

UNIT-II

Cataloging and Indexing: Objectives, Indexing Process, Automatic Indexing, Information Extraction.

Data Structures: Introduction, Stemming Algorithms, Inverted file structures, N-gram data structure, PAT data structure, Signature file structure, Hypertext data structure.

CATALOGING AND INDEXING:

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

- Process can be manual or automatic.
- Creating a direct search in document data base or in direct search through index files.
- Concept based representation: instead of transforming the input into a searchable format some systems transform the input into different representation that is concept-based Search and return item as per the incoming items.

History of indexing: Shows the dependency of information processing capabilities on manual and then automatic processing systems.

- Indexing originally called cataloguing oldest technique to identify the contents of items to assist in retrieval.
- Items overlap between full item indexing, public and private indexing of files.

Objectives:

The public file indexer needs to consider the information needs of all users of library system. Items overlap between full item indexing, public and private indexing of files.

- Users may use public index files as part of search criteria to increase recall.
- They can constrain the research by private index files
- The primary objective of representing the concepts with in an item to facilitate users finding relevant information.
- Users may use public index files as part of search criteria to increase recall.
- They can constrain there search by private index files.
- The primary objective of representing the concepts with in an item to facilitate users finding relevant information.

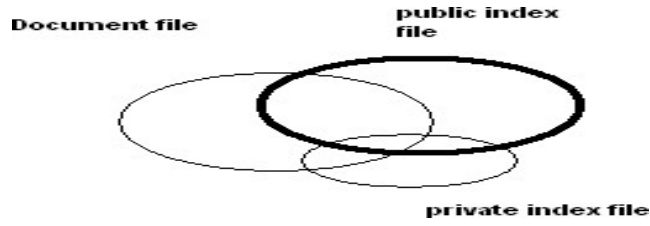
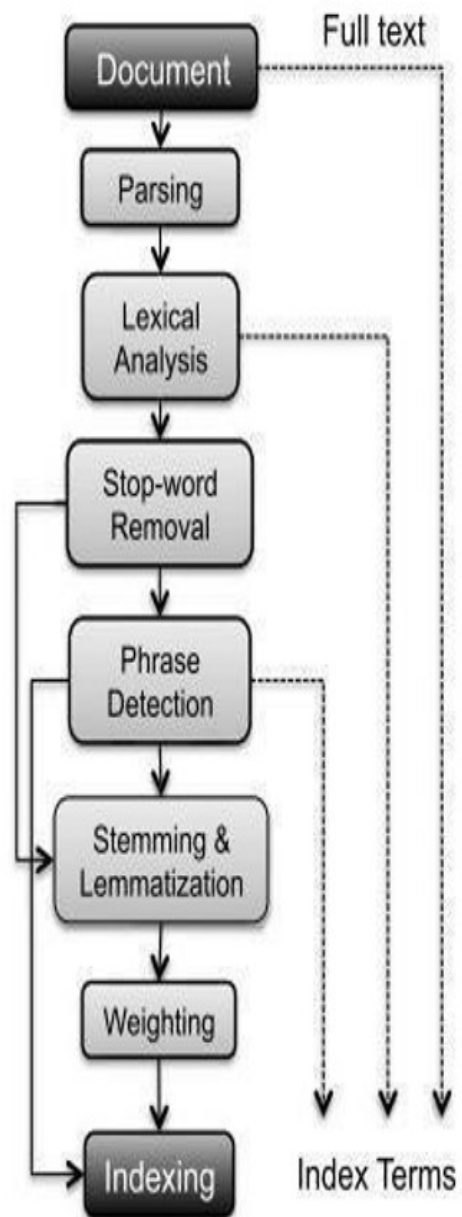


Fig: Indexing process

1. Decide the scope of indexing and the level of detail to be provided. Based on usage scenario of users.
2. Second decision is to link index terms together in a single index for a particular concept.

Fig. 2.3 Text processing phases in an IR system



Document Parsing. Documents come in all sorts of languages, character sets, and formats; often, the same document may contain multiple languages or formats, e.g., a French email with Portuguese PDF attachments. Document parsing deals with the recognition and “breaking down” of the document structure into individual components. In this preprocessing phase, unit documents are created, e.g., emails with attachments are split into one document representing the email and as many documents as there are attachments.

1. **Lexical Analysis.** After parsing, lexical analysis tokenizes a document, seen as an input stream, into words. Issues related to lexical analysis include the correct identification of accents, abbreviations, dates, and cases. The difficulty of this operation depends much on the language at hand: for example, the English language has neither diacritics nor cases, French has diacritics but no cases, German has both diacritics and cases. The recognition of abbreviations and, in particular, of time expressions would deserve a separate chapter due to its complexity and the extensive literature in the field for current approaches.

2. **Stop-Word Removal.** A subsequent step optionally applied to the results of lexical analysis is stop-word removal, i.e., the removal of high-frequency words. For example, given the sentence “search engines are the most visible information retrieval applications” and a classic stop words set such as the one adopted by the Snowball stemmer,¹ the effect of stop-word removal would be: “search engine most visible information retrieval applications”.

3. **Phrase Detection.** This step captures text meaning beyond what is possible with pure bag-of-word approaches, thanks to the identification of noun groups and other phrases. Phrase detection maybe approached in several ways, including rules (e.g., retaining terms that are not separated by punctuation marks), morphological analysis, syntactic analysis, and combinations.

4. **For example,** scanning our example sentence “search engines are the most visible information retrieval applications” for noun phrases would probably result in identifying “search engines” and “information retrieval”.

5. **Stemming and Lemmatization.** Following phrase extraction, stemming and lemmatization aim at stripping down word suffixes in order to normalize the word. In particular, stemming is a heuristic process that “chops off” the ends of words in the hope of achieving the goal correctly most of the time; a classic rule-based algorithm for this was devised by Porter [280]. According to the Porter stemmer, our example sentence “Search engines are the most visible information retrieval applications” would result in: “Search engine are the most visible inform retrieval”.

Lemmatization is a process that typically uses dictionaries and morphological analysis of words in order to return the base or dictionary form of a word, thereby collapsing its inflectional forms (see, e.g., [278]). For example, our sentence would result in “Search engine are the most visible information retrieval application” when lemmatized according to a WordNet-based lemmatizes

Weighting. The final phase of text preprocessing deals with term weighting. As previously mentioned, words in a text have different descriptive power; hence, index terms can be weighted differently to account for their significance with in a document and or a document collection. Such a weighting can be binary, e.g., assigning 0 for term absence and 1 for presence.

SCOPE OF INDEXING

- When perform the indexing manually, problems arise from two sources the author and the indexer the author and the indexer.
- Vocabulary domain may be different the author and the indexer.
- This results in different quality levels of indexing.
- The indexer must determine when to stop the indexing process.
- Two factors to decide on level to index the concept in a item.
- The exhaustively and how specific indexing is desired.
- Exhaustively of index is the extent to which the different concepts in the item are indexed.
- For example, if two sentences of a 10-page item on microprocessors discussion-board caches, should this concept be indexed Specific relates preciseness of index terms used in indexing.
- For example, whether the term “processor” or the term “microcomputer” or the term “Pentium” should be used in the index of an item is based upon the specificity decision.
- Indexing an item only on the most important concept in it and using general index terms yields low exhaustively and specificity.
- Another decision on indexing is what portion of an item to be indexed Simplest case is to limit the indexing to title and abstract (conceptual) zone.
- General indexing leads to loss of precision and recall.

PRE ORDINATION AND LINKAGES

- Another decision on linkages process whether linkages are available between index terms for an item.
- Used to correlate attributes associated with concepts discussed in an item. This process is called preordination.
- When index terms are not coordinated at index time the coordination occurs at search time. This is called post coordination, implementing by “AND” index terms.
- Factors that must be determined in linkage process are the number of terms that can be related.
- Ex., an item discusses the drilling of oil wells in Mexico by CITGO and the introduction of oil refine rise in Peru by the U.S.”

DATA STRUCTURES

- Introduction to Data Structures
- Stemming Algorithms
- Inverted File Structure
- N-Gram Data Structure
- PAT Data Structure
- Signature File Structure
- Hypertext and XML Data Structures

your roots to success...

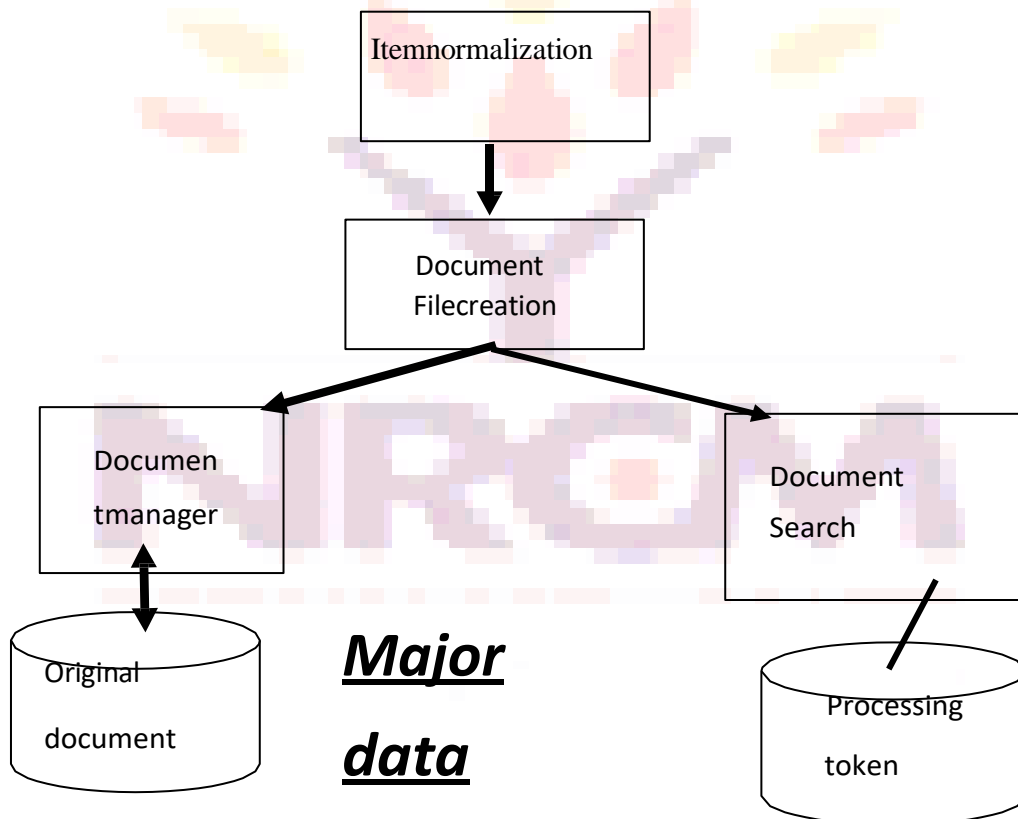
Datastructure:

The knowledge of data structure gives an insight into the capabilities available to the system.

- Each data structure has a set of associated capabilities.
- Ability to represent the concept and their r/s.
- Supports location of those concepts

Introduction Two major data structures in any IRS:

1. One structure stores and manages received items in their normalized form is called document manager
2. The other data structure contains processing tokens and associated data to support search.



Results of a search are references to the items that satisfy the search statement which are passed to the document manager for retrieval.

Focus: on data structure that support search function

Stemming: is the transformation often applied to data before placing it in the searchable data structure.

Stemming represents concept (word) to a canonical (authorized; recognized; accepted) morphological (the patterns of word formation in a particular Language) representation. Risk with

stemming: concept discrimination information may be lost in the process. Causing decrease in performance.

Advantage: has a potential to increase recall. STEMMING ALGORITHMS

- Stemming algorithm is used to improve the efficiency of IR and improve recall.
- Conflation (the process or result of fusing items into one entity; fusion; amalgamation) is a term that is used to refer to mapping multiple morphological variants to a single representation (stem).
- Stem carries the meaning of the concept associated with the word and the affixes (ending) introduces subtle (slight) modification of the concept.
- Terms with a common stem will usually have similar meanings, for example:
 - Ex: Terms with a common stem will usually have similar meanings, for example:
 - CONNECT
 - CONNECTED
 - CONNECTING
 - CONNECTION
 - CONNECTIONS
- Frequently, the performance of an IR system will be improved if term groups such as this are conflated into a single term. This may be done by removal of the various suffixes -ED, -ING, -ION, IONS to leave the single term CONNECT
- In addition, the suffix stripping process will reduce the total number of terms in the IR system, and hence reduce the size and complexity of the data in the system, which is always advantageous.
- ❖ Major usage of stemming is to improve recall.
- Important for a system to categorize a word prior to making the decision to stem.
- Proper names and acronyms (A word formed from the initial letters of a name say IARE ...) should not have stemming applied.
- Stemming can also cause problems for natural language processing NLP systems by causing loss of information.

PORTERSTEMMINGALGORITHM

- Based on a set of conditions of the stem
 - A consonant in a word is a letter other than A, E, I, O or U, some important conditions are
1. The measure of a stem is a function of sequence of vowels (V) followed by a sequence of consonant (C).
 2. C(VC)^mV. m is number of VC repeats. The case m=0 covers the null word.
 3. *<X>-stem ends with a letter X. *v*-stem contains a vowel
 4. *d*-stem ends in double consonant (e.g. -TT, -SS).
 5. *o*-stem ends in consonant vowel sequence where the final consonant is not w, x, y (e.g. -WIL, -HOP).

Suffix condition: `stake the form current_suffix == pattern`
 Actions are in the form `old_suffix -> .New_suffix`

Rules are divided into steps to define the order for applying the rule. Examples of the rules

Step	Condition	Suffix	Replacement	Example
1a	Null	Sses	Ss	Stresses -> stress
1b	*v*	Ing	Null	Making -> mak
1b1	Null	At	Ate	Inflated -> inflate
1c	*v*	Y	I	Happy -> happi
2	m > 0	aliti	al	Formaliti -> formal
3	m > 0	Icate	Ic	Duplicate -> duplie
4	m > 1	Able	Null	Adjustable -> adjust
5a	m > 1	e	Null	Inflate -> inflat
5b	m > 1 and *d	Null	Single letter	Control -> control

2. Dictionarylookupstemmers

- ❖ Useofdictionarylookup.
- ❖ Theoriginaltermorstemmedversionofthetermislookedupinadictiona
ryandreplaced bythestemthatbestrepresentsit.
- ❖ ThistechniquehasbeenimplementedinINQUERYandRetrieval
waresystems-

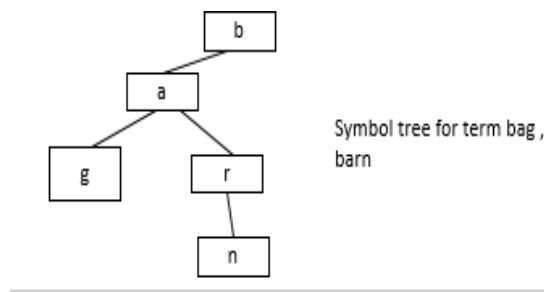
INQUERYsystemusesthetechniquecalledKstem.

- ❖ Kstemisamorphologicalanalyzerthatconflateswordsvariantstoarootform.
 - ❖ Itrequiresawordtobeinthedictionary
 - ❖ Kstemuses6majordatafilestocontrolandlimitthestemmingprocess.
1. Dictionaryofwords(lexicon)
 2. Supplementallistofwordsfordictionary
 3. Exceptionallistofwordsthatshouldretaina,,e“attheend(e.g.,“suites”to“suite
”but“suited”to“suit”).
 4. Direct_conflation-wordpairsthatoverridestemmingalgorithm.
 5. County_nationality_conflation(BritishmapstoBritain)
 6. Propernouns--thatshouldnotbestemmed
- ❖ New words that are not special forms (e.g., dates, phone numbers) arelocatedinthedictionarytodeterminesimplerformsbystrippingoffsuffixes andrespellingpluralsasdefinedinthedictionary.

3. Successorstemmers:

- Basedonlengthofprefixes.
- Thesmallestunitofspeechthatdistinguishesonwordfromanother
- Theprocessusessuccessorvarietiesforaword.

Usesinformationtodivideawordintosegmentsandselectsonofthesegmentstostem.



Successor variety of words are used to segment a word by applying one of the following four methods.

1. Cutoff method: a cutoff value is selected to define the stem length.
2. Peak and plateau: a segment break is made after a character whose successor variety exceeds that of the character.
3. Complete word method: break on boundaries of complete words.
4. Entropy method: use the distribution method of successor variety letters.

1. Let $|Dak|$ be the number of words beginning with k length sequence of letters a .

2. Let $|Dakj|$ be the number of words in Dak with successor j .

3. The probability that a member of Dak has the successor j is given as $|Dakj|/|Dak|$. The entropy of $|Dak|$ is

$$H_{ak} = -\sum_j \left(\frac{|Dakj|}{|Dak|} \right) (\log \left(\frac{|Dakj|}{|Dak|} \right)) p = 1$$

After a word has been segmented these segments to be used as stems must be selected.

Hafer and Weiss selected the following rule

If (first segment occurs in ≤ 12 words in database) First segment is stem Else (second segment is stem)

INVERTED FILE STRUCTURE

Inverted file structure

Most common data structure

Inverted file structures are composed of three

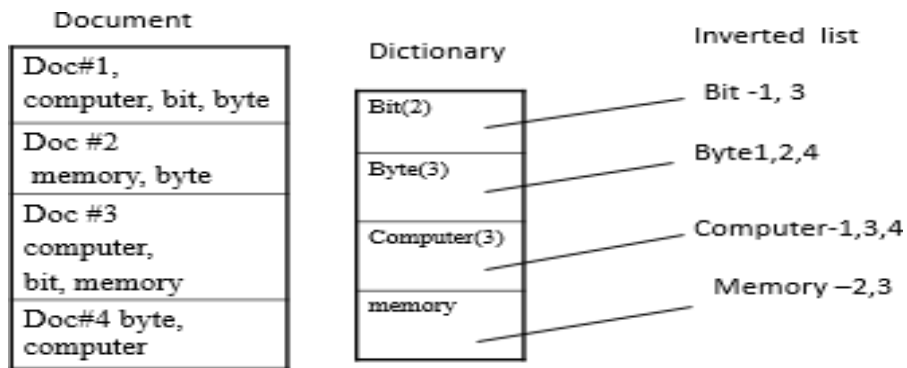
files

- The document file

1. The inversion list (Posting List)
2. Dictionary
3. The inverted file: based on the methodology of storing an inversion of documents.
4. For each word a list of documents in which the word is found is stored (inversion of document)
5. Each document is given a unique numerical identifier that is stored in inversion list. Dictionary is used to locate the inversion list for a particular word.

This is a sorted list (processing tokens) in the system and a pointer to the location of its inversion list.

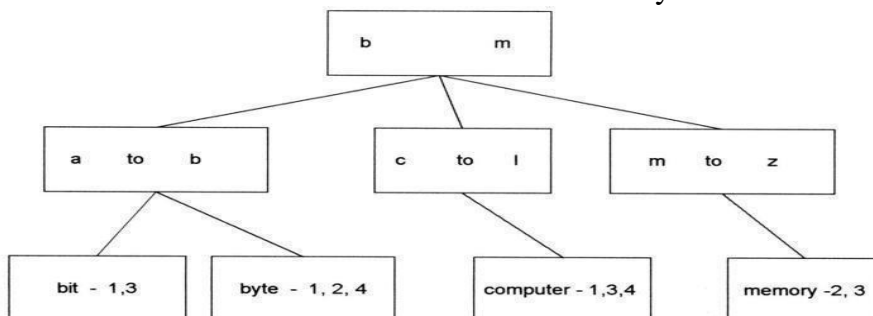
Dictionary can also store other information used in query optimizations such as length of inversion list to increase the precision.



- Use zoning to improve
- Precision and Restrict entries.
- Inversion list consists of document identifier for each document in which the word is found.

Ex: bit 1(10), 1(12) 1(18) is in 10, 12, 18 position of the word bit in the document #1.

- When a search is performed, the inversion lists for the terms in the query are located and appropriate logic is applied between inversion lists.
- Weights can also be stored in the inversion list.
- Inversion lists are used to store concept and their relationship.
- Words with special characteristics can be stored in their own dictionary. Ex: Date... which require dating and numbers.
- Systems that support ranking are re-organized in ranked order.
- B-trees can also be used for inversion instead of dictionary.
- The inversion lists may be at the leaf level or referenced in higher level pointers.
- AB-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 smallest to larger.
 - All leaves are at the same level or differ by at most one level.



N-GRAM DATA STRUCTURE

Transposition of 2 adjacent

computer

chars

- Zamora showed trigram analysis provided a viable data structure for identifying misspellings and transposed characters.
- This impacts information systems as a possible basis for identifying potential input errors for correction as a procedure within the normalization process.
- Frequency of occurrence of n-gram patterns can also be used for identifying the language of a finite m.
- Trigrams have been used for text compression and to manipulate the length of index terms.
- To encode profiles for the Selective Dissemination of Information.
- To store these searchable document files for retrospective search databases.

Advantage:

They place a finite limit on the number of searchable tokens MaxSeg

$n = \left(\frac{m}{l} \right)$

) n maximum number of unique n-grams that can be generated. "n" is

length of n-grams

number of process able symbols Disadvantage:

longer than the size of inversion list increase. Performance has

85% precision.

PAT data structure (practical algorithm to retrieve information coded in alphanumeric)

- PAT structure or PAT tree or PAT array: continuous text input data structures (string like N-Gram data structure).
- The input stream is transformed into a searchable data structure consisting of substrings, all substrings are unique.
- Each position in an input string is an anchor point for a substring.
- In creation of PAT tree each position in the input string is the anchor point for a sub-string that starts at that point and includes all new text up to the end of the input.
- Binary tree, most common class for prefix search, But PAT trees are sorted logically which facilitate range search, and more accurate than inversion file.
- PAT trees provide alternate structure if supporting strings search.

Text EconomicsforWarsawiscomplex.

 sistring1Economics

forWarsawiscomplex.

sistring2

conomicsforWarsawiscomplex.

sistring5omicsforWarsawiscomplex.s

istring 10 for Warsaw is

complex.sistring20w is

complex.sist

ring30ex.

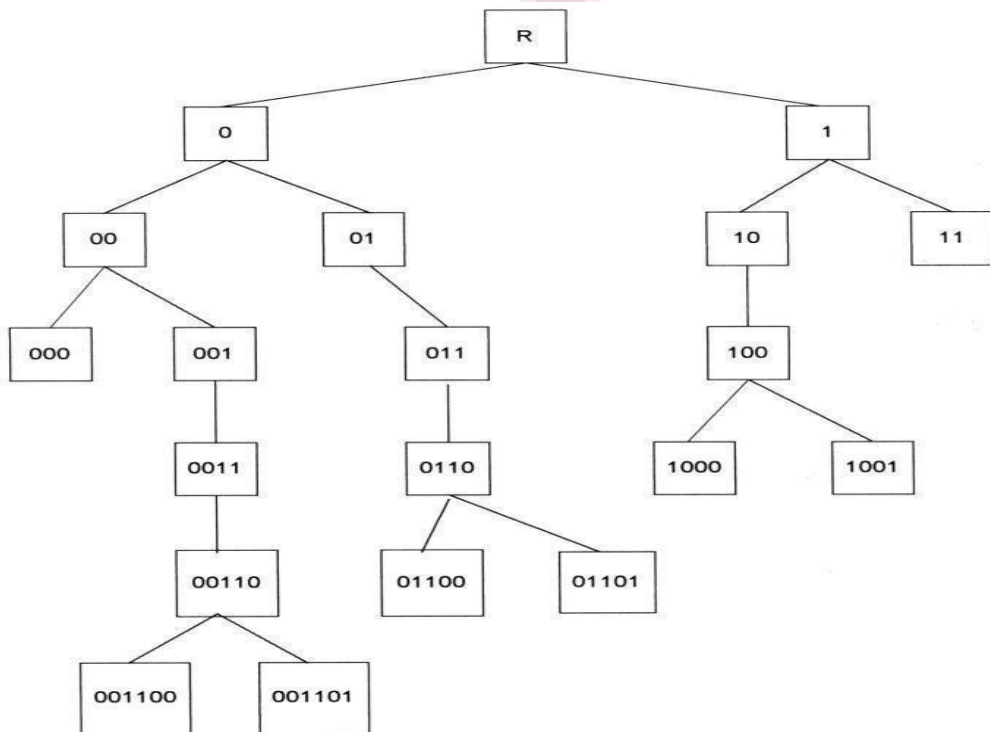
Examplesofsistrings

- Thekeyvaluesarestoredattheleafnodes(bottomnodes)inthePATTree.
- Foratextinputofsize“n”thereare“n”leafnodesand“n-1”atmosthigherlevelnodes.
- Itispossibletoplaceadditionalconstraintsonsistringsfortheleafnodes.
- Ifthebinaryrepresentationsof“h”is(100),“o”is(110),“m”is(01) and“e”is(101)thentheword“home”producestheinput

100110001101 Usingthesistrings.

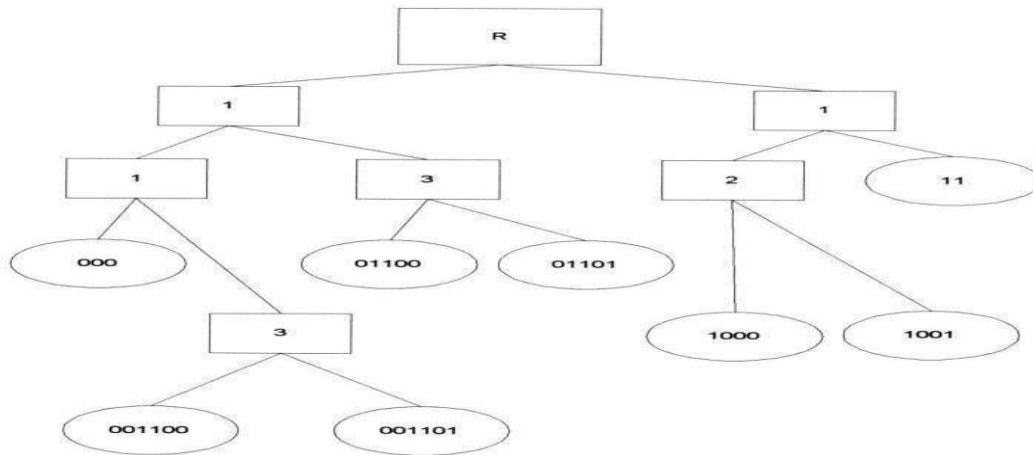
INPUT	100110001101
sistring1	1001....
sistring2	001100...
sistring3	01100....
sistring4	11.....
sistring5	1000...
sistring6	000.....
sistring7	001101
sistring8	01101

ThefullPATbinarytreeis



Thevalueintheintermediatenodes(indicatedbyrectangles)isthenumberofbitstoskipuntilthen ext

bit to compare that causes differences between similar terms.



Skipped final version of PAT tree

Signature file structure

- The coding is based upon words in the code.
- The words are mapped into word signatures.
- A word signature is fixed length code with a fixed number of bits set to 1.
- The bit positions that are set to one are determined via a hash function of the word.
- The word signatures are Ored together to create a signature of an item..
- Partitioning of words is done in block size, which is nothing but set of words, Code length is 16 bits.
- Search is accomplished by template matching on the bit position.
- provide a practical solution applied in parallel processing, distributed environment etc.
- To avoid signatures being too dense with “1”s, a maximum number of words is specified and an item is partitioned into blocks of that size.
- The block size is set at five words, the code length is 16 bits and the number of bits that are allowed to be “1” for each word is five.
- TEXT: Computer Science graduate students study (assume block size is five words)

WORD	Signature
computer	00010110 0000 0110
Science	1001000011100000
graduate	1000010101000010
students	0000011110000100
study	00000110 0110 0100

BlockSignature 1001011111100110

Superimposed Coding

Application(s)/Advantage(s)

- Signature files provide a practical solution for storing and locating information in a number of different situations.
- Signature files have been applied as medium sized databases, databases with low frequency of terms, WORM devices, parallel processing machines, and distributed environments

HYPERTEXT AND XML DATA STRUCTURES

- ❖ The advent of the Internet and its exponential growth and wide acceptance as a new global information network has introduced new mechanisms for representing information.
- ❖ This structure is called hypertext and differs from traditional information storage data structures in format and use.
- ❖ The hypertext is stored in HTML format and XML.
- ❖ Both of these languages provided detailed descriptions for subsets of texts similar to the zoning.
- ❖ Hypertext allows one item to reference another item via an embedded pointer.
- ❖ HTML defines internal structure for information exchange over WWW on the internet.
- ❖ XML: defined by DTD, DOM, XSL, etc.

Document and term clustering

Two types of clustering:

- 1) clustering index terms to create a statistical thesaurus and
- 2) Clustering items to create document clusters. In the first case clustering is used to increase recall by expanding searches with related terms. In document clustering the search can retrieve items similar to an item of interest, even if the query would not have retrieved the item. The clustering process is not precise and care must be taken on use of clustering techniques to minimize the negative impact misuse can have.

Questions:

1. Explain about Cataloging and Indexing?
2. Write about data structures? Explain about Stemming Algorithms?
3. Write about a) Inverted File Structure b) N-Gram Data Structure c) PAT Data Structure?
4. Explain about Hypertext and XML Data Structures?
5. a) Explain about Probabilistic Weighting?
b) What is Vector Space Retrieval Model with an example?
6. Explain about Hidden Markov Models?

