5(a):-

Conditions :-

⊗ (m>1) E → E

⊗ (m=1 & ! * 0) NESS → E

(m>1) E → E

ex:- $\begin{matrix} f & X & f & f & X & f & v \\ \text{P} & \text{o} & \text{r} & \text{b} & \text{a} & \text{t} & e \end{matrix}$ → porbat
↳ e

$m=2 > 1$ then replace E by Null'

($m=1$ & $\text{stem is not ending with "CVC"}$) NESS → E

ex:- $\begin{matrix} f & i & i & f & c & v & c & c \\ \text{G} & \text{o} & \text{o} & \text{d} & \text{n} & \text{e} & \text{s} & \text{s} \end{matrix}$ → Good
↳ e

Here the stem is not ending with CVC so it replaced by E. "NESS" replaced by "E".

Step 5(b) :-

Condition

④ ($m > 1$ & *d & *l) → single letter

Check the condition of ' m ' and the.

The stem is ending with 'd' i.e. double consonant and 'l' is replaced by single letter.

ex(1) :- controll → Control

$\begin{matrix} f & v & c & c & c & v & f \\ \text{C} & \text{o} & \text{n} & \text{t} & \text{r} & \text{o} & \text{l} \end{matrix}$
Controll

$m=2 > 1$ ∴ Condition verified.

The stem is ending with double consonants i.e. 'll' is replaced by single letter, 'L'. Condition verified.

Control \rightarrow Control
↳ single L

ex(2) :- $\begin{matrix} f & v & c & c \\ \text{R} & \text{o} & \text{l} & \text{l} \end{matrix}$ → Roll

Check the condition of ' m ' in this only one combination of 'vc' is there and it is not satisfies the condition and the stem remain same.
Condition not verified: Roll → Roll

Dictionary Look-up Stemmers:-

- A simple stemmer look-up the inflected form in a look-up table.
- The "Original term" (or) stemmed version of the term is looked up in a dictionary and its replaced by stem.
- This technique has been implemented in the "INQUERY" and "Retrieval ware systems".
- The "INQUERY" system uses a stemming technique called "Kstem".
K-STEM is a morphological analyzer that conflates word variants to a root form.
ex:- "memorial" and "memorize" reduced to "memory". But "memorial" and "memorize" are not synonyms and they have very different meanings.