

UNIT - 3

Automatic Indexing:

Automatic Indexing:

1. classes of Automatic Indexing

2. Statistical Indexing

3. Natural Language

4. Concept Indexing

5. Hypertext linkages

Document and Term clustering:

6. introduction to clustering

7. Thesaurus Generation

8. Item clustering

9. Hierarchy of clusters

Automatic Indexing:

It focuses on the process and algorithms to perform indexing. The indexing process is a transformation of an item that extracts the semantics of the topics discussed in the item. The extracted information is used to create the processing tokens and the search-able data structure. The semantic of the item not only refers to the subjects discussed in the item but also in weighted systems, the depth to which the subject is discussed. The index can be based on the full text of the item, automatic or manual generation of a subset of terms to represent the item, natural language representation of the item or abstraction to concepts in the item.

1. Classes of Automatic Indexing:

Automatic Indexing is the process of analyzing an item to extract the information to be permanently kept in an index. This process is associated with the generation of the searchable data structures associated with an item.

The below fig shows is expanded to show where the search process relates to the indexing process.

The left side of the fig including Identify Processing tokens, Apply stop lists, Characterize tokens, Apply stemming & Create Searchable data structure is all part of the indexing process.

All systems go through an initial stage of zoning & identifying the processing tokens used to create the index. Some systems automatically divide the document up into fixed length passages at localities, the indexed filters, such as stop lists & stemming algorithms, are

frequently applied to reduce the number of tokens to be processed. The next step depends upon the search strategy of a particular system.

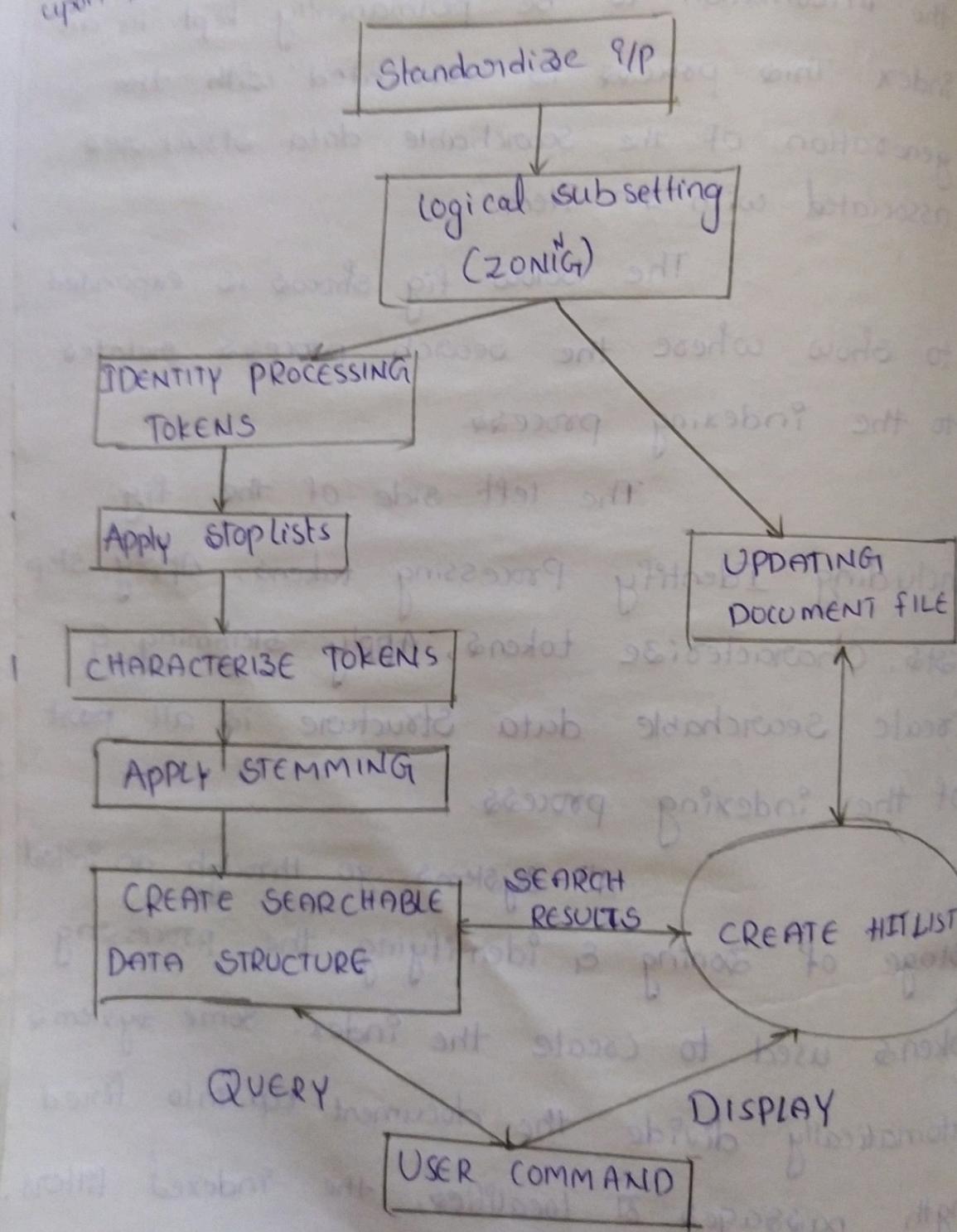


Fig: Data Flow in Information Processing System

Search strategies can be classified as statistical, natural language, & concept. An index is the data structure created to support the search strategy.

Statistical Indexing & Strategies:

It covers the broadest range of indexing techniques & are the most prevalent in commercial systems. The words/phrases are the domain of searchable values.

Natural language: It approaches perform the similar processing token identification as in statistical techniques, but then additionally perform varying levels of natural language parsing of the item eg: present, past, future actions.

Concept index: Concept indexing uses the words within an item to correlate to concepts discussed in the item. This is a generalization of the specific words to values used to index the item.

Finally a special class of indexing can be defined by creation of hypertext linkages. These linkages provide virtual threads of concepts between items vs directly defining the concept within an item.

2. Statistical Indexing:

It uses frequency of occurrence of events to calculate the num to indicate potential relevance of an item. The documents are found by normal Boolean Search and then statistical calculation are performed on the hit file, ranking the o/p (eg: term-frequency algorithm).

1. probability weighting

2. Vector weighting

- > Simple Term Frequency algorithm
- > Inverse Document Frequency algorithm
- > Signal weighting
- > Discrimination value
- > problems with the weighting Schemas
- > problems with the vector model

3. Bayesian Model.

1. Probabilistic weighting:

attempt to calculate a probabilistic value that should be invariant to both calculation method & text corpora (Large collection of written/ spoken texts).

The probabilistic approach is based on direct application of theory of probability to IRS.

Advantage - uses the probability theory to develop the algorithm. This

This allows easy integration of the final results when searches are performed across multiple databases & use different search algorithms.

The use of probability theory is a natural choice - the basis of evidential reasoning (drawing conclusions from evidence).

This is summarized by PRP (Probability Ranking Principle) & its plausible corollary (reasonable result).

PRP - hypothesis - If a preference retrieval system's response to each request is a ranking of the documents in order of decreasing probability of usefulness to the user, the overall effectiveness of the system to the users is best obtainable on the basis of the data available.

Plausible corollary: The techniques for estimating the probabilities of usefulness for dp ranking

in IR is standard probability theory & statistics
probabilities are based on binary condition the
item is relevant or not. But in information sys
the relevance of an item is a continuous funⁿ
from non-relevant to absolutely useful.

Source of problems: In application of probability
theory come from lack of accurate data &

Simplified assumptions that are applied to
mathematical modeling.
cause the results of probabilistic

approaches in ranking items to be less accurate
than other approaches.

Advantage of probabilistic approach is that it
can identify its weak assumptions & work to
strengthen them.

There are many different ideas
in which the probabilistic approach may be
applied. ex: logistical regression.

The approach starts by defining a "model" or
System.

In retrieval system there exist query

Statistics
on the
ion sys
fun?
att 9

probability
&
to
locator
stick)
accurate
P(dadoc9)
that it
work to
model 10
exist query

q_i and a document term d_j which has a set of attributes ($v_1 \dots v_n$) from the query (ex: counts of term frequency in the query), from the document (ex: count \hat{a} of term frequency in the document) and from the database (ex: total num of documents in the database divided by the num of documents indexed by the term).

The logistic reference model uses a random sample of query document terms for which binary relevance judgment has been made from judgment from samples.

Logarithm O is the logarithm of the odds (logodds) of relevance for terms t_k which is present in document D_j and query Q_i

$$\log(O(R|Q_i, D_j, t_k)) = c_0 + c_1 v_1 + \dots + c_n v_n$$

The logarithm that the i^{th} Query is relevant to the j^{th} document is the sum of the logodds for all terms:

$$\log(O(R|Q_i, D_j)) = \sum_{k=1}^q [\log(O(R|Q_i, D_j, t_k)) - \log(O(R))]$$

The inverse logistic transformation is applied to obtain the probability of relevance of a document to a query:

$$P(R|Q_i, D_j) = \frac{1}{1 + e^{-\log(O(R|Q_i, D_j))}}$$

The coefficients of the equation for logodds is derived for a particular database using a random sample of query-document-term-relevance quadruples & used to predict odds of relevance for other query-document pairs.

Additional attributes of relative frequency in the query (QRF), relative frequency in the document (DRF) and relative frequency of the term in all the documents (RFAD) were included, producing the logodds formula:

$$z_j = \log(O(R|t_j)) = c_0 + c_1 \log(QAF) + c_2 \log(QRF) + c_3 \log(DAF) + c_4 \log(DRF) + c_5 \log(IDF) + c_6 \log(RFAD)$$

$$QRF = QAF / (\text{total num of terms in the query})$$

$$DRF = DAF / (\text{total num of words in the document})$$

$$RFAD = (\text{total num of term occurrences in the database}) / (\text{total num of all words in the database}).$$

Logs are used to reduce the impact of frequency information, then smooth out skewed distributions.

The coefficients & $\log(O(R))$ were calculated creating the final formula for ranking for query vector Q , which contains q terms:

$$\log(O(R|Q)) = -5.138 + \sum_{k=1}^q (z_k + 5.138)$$

The logistic inference method was applied to the test database along with the Cornell SMART vector system, inverse document frequency & cosine relevance weighting formulas.

The logistic inference method outperformed the vector method. Attempts have been made to combine different probabilistic techniques to get a more accurate value.

This combination of logarithmic odds has not presented better results. The objective is to have the strong points of different techniques compensate for weaknesses.

To date this combination of probabilities using averages of log odds has not produced better results and in many cases produced worse results.

2. Vector Weighting:

Earliest System that investigated statistical approach is SMART system of Cornell university. The system is based upon a vector model. A vector is one-dimensional set of values, where the position of each value in the set is fixed & represents a particular domain.

Each position in the vector represents a processing token. There are 2 approaches to the domain of values in the vector: binary & weighted. Under binary approach the domain contains a value of 1 or 0. [..1 represented existence of processing tokens in the item].

In the weighted approach, the domain is set of +ve values. The value for each processing token represents the relative importance of the item.

	Petroleum	Mexico	Oil	Taxes	Refineries	Shipping
Binary	(1 , 1 , 1 , 0 , 1 , 0)					
Weighted	(2.8 , 1.6 , 3.5 , 3 , 3.1 , .1)					

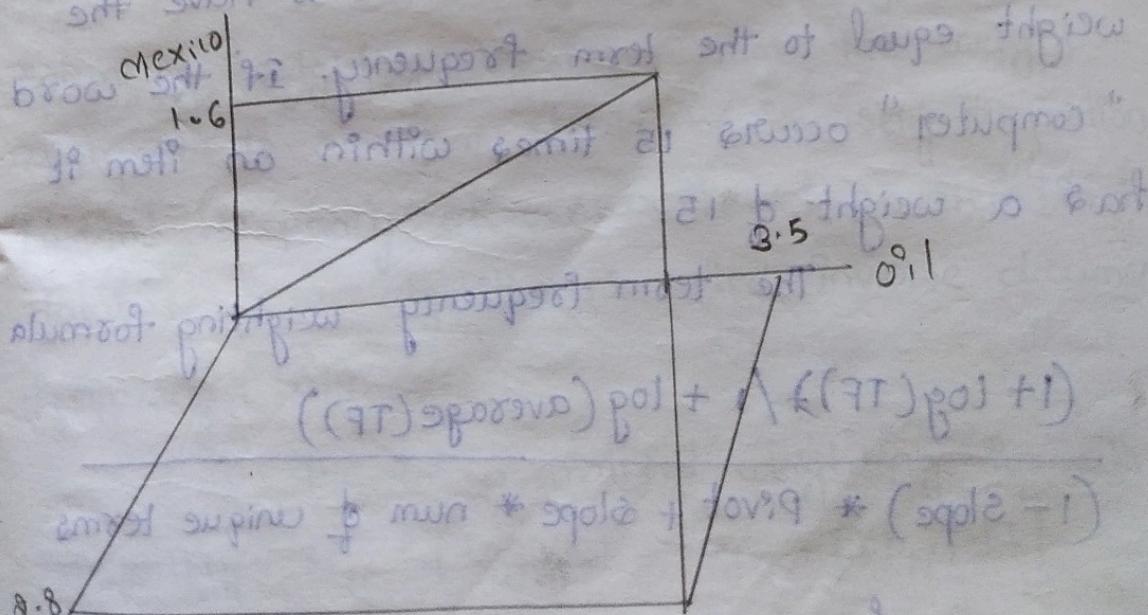
Fig: Binary & Vector Representation of an Item.

The above fig shows how an item that discusses petroleum refineries in mexico would be represented.

A weighted vector acts the same as a binary vector but it provides a range of values that accommodates a variance in the value of relative importance of processing tokens in representing the item.

The use of weights also provides a basis for determining the rank of an item.

The vector approach allows for a mathematical & a physical representation using a vector space model. Each processing token can be considered another dimension in an item representation space. In below fig shows a 3-dimensional vector representation assuming there were only 3 processing tokens, Petroleum, Mexico & Oil.



Petroleum
Oil
Mexico

0.8

0.1

1.6

$(FT)(P0) + A f(FT)P0 + 1)$

sum engine = min * sqols + max * (sqols - 1)

fig: Vector Representation of the word

iii) Simple term frequency algorithms: In both the unweighted & weighted approaches, an automatic indexing process implements an algorithm to determine the weight to be assigned to a processing token for a particular item.

In statistical systems, the data that are potentially available for calculating a weight are the frequency of occurrence of the processing token in an existing item (i.e., term frequency - TF), the frequency of occurrence of the processing token in the existing database (i.e., total frequency - TOTF) & the number of unique items in the database that contain the processing token.

Simplest approach is to have the weight equal to the term frequency. If the word "computer" occurs 15 times within an item it has a weight of 15.

The term frequency weighting formula:

$$\frac{(1 + \log(TF))}{1 + \log(\text{average}(TF))}$$

$$(1 - \text{slope}) * \text{pivot} + \text{slope} * \text{num of unique terms}$$

where slope was set at .2 and pivot was set to the average no of unique terms

both
an automatic
to determin.
processing
data that
a weight
processing
nency - TF),
ing Token
nency - TOTF)
database
e 8000
processing
e the
he word
item it
formula :
terms
nd pivot
ms

occuring in the collection. Slope & pivot are constants for any document/query set. This leads to the final algorithm that weights each term by the above formula divided by the pivoted normalization: $((1 + \log(TF)) / (1 + \log(\text{average}(TF))) / (\text{slope}) (\text{No. unique terms}) + (1 - \text{slope}) * (\text{pivot}))$.

ii) Inverse Document frequency (IDF):
One of the objectives of indexing an item is to discriminate the semantics of that item from other items in the database. If the token computer represents a concept used in an item, but it does not help a user find the specific information being sought since it returns the complete DB.

$\text{IDF} = \log((\text{No. of documents}) / (\text{No. of documents containing the term}))$

This leads to the general stmt enhancing weighting algorithms that the weight assigned to an item should be inversely proportional to the frequency of occurrence of an item in the database. This algorithm is called inverse document frequency (IDF).

The un-normalized weighting formula is:

$$\text{WEIGHT}_{ij} = \text{TF}_{ij} * [\log_2(n) - \log_2(IIDF_{ij}) + 1]$$

where WEIGHT_{ij} is the vector weight that is assigned to term j in item i ,

TF_{ij} (term frequency) is the frequency of term j in item i , n is the number of items in the database, and IF_{ij} (item frequency or document frequency) is the number of items in the database containing term j .

Ex: Assume that the term "oil" is found in 128 items, "Mexico" is found in 16 items and "refinery" is found in 1024 items.

If a new item arrives with all 3 terms in it, "oil" found 4 times, "mexico" found 8 times, and "refinery" found 10 times & there are 2048 items in the total database.

weight calculations using inverse document frequency.

$$\text{weight}_{\text{oil}} = 4 * (\log_2(2048) - \log_2(128) + 1) = 4 * (11 - 7 + 1) = 20$$

$$\text{weight}_{\text{mexico}} = 8 * (\log_2(2048) - \log_2(16) + 1) = 8 * (11 - 4 + 1) = 64$$

$$\text{weight}_{\text{refinery}} = 10 * (\log_2(2048) - \log_2(1024) + 1) = 10 * (11 - 10 + 1) = 20.$$

with the resultant inverse document frequency.

$$\text{Item vector} = (20, 64, 20).$$

The value of "n" and IF vary as items are added and deleted from the database.

iii) Signal
adjusts the
based upon
term in
frequency
the item
ability
for example
DRILL
frequencies

Item

value of
the pro

the fin
instance

(iii) Signal weighting:
Inverse document frequency
adjusts the weight of processing token for an item
based upon the number of items that contain the
term in the existing database.

It does not account for the term
frequency distribution of the processing token in
the items that contain the term - can affect the
ability to rank items.

For example: Assume the terms "SAW" and
"DRILL" are found in 5 items with the following
frequencies:

Item	Distribution	SAW	DRILL
A	10	3	2
B	10	18	10
C	10	18	10
D	10	18	60
E	10	18	60

In information theory, the information content
value of an object is inversely proportional to
the probability of occurrence of the item.

An instance of an event that occurs all
the time has less information value than an
instance of a seldom occurring event.

This is represented as INFORMATION = $-\log_2(p)$, where p is the probability of occurrence of event "P". The information value for an event that occurs 5% of the time is: INFORMATION = $-\log_2(0.05)$

$$= -(-10) = 10$$

50% of the time is: INFORMATION = $-\log_2(.50)$

$$= -(-1) = 1$$

If there are many independent occurring events then the calculation for the average information value across the events is:

$$\text{AVE_INFO} = - \sum_{k=1}^n p_k \log_2(p_k)$$

Its value decreases proportionally to increases in variance in the values of p_k can be defined as TF_{ik} / TOTF_k

Formula for calculating the weighting factor called Signal can be used:

$$\text{Signal}_k = \log_2(\text{TOTF}) - \text{Ave_INFO}$$

It producing a final formula of:

$$\text{weight}_{ik} = TF_{ik} * \text{Signal}_k$$

$$\text{Weight}_{ik} = \text{TF}_{ik} * [\log_2(\text{TOTF}_k) - \sum_{i=1}^D \text{TF}_{ik} / \text{TOTF}_k \log_2 \frac{\text{TF}_{ik}}{\text{TOTF}_k}]$$

Based on example of item distribution for SAW and DRILL is

$$\text{Signal}_{SAW} = \log_2(50) - [5 * \{10/50 \log_2(10/50)\}] = 3.35$$

$$\begin{aligned} \text{Signal}_{DRILL} &= \log_2(50) - [2/50 \log_2(2/50) + 2/50 \log_2(2/50) + \\ &\quad 18/50 \log_2(18/50) + 10/50 \log_2(10/50) + 18/50 \log_2(18/50)] \\ &= 3.75 \end{aligned}$$

The weighting factor for term "DRILL" that does not have a uniform distribution is larger than that for term "SAW" and gives it a higher weight.

This technique could be used by itself or in combination with inverse document frequency or other algorithms. The overhead of the additional calculations required to get the values have not been demonstrated to produce better results.

It is a good example of use of information theory in developing information retrieval algorithms.

It is a good example of use of information theory in developing information retrieval algorithms.

It is a good example of use of information theory in developing information retrieval algorithms.

(iv) Dissemination Value:

Particulars
marked symbolic

Another approach to creating a weighting algorithm is based on the dissemination of value of the term. The all items appear the same, the harder it is to identify those that are needed.

Salton & Yang proposed a weighting algorithm that takes into consideration the ability for a search term to discriminate among items. They proposed use of a discrimination value for each term "i".

$$\text{DISCRIM}_i = \text{AVESIM}_i - \text{AVESIM}_{\text{all}}$$

where AVESIM is the average similarity of every item in the database & AVESIM_i is the same calculation except that term "i" is removed from all items.

DISCRIM_i value being +ve, close to 0/-ve.

» A +ve value indicates that removal of term "i" has increased the similarity of items. In this case, leaving the term in the database assists in discriminating items and is of value.

» A value close to zero, implies that the term's removal or inclusion does not change the there is no change in similarity.

Similarity

» If the the database more similar decreased

as a positive standard

out weight

based

Problems
set of 3
Schemes

missing to

information

of items value

items being

are char

between too

Compensat

tion p

» Ign

based

chang

Search

compar

Similarity b/w items.

» If the value is negative, the terms effect on the database is to make the items appear more similar since their average similarity decreased with its removal.

Once the value of DISCRIM is normalized as a positive number, it can be used in the standard weighting formula as:

$$\text{Weight}_{ik} = \text{TF}_{ik} * \text{DISCRIM}_k$$

↓ Problems with Weighting Schemes:

Often weighting schemes use information that is based upon processing token distributions across the database. The information database tends to be dynamic with new items always being added & to a lesser degree old items being changed or deleted. Thus these factors are changing dynamically.

There are a num of approaches to compensate for the constant changing values.

» Ignore the variances & calculate weights based upon current values, with the factors changing over time. Periodically rebuild the complete search database.

similarity

- » Use a fixed value while monitoring changes in the factors. When the changes reach a certain threshold, start using the new value & update all existing vectors with the new value.
- » Store the invariant variables (ex: term frequency within an item) & at search time calculate the latest weights for processing taken in items needed for search terms.

In the 1st approach, Periodically the database & all term weights are recalculated based

upon the most relevant recent updates to the database based on frequent updates to the database. In the 2nd approach, there is a recognition that for the most frequently occurring items, the aggregate values are large. As such, minor changes in the values have negligible effect on the final weight calculation.

The 3rd approach is the most accurate. The weighted values in the database only matter when they are being used to determine items to return from a query.

Solutions: If the system is using an inverted file search structure, this overhead is very minor.

The best
a query ag
different
weighting
the result
vis Problem
comes in
topics bei
for exam

discussion
Pensylva
mechanis
it's part

corrrelation
in a v
Thus

the search

Space

store with

no no

The best environments would allow a user to run a query against multiple different time periods & different databases that potentially use different weighting algorithms, & have the system integrate the results into a single ranked hit file.

vi) Problems with the vector model:

A major problem comes in the vector model when there are multiple topics being discussed in a particular item. For example, assume that an item has an in-depth discussion of "oil" in "Mexico" and also "coal" in "Pennsylvania". The vector model does not have a mechanism to associate each energy source with its particular geographic area.

There is no way to associate correlation factors b/w terms since each dimension in a vector is independent of the other dimension.

(Thus the item results in a high value in a search for "coal in Mexico".)

Another major limitation of a vector space is in associating positional information with a processing item term.

The concept of a vector space allows only one scalar value to be associated with each processing term for each item.

3. Bayesian model:

One way of overcoming the restrictions inherent in a vector model is to use a Bayesian approach to maintaining information on processing tokens. The Bayesian model provides a conceptually simple yet complete model for information systems. The Bayesian approach is based upon conditional probabilities.

This general concept can be applied to the search function as well as to creating the index to the database. The objective of information systems is to return relevant items.

Thus, the general case, using the Bayesian formula is $P(\text{REL} | \text{DOC}, \text{Query})$ which is interpreted as the probability of relevance (REL) to a search statement given a particular document and query.

In addition to search, Bayesian formulas can be used in determining the weights associated with a particular processing token in an

item. The objective of creating the index to an item is to represent the semantic information in the item. A Bayesian net can be used to determine the final set of processing tokens (called topic) & their weights.

In below fig shows a simple view of the process where T_i represents the relevance of topic "i" in a particular item and P_j represents a statistic associated with the event of processing token "j" being present in the item.

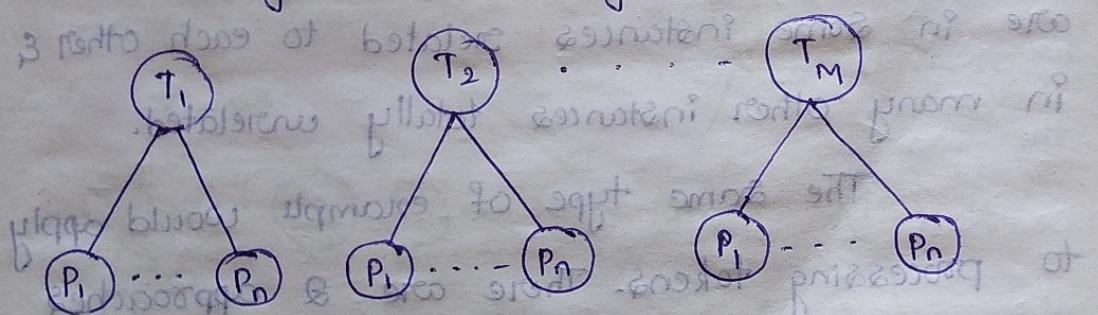


Fig: Bayesian Team Weighting

"topics" would be stored as the final index to the item. There is one major assumption made in this model:

Assumption of Binary Independence:

The topics & the processing token statistics are independent of each other.

- » The existence of the one topic is not related to the existence of the other topics.
- » The existence of one processing token is not related to the existence of other processing tokens.

In most cases this assumption is not true. Some topics are related to other topics & some processing tokens related to other processing tokens.

for example: The topics of "Politics" & "Economics" are in some instances related to each other & in many other instances totally unrelated.

The same type of example would apply to processing tokens. There are 2 approaches to handling this problem.

- » The 1st is to assume that there are dependencies, but that the errors introduced by assuming the mutual independence do not noticeably effect the determination of relevance of an item nor its relative rank associated with other retrieved items.
- » A 2nd approach can extend the rule to additional

layers to handle interdependencies. Thus an additional layer of Independent Topics (ITs) can be placed above the Topic layer & a layer of Independent Processing Tokens (IPs) can be placed above the processing token layer. The below fig shows the extended Bayesian n/w.

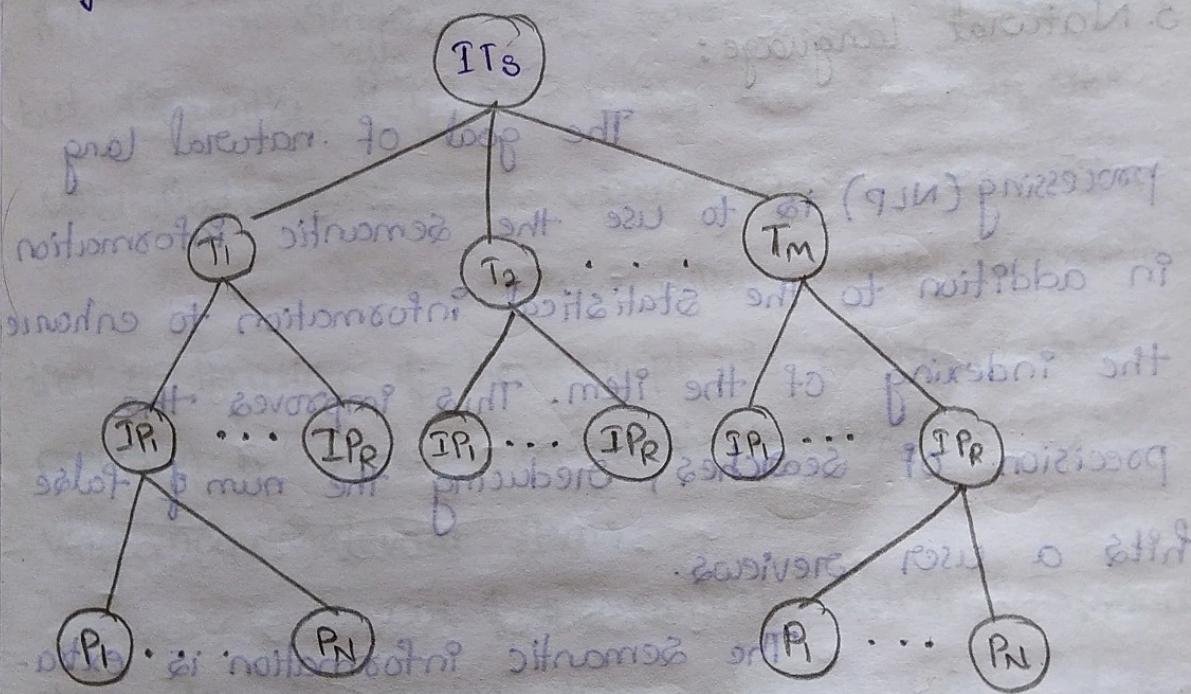


Fig: Extended Bayesian n/w.

Extending the n/w creates new processing tokens for those cases where there are dependencies b/w processing tokens. The new set of Independent processing Tokens can be used to define the attributes associated with the set of

topics selected to represent the semantics of an item. To compensate for dependencies b/w topics the final layer of Independent Topics is created. The degree to which each layer is created depends upon the error that could be introduced by allowing for dependencies b/w Topics or Processing Tokens.

3. Natural language:

The goal of natural language processing (NLP) is to use the semantic information in addition to the statistical information to enhance the indexing of the item. This improves the precision of searches, reducing the num of false hits a user reviews.

The semantic information is extracted as a result of processing the lang rather than treating each word as an independent entity. The simplest o/p of this process results in generation of phrases that become indexes to an item. More complex analysis generates thematic representation of events rather than phrases.

Natural
the conne
referred

g, Index

is to g

the info
information

but freq
user fir

tion &
precision
ving dr

2010

chea

COHES

JUSTO

COHE

OPEN

di

base

PAIR

Natural Language Processing can also combine the concepts into higher level concepts sometimes referred to as thematic representations.

i) Index phrase generation:

The goal of indexing is to represent the semantic concepts of an item in the information system to support finding relevant information. Single words have conceptual context, but frequently they are too general to help the user find the desired information.

Term phrases allow additional specification & focusing of the concept to provide better precision & reduce the user's overhead of selecting non-relevant items.

One of the earliest statistical approaches to determining term phrases using of a

COHESION factor b/w terms.

$$\text{COHESION} = \text{SIZE-FACTOR} * (\text{PAIR-FREQ}_{k,h} / \text{TOTF}_k * \text{TOTF}_h)$$

where SIZE-FACTOR is a normalization factor

based upon the size of the vocabulary &

PAIR-FREQ_{k,h} is the total frequency of co-occur-

- co-occurrence of the pair term_k, Term_h in the item collection. Co-occurrence may be defined in terms of adjacency, word proximity, sentence proximity etc.

This initial algorithm has been modified in the SMART system to be based on the following guidelines.

- » Any pair of adjacent non-stop words is a potential phrase.
- » Any pair must exist in 25 or more items.
- » Phrase weighting formula is a modified version of the SMART system single term algorithm.
- » Normalization is achieved by dividing by the length of the single term sub-vector.

The tagged text parser (TTP), based upon the linguistic string grammar, produces a regularized parse tree representation of each sentence reflecting the predicate - argument structure.

The tagged text parser contains over 400 grammar production rules. Some examples of

CLASS	EXAMPLES
determiners	a, the
singular nouns	paper, notation, language

CLASS	EXAMPLES
plural nouns	operations, data, processes
preposition	in, by, of, for
adjective	high, concurrent
present tense verb	presents, associates
present participle	multiprogramming.

To determine if a header-modifier pair warrants indexing, Stozalkowski calculate a value for Informational Contribution (IC) for each element in the pair. The basis behind the IC formula is a conditional probability b/w the terms. The formula for IC b/w x & terms (x_iy) is:

$$IC(x_i[x_iy]) = f_{x_iy} / N_x + D_{x_iy} + 1$$

where f_{x_iy} is the frequency of (x_iy) in the database, N_x is the num of pairs in which x occurs at the same position as in (x_iy) & D_{x_iy} is the dispersion parameter which is the num of distinct words with which x is paired.

When $IC=1$, x occurs only with the

$$y (f_{x_iy} = n_{x_iy} \text{ and } D_{x_iy} = 1).$$

following formula for weighting phrases:

$$\text{weight}(\text{Phrase}_i) = (C_1 * \log(\text{term f}) + C_2 * \omega(N, i)) * \text{IDF}$$

where $\omega(N, i)$ is 1 for $i < N$ and 0 otherwise

C_1, C_2 normalizing factors.

IDF is inverse document frequency.

(ii) Natural Language Processing (NLP):

only produces more accurate term phrases, but

can provide higher level semantic information

identifying relationships b/w concepts. System adds

the functional processes, Relationship concept Detection

&, conceptual Graph Generators & Conceptual Graph

matchers that generate higher level linguistic

relationships including semantic & coarse level

relationships.

During the 1st phase of this approach,

the processing tokens in the document are mapped

to Subject codes. These codes equate to index

term assignment & have some similarities to the

concept-based systems.

The next phase is called the Text Structure, which attempts to identify

general discou
The next level
blw the conc

weights to the
weight & are

information &

used in estal

4. Concept Ind

the terms

information

level concepts

most out w

a level fur

of terms

a reduced

number of

information

classes

general discourse level areas within an item.

The next level identifies interrelationships b/w the concepts.

The final step is to assign final weights to the established relationships. The weights are based upon a combination of statistical

information & values assigned to the actual words used in establishing the linkages.

4. Concept Indexing:

NLP starts with a basis of the terms within an item and extends the

information kept on an item to phrases & higher level concepts such as the relationships b/w concept

Concept indexing takes the abstraction a level further. Its goal is use concepts instead of terms as the basis for the index, producing a reduced dimension vector space.

Concept indexing can start with a number of unlabeled concept classes & let the information in the items define the concepts

classes created.

A term such as "automobile" could be associated with concepts such as vehicle, transportation, "mechanical device", "fuel" and "environment". The term "automobile" is strongly related to "vehicle", lesser to transportation and much lesser the other terms.

Thus a term in an item needs to be represented by many concept codes with different weights for a particular item. The basis behind the generation of the concept approach is a neural network model.

Special rules must be applied to create a new concept class. Example demonstrates how the process could work for the term "automobile".

TERM: automobile	Weights for associated concepts:
Vehicle	.65
Transportation	.60
Environment	.35
Fuel	.15
Mechanical Device	.15

Vector representation Automobile: { .65, ..., .60, ..., .35, .33, ..., .15 }

5. Hypertext Links

representation be generated retrieval viewed as a

Text of references as

al electronic

the items

linkages all

the linked

issue: How

locate related

mechanism

1. manually

were info

indexed

user nav

the hyp

5. Hypertext Linkages:

It's a new class of information representation is evolving on the Internet. Need to be generated manually creating an additional retrieval dimension. Traditionally the document was viewed as 2 dimensional.

Text of the item as one dimension and references as 2nd dimension.

Hypertext with its linkages to additional electronic items, can be viewed as n/w b/w the items that extend contents, ie by embedding linkages allows the user to go immediately to the linked item.

Issue: How to use this additional dimension to locate relevant information.

At the internet we have 3 classes of mechanism to help find information.

1. manually generated Indexes - ex: www.yahoo.com
where information sources on the home page are indexed manually into hyperlink hierarchy. The user navigates through hierarchy by expanding the hyper link. At some point the user starts

to see the end of items.

2. Automatically generated indexes - sites like lycos.com and altavista.com automatically go to other internet sites and return the text, google.com.
3. Web Crawler's : A web crawler (also known as a web spider or web robot) is a program or automated script which browses the world wide web in a methodical, automated manner.

Web crawlers (e.g. OpenText, Path Finder) are tools that allow a user to define items of interest & they automatically go to various sites on the Internet searching for the desired information.

What is needed is an index algorithm for items that include hypertext linkages as an extension of the concept where in the link exists w.r.t current concepts defined by the proximity of information.

Approaches in viewing the hyperlinks:

Approach 1 :

In this approach the hyperlink is viewed as an extension of the item or item in another dimension. It doesn't categorize rows & columns. The index values of hyperlinked item has

reduced weighted value from contiguous text
is biased by the type of linkage.

The weight of processing token could be.

$$\text{weight}_{i,j,k,l} = (\alpha * \text{weight}_{i,j} + \beta * \text{weight}_{k,l}) * (\gamma * \text{link}_{i,k})$$

where.

$\text{weight}_{i,j,k,l}$ → weight associated with token j in item i & token l in item k . (all these are related via hyperlink).

$\text{link}_{i,k}$ → weight related with strength of the link.

α, β, γ → weighting / normalization factors.

The link can be one-level link and it can be either strong or weak or it would be a multi-level transitive link.

The values can be stored in an expanded index structure or calculated dynamically if only the hyperlink relationships b/w items are available.

Approach 2:

→ In this approach, the system could

generate the hyperlinks b/w the items automatically.

→ Several attempts have been made to achieve this capability but they face problems with

Static vs dynamic growing databases & the efficiency needed for an operational environment is ignored.

- A new solⁿ on the basis of documents segmentation & clustering have been proposed by kelog & subhas.
- The links b/w the document pairs for each cluster are generated with the linkage at both the document & document sub-part level using cover coefficient based incremental Clustering method (~~cc²(m)~~) → This technique is referential similarity.
- The clustering & automatic link generation is performed in parallel.
- The candidates for hyperlinking are the item pairs in the same cluster & they should have the similarity above a given threshold.
- The process can be completed in 2 phases.
 - > The estimation of num of clusters is calculated
 - > In the first phase.
 - > The items are clustered & links are created in the 2nd phase.
- All of the link information is stored externally instead of storing it within an item or storing persistent link ID with in an item etc.

on demand,
when the mi
the 3 commo
> Mis-spelling &
> parser prob
> problems o
item was
→ The mode
→ This techn
referential
similarity
Document

- ## 6. Introduction
- has been
 - libraries
 - > clustered
 - the topic
 - > The de
 - of organa
 - are se

on demand, HTML items are created.

when the missed links are analyzed by the alg,
the 3 common problems discovered are,

- > misspellings & representations of multiple words.
 - > parser problems.
 - > problems occurred when definition of subparts of item was attempted.
- The more num of errors are from parsing.
- This technique has maximum effectiveness for referential links which naturally have higher similarity measures.

Document and Term Clustering.

6. Introduction to "clustering":

- > The concept of clustering has been in most usage as long as there are libraries.
- > clustering is used in the process of combining the topics on the same subject.
- > The definition of clustering could be the process of organizing objects into groups whose members are similar in some way to other.

- > The main goal of clustering would be in locating the information required.
- > As the clustering helps to locate the information, it needs to be indexed in organization of items in the libraries & it also need to follow the standards of electronic indexes.
- > The origin of clustering started with the "generation of thesauri".
- > The Thesaurus coming from latin word which means "treasure". The thesauri is similar to that of a dictionary.
- > The thesauri provides the synonyms & antonyms of the words instead of their definitions.
- > The main purpose of thesaurus is to help the authors in vocabulary selection.
- > The goal of clustering is to combine the similar objects into group.
- > There can be linkages b/w different clusters in clustering.

The process of clustering follows the following

steps: They are

Step 1: > The clustering effort should be defined

- with a d
- > The defini determin
- If thesea can be ta
- The determ should b clustering
- > with identify can find

Step 2:

- > The n determin be cluster
- > In th should
- > In c zones are us title, m
- > Here econo

with a domain.

- > The definition of a domain can be considered as determining the scope of the process.
- > If thesaurus is considered, the domain or scope can be taken as "medical terms", "Scientific terms" etc.
- > The determination of the collection of items that should be clustered can be treated as document clustering. This can be a part of complete database.
- > With the defining of a domain, it is easy to identify the objects in clustering process and it can reduce the erroneous data.

Step 2: To determine the attributes of the objects.

- > The next step after defining domains is to determine the attributes of the object that should be clustered.
- > In thesauruses, the specific words in the objects should be determined.
- > In case of document clustering, the specific zones with in the item are focused. These zones are used to determine the similarity such as title, main body but not references.
- > Here the main purpose is to reduce the erroneous relationships.

Step 3: To determine the strength of relationships b/w attributes.

- > The next step is to find the relationships b/w the attributes of the object.
- > This would help in deciding whether the object should be in this cluster or other cluster.
- > In this step, this would help in determining the words which are synonyms & the strength of their relationships.
- > In document clustering, this step may define a function known as "similarity function". This function is defined on the basis of co-occurrences of words which would determine the similarity b/w 2 items.

Step 4: To determine the total set of objects & their relationships.

- > The total set of object determination is the final step.
- > In this step, some kind of algorithm is applied to find the clusters & the items to be assigned to the clusters.

Guidelines for cluster characteristics:

There are some guidelines to be considered on the characteristics of the clusters. They are,

- > Cluster size
- > For each cluster there is a defined size.
- > There is a relationship to the defined size. Sometimes it is to reduce with the size.
- > For example, class can have different sizes that it can have. But the size of one class is another.

2) There

the size

> One

The size

size is

> If size

object

> So on

the size

3) The

object

relationships

is blue

object

determining
strength

define a
un in
words
below 2 items

as & their

the final

applied to
signed to

here are
characteristi-

1) Cluster should have well defined Semantic definition:

> for each & every cluster, there should be a well defined Semantic definition.

> There is a chance that the name that is assigned to the definition of the cluster leads to confusion sometimes.

> To reduce the confusion, the clusters are named with the numbers in some systems.

> for example, if some items are clustered into a class called "computers" can lead to a confusion that it also contains the items on main memory. But the items on main memory are clustered into another class known as "hardware".

2) There should be same order of magnitude for the size of class & cluster:

> One of the main purpose of cluster is expansion. The queries can be expanded to the set of items retrieved can also be expanded.

> If suppose 95% of the cluster is filled with the objects, then it can't be used for expansion.

> So order of magnitude should be maintained for the cluster size.

3) There should be no Domination Among the objects of the clusters: One object in the

cluster should not dominate the other object
in the cluster.

Example:

- > Assume a thesaurus class called "computer" contains the objects called "microprocessor", "286 processor", "386 - processor", "i3 processor" & "pentium".
 - > If the term "microprocessor" is found 94% & other one found 2% each.
 - > Then, if we use "microprocessor" as a synonym for "286 processor" then it would result in many errors.
 - > It would be better if we do not replace the place of "microprocessor".
- 4) The object can be assigned to multiple classes
- It is not:
- > The object once created should be only in the same cluster. Can it be assigned to other multiple clusters?
 - > This question would be treated as a trade off based on specificity & partitioning capability of the object semantics.
 - > If there is some ambiguity of the language, then the objects can be assigned to multiple classes instead of restricting it to single cluster.
 - > But this would increase the complexity in cluster creation & cluster maintenance.

The are additional important decisions associated with the generation of thesauri that are not part of item clustering.

Important Decisions in Thesauri Generation:

They are,

- 1) Word coordination approach
- 2) Word relationships
- 3) Homograph Resolution
- 4) Vocabulary Constraints.

1) Word coordination Approach:

- > This approach is used to specify the phrases & individual terms that are to be clustered.
- > The process of creating the linkages at the time of creating index is known as pre coordination.
- > If the coordination occurs at the time of search, it is known as Post Coordination.

2) Word Relationships:

- > The thesaurus generation would include a lot of relationships b/w the words that are possible. Three types of relationships are specified by Aitchison & Gilchrist. They are,

- a) Equivalence
- b) Hierarchical
- c) Non hierarchical

a) Equivalence Relationship:

- > The relationships that represents synonyms are known as equivalence relationships.
 - > They are most common in nature.
 - > In the creation of thesaurus, the words with significant overlapping but little differences are allowed.
- Example: > The words "points" and "photographs" are taken as synonyms even though the "points" can also include lithography.

> The words that are having same role are also considered as synonyms. They need not have same meaning necessarily.

> Example: The words "genius" and "mason" can be synonyms with respect to the class of "intellectual capability".

b) Hierarchical Relationship:

- > A technique where the class is general term & the entries are the examples of the general term. The relationship b/w the general & specific (example) can be treated as hierarchical relationship.

Example: The relationship b/w "company" & "employee" can be hierarchical.

c) Non Hierarc

- > The other can be, part synonymy and
- > The 2 words can be treated as same sense
- > The other one child one contrast

3) Homographs

- > If the meanings, t

Example:

field of

- > The elin unique

> The sys hz are & the

> When + help of

c) Non Hierarchical:

- > The other relationships that are non hierarchical can be, part wholes, collocation, paradigmatic, Taxonomy, Synonymy and Antonymy.
- > The 2 words related to each other with proximity closeness can be treated as collocation.
- > If the words are related on the basis of same semantics, then they are paradigmatic.
- > The other relationships included in semantic n/w are child of, parent of, part of All these words are contrasted.

3) Homograph Resolution:

- > If the word is having multiple different meanings, then it is called as "homograph".
Example: The term field can mean either "electric field" or "magnetic field".
- > The elimination of homographs by providing a unique meaning would be difficult.
- > The system would be in such a way that homographs are allowed by providing a unique meaning & the user selects the meaning that he desires.
- > When the user enters multiple terms, with the help of those terms, the correct meaning of the

homograph can be identified.

4) Vocabulary constraints:

- > The guidelines on normalization & specify that is based on vocabulary are included in the vocabulary constraints.
- > The constraints based on normalization can be stems vs complete words.
- > The specificity would eliminate the specific terms and it uses more general terms for class identifiers.

5. Thesaurus Generation:

- > From hundreds of years, the generation of clusters manually focused mainly on generating thesaurus.
- > But now a days, all the items are available in electronic form. Due to this the automated statistical clustering techniques are made available.
- > The thesauri which is automatically generated would contain the classes that has the usage of words in the corpora ie., database.
- > The classes of clusters does not have any name but are just a group of similar terms.

- > If the class then it requires
- > The other with the exception computation classes.

Basic methods

a thesaurus

i, Hand c

ii, co-occ

iii, Header

ii, Hand c

> The t

help in

domain a

> The ge

because

meanings

ii, co-occ

> There

co-oc

> The t

related

- > If the clusters are to be generated optimally, then it requires intensive computation.
- > The other techniques of clustering which starts with the existing clusters can reduce the computation but may not generate the optimum classes.

Basic methods of Generating a Thesaurus:

There are 3 basic methods for generating a thesaurus. They are based on,

- i) Hand crafted
 - ii) co-occurrence
 - iii) Header modified.
- (i) Hand crafted:
- > The thesauri that is manually generated would help in expanding the query if the thesauri is domain specific.
 - > The general thesaurus would not help much because of one word having many different meanings.
- (ii) Co-occurrence:
- > There are several techniques based on co-occurrence.
 - > The thesaurus is generated based on the related or associated words.

(iii) Head modified:

- > The relationships b/w the terms in the thesaurus are found based on the linguistic relationships.
- > The words that are similar in grammatical contexts are considered to be similar.
- > The syntactical structures that are discovered by the linguistic parsing are:
 - i, Subject verb
 - ii, Verb object
 - iii, Adjective noun
 - iv, Noun noun.
- > A mutual value is calculated for each noun by using a log function & the noun would contain a set of verbs, adjectives and nouns.
- > Using the mutual information, a final similarity b/w the words is calculated.

Manual clustering:

The manual clustering process would follow the steps that are described in as steps in the process of clustering.

- i) The 1st step in the process would be the determination of domain for clustering.
 - > The determination of domains would help in reducing the ambiguities that can be caused by homographs & others.

2, The starting generation of in the new thesaurus. In from items + are used.

- > A concord in alphabetical references found.

3, The art lie in the included q.

- > core sha be includ
 - undela
 - The wo values

> The other useful the oth

a) Kwo

> The

- the
thesau,
ships,
ical
vered
- oun by
ain a
plarity
- process
in as
- the
- ip in
caused
- 2, The starting points for the words for the generation of set of words that are to be included in the new thesaurus are taken from the existing thesauri. In the existing thesauri, the concordance from items that covers the domain & dictionaries are used.
- > A concordance would be the listing of words in alphabetical order along with their frequency & references of the items in which these words are found.
- 3, The art of construction of manual thesaurus would lie in the selection of the words that should be included in it.
- > care should be taken that some words should not be included such as,
 - unrelated to the domain of thesaurus.
 - The words with very high frequency & no informant values.
 - > The other tools may be helpful in determining the useful words if a concordance is used.
the other tools include KWOC, KWIC & KWAC.
- a) KWOC: Key word out of Context
- > The another term for concordance is KWOC.

b) KWIC: keyword in context

- > The possible term in the context of phrase is displayed by KWIC.
- > KWIC is structured in a way that it can easily identify the location of the term that is considered in a sentence.

c) KWAC: key word and context

- > KWAC would display the keywords on the basis of their context.

Example:

KWOC

TERM

FREQ

ITEM IDs

chip

2

doc2, doc4

computer

3

doc1, doc4, doc10

design

1

doc4

memory

4

doc3, doc4, doc8, doc12

KWIC

chip / computer design contains memory

computer / design contains memory chip /

design / contains memory chip / computer

chip / computer design contains

KWAC

chips

computer

design

memory

fig: example

> In the

KWIC to in

> The KWIC

the meaning

> In both

editor of

fragment

meaning

> The kwo

resolving

4, After the

clustered

lines

> With

human

KWAC

chips computer design contains memory chips
computer computer design contains memory chips
design computer design contains memory chips
memory computer design contains memory chips

fig: examples of KWOC, KWIC and KWAC.

- > In the above fig. character "I" is used in KWIC to indicate the end of the phrase.
- > The KWIC and KWAC are useful in determining the meaning of homographs.
- > In both the KWIC and KWAC displays, the editor of the thesaurus can read the sentence fragment associated with the term and its meaning is determined.
- > The KWOC does not provide any data in resolving this ambiguity.
- 4. After the selection of words, the words should be clustered on the basis of the relationship guidelines & interpretation of the relationship strength.
 - > With the help of human judgement & human analysis, this step would also be a

part of art of manual creation of thesaurus

- > The thesaurus finally would undergo several quality assurance reviews and the guidelines & then it is finalized.

Automatic Term clustering:

- > There are many techniques for the automatic generation of term clusters to create statistical thesauri.
- > The basis of all those techniques is that if the terms co-occur in the same item frequently then the terms are about the same concept.
- > The terms would be differed by the completeness by which the terms are related.
- > If the correlation is more complete then the time & computational overhead would also be higher.
- > The relationship b/w all the combinations of n unique words with an overhead of $O(n^2)$ is computed by the most complete process.
- > All the other techniques would be initiated by the arbitrary set of clusters & the process would be repeated on the assignment of terms to these clusters.

> A hierarchical clustering by making point.

steps for

of clustering methods:

automatic term

(i) complete

(ii) clusters

(iii) One pa

i. Complete

clustering

& simil

uses ve

> The v

- this m

> The vect

a matrix

> The no

> A hierarchy can be created when the num of clusters generated are very large. It can be done by making the initial cluster as the starting point.

Steps for processing:
Refer steps in the process of clustering.

Methods:

The different methods used for automatic term clustering are,

- (i) Complete Term Relation method.
- (ii) clustering using existing clusters.
- (iii) One pass Assignment.

i. Complete Term Relation method:

> In this method of

clustering, the clusters are formed on the basis of similarity b/w every pair of terms. This uses vector model for understanding.

- > The vector model can be used for understanding this method easily.
- > The vector model is indicated in the form of a matrix.
- > The rows of the matrix represents individual

items & the columns of the matrix represents unique words of processing tokens.

- > The values of the matrix would represents how strongly a particular word would represent the concepts in the item.
- > The example of a database is shown in below table.

	Term 1	2	3	4	5	6	7	8
Item 1	0	4	0	0	0	2	1	3
2	3	1	4	3	1	2	0	1
3	3	0	0	0	3	0	3	0
4	0	1	0	3	0	0	2	0
5	2	2	2	3	1	4	0	2

(Table 1) fig: Example of vector

- > There should be a measure to calculate the similarity b/w 2 terms.
- > The simple measure could be,

$$SIM(Term_i, Term_j) = \sum (Term_{k,i})(Term_{k,j})$$

where k is summed across set of all items.

- > The effect of the formula would be a column of 2 terms are being analyzed by multiplying & accumulating the values in each row.

- > The rows matrix.
- > The term in the
- > As the matrix
- > other non-
- > For the matrix

1	0
2	0
3	0
4	0
5	0
6	0
7	0
8	0

- > The results would be stored in $m \times m$ square matrix. That square matrix is known as Term Term matrix.
- > The value of m would be the num of columns in the original matrix.
- > As the formula is reflexive, the generated matrix would be "Symmetric".
- > other similarity formulas could produce a non-Symmetric matrix.
- > For the data in table 1 the term term matrix would be as shown in below table.

	Term 1	2	3	4	5	6	7	8
Item 1	7	16	15	14	14	9	7	
2	7	8	12	3	18	6	17	
3	16	8						
4	15	12	18					
5	14	13	6	6	6	18	6	9
6	14	18	16	18	6	1	12	16
7	9	6	0	6	9	2	0	3
8	7	17	8	9	3	16	3	

Table 2 : Term Term Matrix.

- > In the above table 2 the diagonal of matrix is empty as it represents the correlation of a word to itself.
- > Now, the threshold value should be selected to check whether the 2 values belong to the same cluster or not.
- > Suppose if the threshold value is 10 for the table 2, the terms with the value of 10 or greater are said to be similar. The resultant would be "Term relationship matrix" as shown in table 3.

Item	Term 1	2	3	4	5	6	7	8
Item 1	0	1	1	1	1	0	0	0
2	0	0	1	0	1	0	1	
3	1	0	1	0	1	0	0	
4	1	1	1	0	1	0	0	
5	1	0	0	0	0	0	0	
6	1	1	1	1	0	1	0	1
7	0	0	0	0	0	0	0	
8	0	1	0	0	0	1	0	

[Table] fig : Term Relationship Matrix.

- > In the above table the value 1 in term specified is the same class.
- > The term & clusters as the term &
- > The last step is to find the clusters in the same cluster.
- > There are several objects. They
 - 1) cliques
 - 2) stars
 - 3) complete graphs
- > Using any of them can be solved.
- > Thus clustering is done.
- 1) cliques method
 - in an cluster
 - sort all other
 - It can be
 - 1, let $i=1$
 - 2, Select the
 - 3, start with

In the above table represents a binary matrix. The value 1 in the matrix indicates that the term specified by the column & row can be in the same class.

The term f can't be included in any of the clusters as there are no similar terms for term f.

The last step in the process would be to determine the clusters when two objects are in the same cluster.

There are several methods to determine such objects. They are,

1) cliques 2) single link

3, start 4, connected components

Using any of the above methods, the ambiguity can be solved.

Thus clusters are created.

1) Cliques method:

In cliques method, all items in a cluster should be within the threshold of all other items.

It can be explained as,

1, let $i=1$

2, Select term i and place it in new class.

3, Start with term k where $k = i+1$ and $\sigma = i+1$.

4. Validate if term_k is within the threshold of all terms within the current class. If not then
 5. else let $k = k + 1$
6. if $k > m$ (num of words)
 then $\delta = \delta + 1$
 if $\delta = m$ then go to 7
 else create a new class with term_i in it
 go to 4
 else goto 4.
7. If the current class has only term_i in it &
 there are other classes with term_i in them then
 delete current class
 else
 $i = i + 1$
8. If $i = m + 1$ then goto 9.
 else goto 2
9. Eliminate any classes that duplicate & that are subsets of other classes.
- > By applying the above algorithm to the table 3. the clusters that would be created are, $\delta = 1$ and $i = 1$ which most often treat

cluster 1

cluster 2

cluster 3

cluster 4

> The team

One cluster

2) Single

Similar to

added to

Algorithm

Step 1:

Step 2: All

class

Step 3: fo

elations

Step 4: w

> A team

clusters

> Applying

can be

- shold
 alor soft
 is most
 ngs soft
 of soft
 eadsoft
 -F most
 ed soft
 it soft
 cts
 soft
 -soft
 in it & G
 ffects them then
 a priso
 ed nis
 o soft
 compo
 b no si
 g that
 us to
 to the
 created
 ffects
 shold
 alor soft
 is most
 ngs soft
 of soft
 eadsoft
 -F most
 ed soft
 it soft
 cts
 soft
 -soft
 in it & G
 ffects them then
 a priso
 ed nis
 o soft
 compo
 b no si
 g that
 us to
 to the
 created
 ffects
- | | | |
|-----------|--|----------|
| Cluster 1 | $\{ \text{term} 1, \text{term} 3, 4, 6 \}$ | clusters |
| Cluster 2 | $\{ \text{term} 1, 5 \}$ | clusters |
| Cluster 3 | $\{ \text{term} 2, 4, 6 \}$ | clusters |
| Cluster 4 | $\{ \text{term} 2, 6, 8 \}$ | clusters |

The term 1 & term 6 are present in more than one cluster.

2) Single Link Method:

In this method, any term

similar to any term in the cluster can be added to the cluster.

Algorithm

Step 1: Identify the item that is not in the class & place it in new class.

Step 2: All the other terms related to new class must be placed in it.

Step 3: for every term that is entering into the class, perform step 2.

Step 4: When there are no new terms in step 2, goto step 1.

- > A term can't exist in 2 or multiple different clusters.
- > Applying this algorithm, the matrix in tables can be divided as,

Cluster 1 {term 1, 2, 3, 4, 5, 6, 8}

1st cluster

Cluster 2 {term 7}

2nd cluster

3. Star method:

> In this technique, a term is selected & it is placed in a class & all terms related to the term are placed in that class.

> The terms that are not placed in the class are selected & they are placed into classes with the same criteria.

> There can be many different classes created based on star technique.

The clusters of table 3 can be,

Cluster 1 {term 1, term 3, term 4, term 5, term 6}

Cluster 2 {term 2, term 7, term 8}

Cluster 3 {term 6}

> The terms can be in multiple clusters in this method.

This can be eliminated by expanding the constraints to exclude any term that has already been selected from a previous cluster.

4. String method:

> In this method, it starts with the term & the term is included in a cluster & the additional term is added if it is similar to other term if that is not present in the cluster already.

> The new process will

terms to be

> The new to
clusters

> The clusters
in table 3

Cluster 1

Cluster 2

Cluster 3

N/W Diagnosis

technique

clusters

> Each term

the nodes

spindle

to reduce

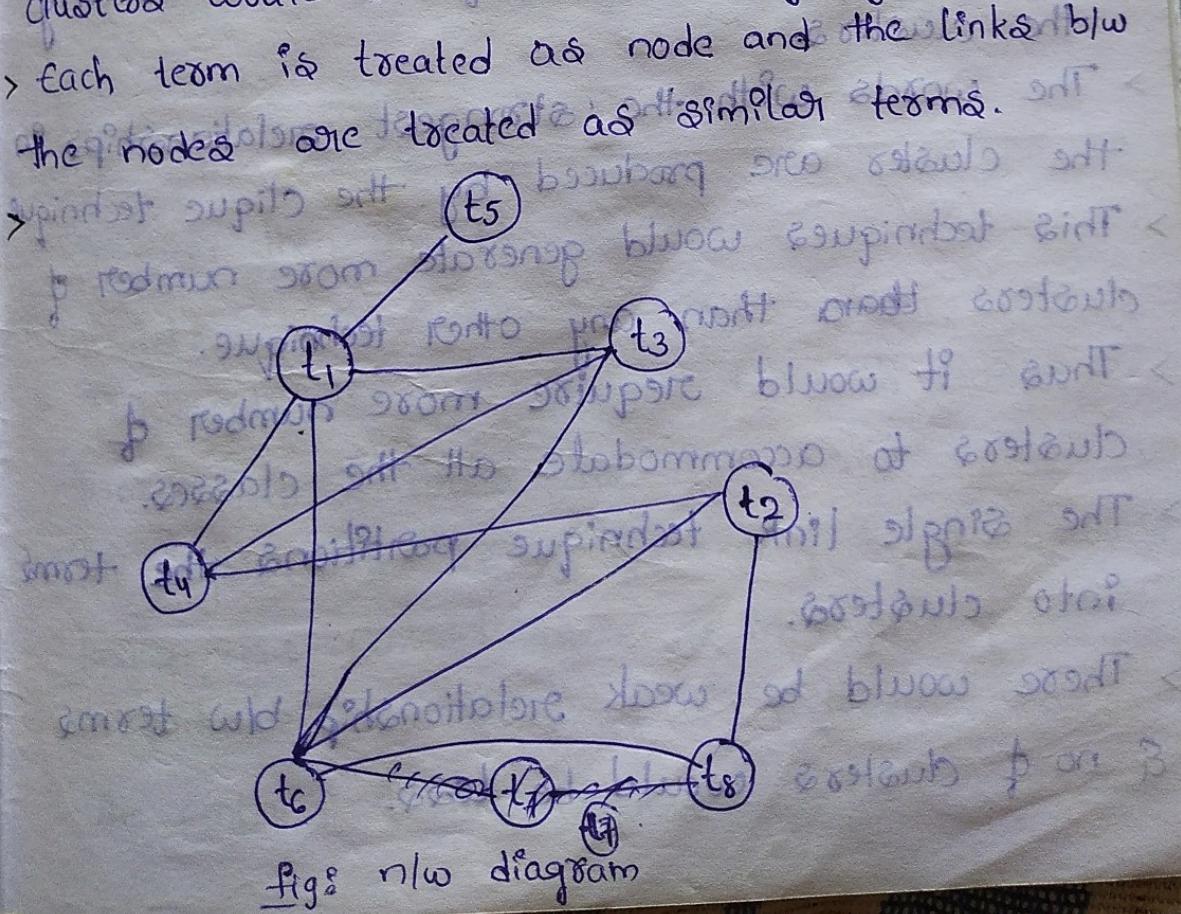
b

smart

44

smart

- > The new term is treated as a new node & the process will be repeated till there are no new terms to be added to the clusters.
 - > The new term should not exist in any of the clusters already created.
 - > The clusters that are created for the matrix in table 3 would be,
- Cluster 1 {term 1, term 3, term 4, term 2, term 6, term 8}
- Cluster 2 {term 5}
- Cluster 3 {term 7}
- N/w Diagram of the terms
- > The n/w diagram is a technique that validate the clusters generated by the clusters would be understood easily by this.



- > To identify cliques, the sub nets are selected where all the items are connected by the links.
 - > From table 2, term 7 (t_7) is in a class and t_5 is in cluster without blocks most used off.
 - > The focus should be on triangles & polygons of 4 sides with diagonals.
 - > To find all the clusters for an item, there is need to find all the subnetworks & there should be maximum num of nodes in each sub n/w.
 - > For t_1 , it is in sub n/w ($\{t_1, t_3, t_4, t_6\}$).
 - > The term t_2 has 2 sub n/w, one is $\{t_2, t_4, t_6\}$ & other is $\{t_2, t_5, t_8\}$.
 - > The n/w diagram is a simple visual tool where there would be less number of nodes to identify the clusters.
 - > The words with the strongest relationships in the cluster are produced by the clique technique.
 - > This techniques would generate more number of clusters than any other technique.
 - > Thus it would require more number of clusters to accommodate all the classes.
 - > The single link technique partitions the terms into clusters.
 - > There would be weak relationship b/w terms & no of clusters would be less.
- These can if the value
cliques prove
- The single link recall so
- It requires
- ## Clustering
- > One mode with the same clusters.
 - > As the process there is no the clusters.
 - > The terms are revised the clusters.
 - > When the n/w clusters, the
- Centroid:**
- > The centroid in order to
 - > In physics, known as centroid is -
 - > The centroid

- > There can be 2 terms in the same cluster if the value of similarity is zero.
- > Cliques provide highest precision.
- > The single link process would maximize the recall.
- > It requires the overhead of $O(n^2)$ comparisons.

Clustering Using Existing Clusters:

- > One more method of creating clusters is to start with the set of existing clusters.
- > As the process starts with already existing clusters, there is no need to calculate the similar terms in the clusters.
- > The terms that are assigned initially to the clusters are revised by validating each & every term in the clusters again.
- > When the minimal movement is detected b/w the clusters, then the process would be stopped.

Centroid:

- > The centroids are calculated for each cluster in order to minimize the calculations.
- > In physics, the centre of mass of set of objects is known as Centroid where as in context of vectors, centroid is the average of all vectors.
- > The centroid of the clusters can be viewed as

- another point in n dimensional space where n is the number of items.
- > The centroids of the initial clusters are not same as the final clusters.
 - > The relation of similarity b/w all the terms that are existing and the centroids can be calculated.
 - > The term is moved to the cluster with highest similarity. The process will be continued till the clusters stabilize.
 - > The order of the calculation in this process is $O(n)$.
 - > As the process is iterative & the terms change their clusters till it stabilizes, the initial assignment of terms to clusters is not that important.

Graphical Representation:

The graphical representation of centroids is shown as below

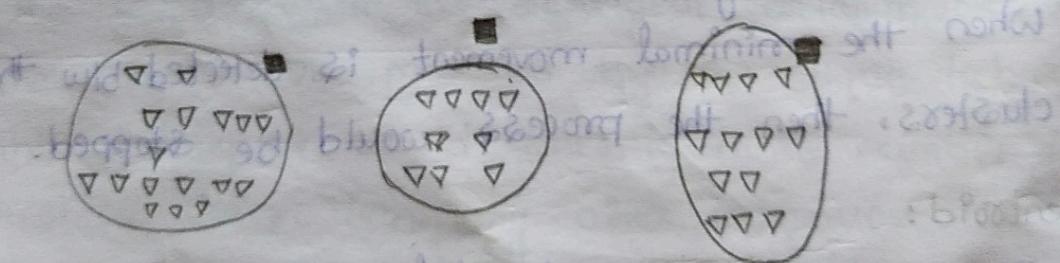


fig (a) initial Centroids for clusters.

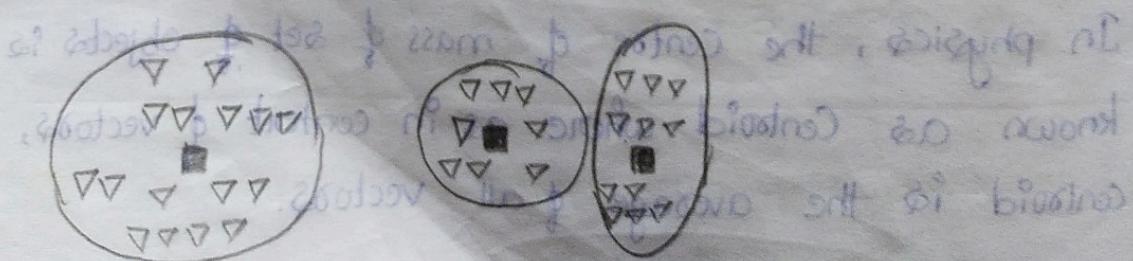


fig (b) Centroids after reassessing terms.

- > fig shows initial as
- > The solid cluster.
- > In fig (a) clusters c assignmen
- > After p E new one almo

Item	1	2	3	4	5
cluster	2	3	3	4	5

It is

cluster

Cluster

Cluster

The c

cluster

- > fig shows how the clusters are moved after initial assignments.
- > The solid box represents the centroid for each cluster.
- > In fig (a) represents the centroids of the initial clusters and the ovals represents the ideal cluster assignments.
- > After performing the re-assignment, the new centroid & new clusters are represented in fig (b). They are almost similar to the ideal clusters.

Example:

	Term1	T2	T3	T4	T5	T6	T7	T8
Item1	0	4	0	0	0	2	1	3
Item2	3	1	4	3	6	2	0	1
Item3	3	0	0	0	3	0	3	0
Item4	0	1	0	3	0	0	2	0
Item5	2	2	2	3	1	4	0	2

table: vector

It is assumed that the clusters formed are,

$$\text{Cluster 1} = \{\text{Term 1, Term 2}\}$$

$$\text{Cluster 2} = \{\text{Term 3, Term 4}\}$$

$$\text{Cluster 3} = \{\text{Term 5, Term 6}\}$$

The centroids for each cluster would be

$$\begin{aligned} \text{Cluster 1} &= \left(\frac{0+4}{2}, \frac{3+1}{2}, \frac{3+0}{2}, \frac{0+1}{2}, \frac{2+2}{2} \right) \\ &= 4\frac{1}{2}, 4\frac{1}{2}, 3\frac{1}{2}, 1\frac{1}{2}, 4\frac{1}{2}. \end{aligned}$$

$$\text{Cluster 2} = \frac{(0+0)}{2}, \frac{(4+3)}{2}, \frac{(0+0)}{2}, \frac{(0+3)}{2}, \frac{(2+3)}{2}$$

$$= \frac{0}{2}, \frac{7}{2}, \frac{0}{2}, \frac{3}{2}, \frac{5}{2}$$

$$\text{Cluster 3} = \frac{(0+2)}{2}, \frac{(1+2)}{2}, \frac{(3+0)}{2}, \frac{(0+0)}{2}, \frac{(1+4)}{2}$$

$$= \frac{2}{2}, \frac{3}{2}, \frac{3}{2}, \frac{0}{2}, \frac{5}{2}$$

> The value of centroid in each cluster is the avg of the weights of the terms in the cluster.

> The measure of similarity applied here is

$$\text{SIM}(\text{term}_p, \text{term}_q) = \sum (\text{term}_{k,p}) (\text{term}_{k,q})$$

This method is applied b/w 18 terms & 3 centroids & the result would be as shown in below

table.

	Term 1	2	3	4	5	6	7	8
cluster 1	$\frac{29}{2}$	$\frac{29}{2}$	$\frac{21}{2}$	$\frac{27}{2}$	$\frac{17}{2}$	$\frac{32}{2}$	$\frac{15}{2}$	$\frac{24}{2}$
2	$\frac{31}{2}$	$\frac{20}{2}$	$\frac{38}{2}$	$\frac{45}{2}$	$\frac{12}{2}$	$\frac{31}{2}$	$\frac{6}{2}$	$\frac{17}{2}$
3	$\frac{28}{2}$	$\frac{21}{2}$	$\frac{22}{2}$	$\frac{24}{2}$	$\frac{17}{2}$	$\frac{30}{2}$	$\frac{11}{2}$	$\frac{19}{2}$
Assign cluster	2	2	2	2	3	2	1	1

table 5

Fig: Iterated Assignments of clusters.

- > In case of similarity
- > This tie weights of
- > One that
- > Most of term is
- > In the

cluster 1

cluster 2

cluster 3

> The components

cluster

privately

2

3

Assign

cluster

> By

term

> The

Step

- > In case of term 5, the cluster 1 & cluster 3 have same similarity. & clusters can be assigned at ~~dearable~~ ~~dearable~~
 - > This tie can be broken by considering the similarity weights of other items in the cluster & assign the one that is having most similar weights.
 - > Most of the terms in cluster 1 are above $\frac{1}{2}$. So the term is assigned to cluster 3.
 - > In the same way the next iteration would be
- Iteration
- $$\text{cluster 1} = \left\{ \frac{8}{3}, \frac{2}{3}, \frac{3}{3}, \frac{3}{3}, \frac{4}{3} \right\} \quad \text{class 1} = \{1, 17, 18\}$$
- $$\text{cluster 2} = \left\{ \frac{2}{4}, \frac{12}{4}, \frac{3}{4}, \frac{3}{4}, \frac{1}{4} \right\} \quad \text{class 2} = \{1, 3, 4, 6\}$$
- $$\text{cluster 3} = \left\{ \frac{0}{1}, \frac{1}{1}, \frac{3}{1}, \frac{0}{1}, \frac{1}{1} \right\} \quad \text{class 3} = \{5\}$$
- > The above are new centroids & the cluster assignments would be as shown in table below.

	Term 1	2	3	4	5	6	7	8
cluster 1	$\frac{23}{3}$	$\frac{45}{3}$	$\frac{16}{3}$	$\frac{27}{3}$	$\frac{15}{3}$	$\frac{36}{3}$	$\frac{23}{3}$	$\frac{34}{3}$
cluster 2	$\frac{67}{4}$	$\frac{45}{4}$	$\frac{70}{4}$	$\frac{78}{4}$	$\frac{33}{4}$	$\frac{72}{4}$	$\frac{17}{4}$	$\frac{40}{4}$
cluster 3	$\frac{12}{1}$	$\frac{3}{1}$	$\frac{6}{1}$	$\frac{6}{1}$	$\frac{11}{1}$	$\frac{6}{1}$	$\frac{9}{1}$	$\frac{3}{1}$
Assign clusters	2	1	2	2	3	2	3	1

table 6: Cluster Assignment

- > By comparing table 5 & 6, the only change is term 1 is changed from cluster 1 to cluster 3.
- > The reason would be the term 7 is not strongly related to cluster 1.

Limitation:

- > Although there are less calculations in this method when compared to Complete Term Relationship method, this method has few limitations.
- > The number of clusters were defined at the starting of the process & they can't be expanded or increased.
- > It is also possible that there would be less number of clusters at the end of the process.
- > even if the similarity is relatively weaker, the term would be assigned to that cluster because all the items should be assigned to some other clusters.

One Pass Assignment:

- > This technique of clustering has minimum overhead because one pass of all the terms is used to assign the terms to the clusters.
- > In this technique, the 1st term is assigned to the 1st cluster.
- > The additional terms are compared with the centroids of the existing clusters. A threshold value is chosen & if the item has greater weight than the threshold, then it is assigned with the highest similarity one.

> For other clusters
of centroid
> If the similarity is less than the item in a

Example:

- > Consider a table.

Item	Term
1	9
2	7
3	16
4	15
5	14
6	14
7	9
8	7

the cluster

Cluster 1

Cluster 2

Cluster 3

Cluster 4

- > for the cluster that is modified, a new value of centroid have to be found.
 - > if the similarity of all the centroids existing is less than that of threshold, then item is the 1st item in a new class.
- Example:

> Consider the matrix that is represented in below table. & let the threshold value be 10, then

Item	Term 1	2	3	4	5	6	7	8
Term 1	7	16	15	14	14	9	7	
Term 2	7	8	12	13	18	6	17	
Term 3	16	8		18	6	16	0	8
Term 4	15	12	18		6	18	6	9
Term 5	14	3	6	6		6	9	3
Term 6	14	18	16	18	6		2	16
Term 7	9	6	0	6	9	2		3
Term 8	7	17	18	9	3	16	3	

table : term term matrix most 3
columns terms are 3

the clusters are,

Cluster 1 = {term 1, term 3, term 4}

Cluster 2 = {term 2, term 6, term 8}

Cluster 3 = {term 5}

Cluster 4 = {term 7}

> During one pass process, the centroids values used are,

$$\text{cluster } 1 \{ \text{term} 1, 3 \} = 0, \frac{7}{2}, \frac{3}{2}, 0, \frac{4}{2}$$

$$\text{cluster } 2 \{ \text{term} 1, 3, 4 \} = 0, \frac{10}{3}, \frac{3}{3}, \frac{3}{3}, \frac{7}{3}$$

$$\text{cluster } 2 \{ \text{term} 2, 6 \} = \frac{6}{2}, \frac{3}{2}, \frac{0}{2}, \frac{1}{2}, \frac{6}{2}$$

Limitations:

- > Even though, this process has less computational overhead of order $O(n)$, it can produce optimal results i.e., optimum clustered classes of clusters.
- > If the order of the items in analysis are changed then different clusters would be produced.
- > Due to the averaging nature of centroids, the items in same cluster would be appeared as they are in different clusters.

Item clustering:

- > In the sawpi generation, the item clustering & term clustering are almost similar.
- > Any library or filing system uses manual item clustering.
- > In manual clustering, one person would read the item & determine the category into which the item belongs.

> The item is categorized.
> But with would be categories assigned.
> Because automatically developed
> the machine be used
> The machine, complies
(i), clustering
(ii), One

> The clustering on the basis of features

> The clustering based on the column

> The clustering based on the row

Six

- > The item is generally assigned to a single category in physical clustering.
- > But with the introduction of indexing, the item would be stored physically in the primary category & the item is found in all the other categories that are defined by index items assigned to it.
- > Because all the items are stored electronically, the automatic clustering techniques have also been developed.
- > The methods of automatic term clustering can also be used for clustering the items.
- > The methods include,
 - (i) Complete term Relation Method.
 - (ii) Clustering using existing clusters.
 - (iii) One Pass Assignments.
- > The similarity b/w the documents can be known on the basis of items having the common terms v/s terms with items in common.
- > The similarity funⁿ is performed b/w rows & columns of the item matrix.
- > The measure of similarity would be calculated by,

$$\text{SIM}(\text{item}_i, \text{item}_j) = \sum (\text{term}_{i,k})(\text{term}_{j,k})$$

Example: Consider the vector table,

	Term 1	2	3	4	5	6	7	8
Item 1	0	4	0	0	0	2	1	3
Item 2	3	1	4	3	1	2	0	1
Item 3	0	0	0	3	0	0	3	0
Item 4	0	1	0	3	0	0	2	1
Item 5	2	3	0	4	1	0	0	2

the value of k goes from 1 to 8, the item-item matrix would be as shown in below table.

	Item 1	2	3	4	5
Item 1	11	3	6	22	33
2	11	12	10	36	29
3	3	12	6	9	11
4	6	10	6	11	10
5	22	36	9	11	20

table: Item - item matrix

> If the threshold value is 10, then the item relationship matrix could be as shown in below table.

$$(A_{ij})_{(1,1)(1,2)} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} M_{12}$$

	Item 1	Item 2	Item 3	Item 4	Item 5
Item 1	1	0	0	1	0
2	1	0	1	0	1
3	0	1	0	0	0
4	0	1	0	0	1
5	1	1	0	1	1

table: Item Relationship matrix.

- > The clusters using cliques technique would be for above table.
- so clustered 1 = {item 1, item 2, item 5}
- cluster 2 = {item 2, item 3}
- cluster 3 = {item 2, 4, 5}
- > using single link technique, the clusters created would be,
- cluster 1 = {item 1, 2, 3, 4, 5}
- > The Star technique would generate the following clusters
- cluster 1 = {item 1, 2, 5}
- cluster 2 = {item 3, 2}
- cluster 3 = {item 4, 2, 5}
- > The string technique would generate
- cluster 1 = {item 1, 2, 3}
- cluster 2 = {item 4, 5}

- > The ambiguities & erroneous hits are caused by vocabulary domain homographs.
- > The process of clustering can also be started with the existing clusters.
- > If,

$$\text{cluster 1} = \{\text{item 1, item 3}\}$$

$$\text{cluster 2} = \{\text{item 2, item 4}\}$$

then centroids are,

$$\text{cluster 1} = \frac{3}{2}, \frac{4}{2}, \frac{0}{2}, \frac{0}{2}, \frac{3}{2}, \frac{2}{2}, \frac{4}{2}, \frac{3}{2}$$

$$\text{cluster 2} = \frac{3}{2}, \frac{2}{2}, \frac{4}{2}, \frac{6}{2}, \frac{1}{2}, \frac{2}{2}, \frac{1}{2}, \frac{1}{2}$$

- > After calculating again the results would be as in below table

	Item 1	2	3	4	5
cluster 1	$\frac{3}{2}$	$\frac{2}{2}$	$\frac{3}{2}$	$\frac{8}{2}$	$\frac{3}{2}$
cluster 2	$\frac{7}{2}$	$\frac{5}{2}$	$\frac{18}{2}$	$\frac{9}{2}$	$\frac{4}{2}$
Assign cluster	Storage	Know	Upwind	Not	Off

table: Clustering of items with existing clusters

- > Now the cluster 2 contains 4 items and cluster 1 only one item. So recalculating does not result in any re-assignment.

$$\{E, G, L, M\} = \text{Cluster 2}$$

$$\{A, H, M\} = \text{Cluster 1}$$

Hierarchy of Clusters:

- > In Information Retrieval, the hierarchical clustering concentrates mainly on the concept of hierarchical agglomerative clustering methods (HACM).
- > In agglomerative clustering, the process would start with unclustered items & then similarity measures are performed pair wise to determine the clusters.

objectives:

The objectives of hierarchy of clusters

would be,

- i) Reduction of search overhead.
- ii) visual representation of information space.
- iii) Relevant item retrieval expansion.

ii) Reduction of search overhead:

The overhead in

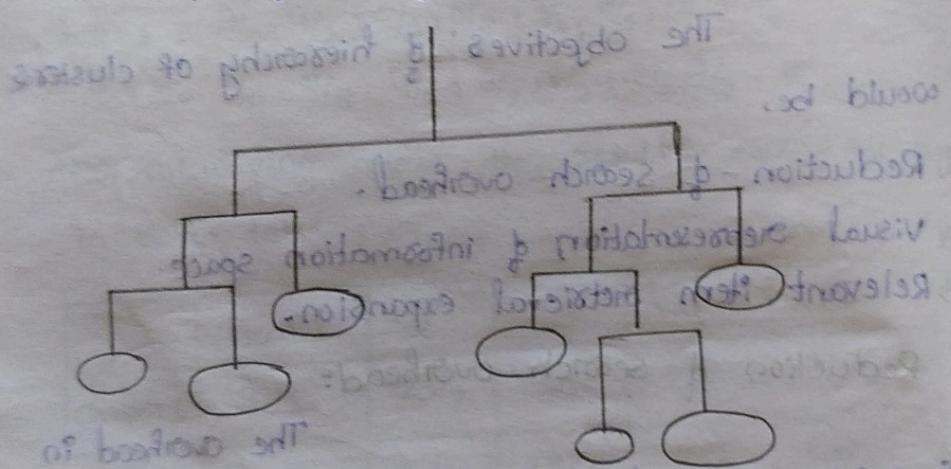
searching is reduced by performing top down searches of the centroid & of clusters present in the hierarchy & by trimming (cut) the branches of hierarchy that are not relevant.

iii) Visual representation of information space:

- > It would be very difficult task to display the complete item space visually.
- > The pictorial representations known as Dendograms would help the user to determine the clusters to

be reviewed & likely to find the relevant items.

- > The size of clusters is defined by size of ellipse in dendograms.
- > The linkage strengths are indicated by lines. Dashed lines indicated reduced similarity.
- > It allows the user to browse the alternate paths in database.
- > The example of dendrogram is as shown in below



- Relevant Items Retrieval:
- > The items that are relevant can be browsed by using logical hierarchy instead of visual display hierarchy.
 - > Once the relevant item is identified, the user can send the request to other items in the cluster.
 - > The specification of the item can be increased by passing through child clusters by increasing the generality by reviewing the parent clusters.