

MODULE III

Document and Term clustering

Introduction :-

Clustering :-

- ↳ Goal - to assist in the location of info.
- ↳ of words originated with the generation of thesauri.
- Thesaurus provides the synonyme and antonyme of words.
- ↳ Aim is to provide a grouping of similar objects into a class under a general title.
- ↳ Allows linkages b/w clusters to be specified.

The steps of clustering :-

- (a) Define the domain for the clustering effort:-
 - ↳ identify those objs to be used in the clustering process and reduce the potential for erroneous data that could induce error in the clustering process.
- (b) Once the domain is determined, determine the attributes of the obj to be clustered.
- (c) Determine the strength of the relationship b/w the attributes whose co-occurrence in objs suggest those objs should be in the same class.

(d) At this point, the total set of objs and the strengths of the relationships b/w the objs have been determined. The final step is applying some alg to determine the class to which each item will be assigned.

- Guidelines on the characteristics of classes :-
- * A well-defined semantic def should exist for each class.
 - * The size of the classes should be within the same order of magnitude.
 - * Within a class, one obj should not dominate the class.
 - * Whether an obj can be assigned to multiple classes or just one must be decided at creation time.

→ Decisions associated with the generation of thesauri and not part of item clustering :-

- * Word coordination approach:- specify if phrases as well as individual terms are to be clustered.
- * Word relationships:- when the generation of a thesaurus includes a human interface, a variety of relationships b/w words are possible.
 - ↳ equivalence
 - ↳ hierarchical
 - ↳ non-hierarchical.
- * Homograph resolution
- * Vocabulary constraints.

- I. Thesaurus Generation: -
- ← Manual Clustering
 - Automatic Term Clustering
 - (i) complete Term Relation Method.
 - (ii) clustering using existing clusters.
 - (iii) One Pass Assignment

Manual clustering:

→ 1st step is to determine the domain for the clustering.

Defining the domain assists in reducing ambiguities caused by homographs and helps focus the creator.

A concordance is an alphabetical listing of words from a set of items along with their frequency of occurrence and references of which items in which they are found.

→ Main focus in manual thesaurus construction is the selection of the set of words to include, where words that are not part of the domain are not included.

Also words that have less freq of occurrence are also not included.

→ A keyword out of context (KWOC) is another name of concordance.

keyword in context (KWIC) displays a possible term in its phrase context

→ used to identify easily the local

Key word and context (KWAC) displays the words followed by their context

Eg: "Computer design contains memory chips".

KWOC

Term	Freq	Item Ids
chips	2	doc 2, 4
computer	3	doc 1, 4, 10
design	1	doc 4
memory	3	doc 3, 4, 8, 12

KWIC

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

KWAC

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

KWIC & KWAC - determine the meaning of homographs.

Once the terms are selected they are clustered based upon the word relationship guidelines and the interpretation of the strength of the relationship.

Automatic Term clustering :-

This is based on the concept that the more frequently two terms co-occur in the same item, the more likely they are about the same concept.

The most complete process computes the strength of relationships b/w all the combinations of the 'n' unique words with an overhead of $O(n^2)$.

The simplest case employs one pass of the data in creation of the clusters. When the number of clusters created is very large, the initial clusters may be used as a starting point to generate more abstract clusters creating a hierarchy.

→ The basis for automatic generation of a thesaurus is a set of items that represents the vocabulary to be included in the thesaurus. The processing tokens in the set of ^{selected} items are the attributes to be used to create the clusters.

→ The ~~not~~ automated method of clustering does is based upon the polythetic clustering, where each cluster is defined by a set of words and phrases.

(i) Complete Term Relation Method :-

Here, the similarity b/w every term pair is calculated as a basis for determining the clusters, using the vector model.

The vector model is represented by a matrix where the rows are individual items and the cols are the unique words in the items.

The values in the matrix represent how strongly that a particular word represents the concept in the terms item.

Eg:	T1	T2	T3	T4	T5	$\text{Sim}(T_1, T_2) = \frac{(0+3+3)}{5} = \frac{6}{5}$
I1	0	4	0	0	0	$+ 3 \times 0 + 0 \times 0 +$
I2	3	1	4	3	1	2×2
I3	3	0	0	0	3	$= \frac{6}{5}$
I4	0	0	0	3	0	
I5	2	2	2	3	1	

To determine the relationship b/w terms, a similarity measure is used.

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

where "k" summed across the set of all items. These results can be placed in Term-Term matrix.

Using the Item-to-Term matrix, Term-Term matrix is produced, as shown below:

	Term1	2	3	4
Term1	No similarity from i to i	7	16	15
2	7		8	12
3	16	8		18
4	15	12	18	

$\text{Sim}(T_1, T_2) = \frac{1}{5}$

To obtain this value, consider the 2 cols of T_1 & T_2 from the above matrix & multiply them corresponding freq.

There are no values on the diagonal since that represents the auto correlation of a word to itself.

The next step is to select a threshold that determines if two terms are considered similar enough to each other to be in the same class. This produces a new binary matrix called Term-Relationship matrix, defining similar terms.

This is shown below:-

T1	2	3	4	5	
T1	0	1	1	1	1 - similar
2	1	0	1	0	0 - dissimilar
3	1	0			
4	1				
5	1				

The final step in creating clusters is to determine when two obj are in the same cluster from the above matrix.

The foll alg are the most common:-

- (a) cliques, (b) single link, stars and (c) connected components.

(a) cliques require all items in a cluster to be within the threshold of all other items.

Algorithm for cliques:-

1. Let $i = 1$
2. Select term_i and place it in a new class
3. Start with term_k where $i = k = i + 1$
4. Validate if term_k is within the threshold of all the terms within the current class.

5. If not, let $k = k+1$
6. If $k > m$ (# of words)
 - then $M = M+1$
 - If $M = m$ then goto 6 else
 - $K = M$
 - Create a new class with term $_i$ in it
 - goto 3
 - else goto 8
7. If current class only has term $_i$ in it and 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 or are subsets of other classes.

By applying the alg over the Term-Term matrix, it can be found that the same term can be found in multiple classes.

(b) In single link clustering, the strong constraint that every term in a class is similar to every other term is relaxed. The rule to generate single link cluster is that any term that is similar to any term in the cluster can be added to the cluster. Now, it becomes impossible for a term to be in two different clusters.

The algorithm for single-link is as follows:-

1. Select a term that is not in a class and place it in a new class.
2. Place in that class all other terms that are related to it.
3. For each term entered into the class, perform step 2.
4. When no new terms can be identified go to 1.

(c) The star technique selects a term and then places in the class all terms that are related to that term.

Terms not yet in classes are selected as new seeds until all terms are assigned to a class. This technique also allows terms to be in multiple clusters.

The string technique starts with a term and adds in the class one additional term that is similar to the term selected and not already in a class. The new term is then used as the new node and the process is repeated until no new terms can be added bcoz the term being analyzed does not have another term related to it or the terms related to it are already in the class.

(ii) Clustering Using Existing clusters:-

This methodology reduces the number of similarity calculations required to determine the clusters. The initial assignment of terms to the clusters is revised by revalidating every term assignment to a cluster.

The process stops when minimal movement b/w clusters is detected.

To minimize calculations, centroids are calculated for each cluster, which is the average of all of the vectors in a cluster.

The similarity b/w all existing terms and the centroid of the cluster can be calculated. The term is reallocated to the cluster that has the highest similarity.

Each value in the centroid is the avg of the wts of the terms in the cluster for each item in the database.

The problem is that the no. of classes is defined at the start of the process and can not grow. // ①

(iii) One Pass Assignments:-

Here, The 1st term is assigned to the 1st class. Each additional term is compared to the centroid of the existing classes. A threshold is chosen. If the item is greater than the threshold, it is assigned to the class with highest similarity.

A new centroid has to be calculated for the modified class. If the similarity to all of the existing centroids is less than the threshold, the term is the first item in a new class. This process continues until all items are assigned to classes.

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

From the above Term x Term matrix, with a threshold to the boy class would be generated :-

class 1 = Term 1, 3, 4

class 2 = Term 2, 6, 8

class 3 = Term 5

class 4 = Term 7

class 1 (Term 1, 3) = 0, 7/2, 3/2, 0, 4/2

class 1 (Term 1, 3, Term 4) = 0, 10/3, 3/3, 3/3, 7/3

class 2 (Term 2, Term 6) = 6/2, 3/2, 0/2, 4/2, 6/2

II. Item clustering:-

In this case, someone reads the item and determines the category or categories to which it belongs. When physical clustering occurs, each item is usually assigned to one category.

With the advent of indexing, an item is physically stored in a primary category, but it can be found in other categories as defined by the index terms assigned to the item.

Similarity b/w doc is based upon two items that have terms in common vs terms with items in common.

Thus, the similarity fn is performed b/w rows of the item matrix.

$$\text{SIM}(\text{Item}_i, \text{Item}_j) = \sum_{k=1}^8 (\text{Term}_{i,k}) (\text{Term}_{j,k})$$

\uparrow k terms

k - ranges from 1 to 8 for eight terms.

Now an Item-Item matrix is created.

	I1	I2	I3	I4	I5	
I1	.	11	3	6	22	$\text{Term}_1 \text{ in Item}_1$
I2	11	.	12	10	86	$\text{Term}_1 \text{ in Item}_2$
I3	3	12	.	6	9	$\text{Term}_2 \text{ in Item}_1$
I4	6	10	6	.	11	$\text{Term}_2 \text{ in Item}_2$
I5	22	36	9	11	.	$\text{Term}_3 \text{ in Item}_1$

$\text{S}(I_1, I_2) = (T_{11} \times T_{21}) + (T_{12} \times T_{22}) + (T_{13} \times T_{23}) + (T_{14} \times T_{24}) + \dots + (T_{18} \times T_{28})$
 $= 0 + 4 + 0 + 0 + 0 + 4 + 0 + 3$
 $= 11.$

From the above matrix, let us now create an Item Relationship matrix, as shown below by applying a user defined threshold of '10', as shown below:-

	I ₁	I ₂	I ₃	I ₄	I ₅
I ₁		1	0	0	1
I ₂	1		1	1	1
I ₃	0	1		0	0
I ₄	0	1	0		1
I ₅	1	1	0	1	

Using clique alg for assigning items to classes produces the foll classes based upon the above matrix :-

Class 1 = Item 1, Item 2, Item 5

Class 2 = Item 2, Item 3

Class 3 = Item 2, Item 4, Item 5

Applying single link technique produces :-

Class 1 = Item 1, Item 2, Item 5, Item 3, Item 4.

Here Item 3 & 4 are added because of their similarity to Item 2.

The star technique produces :

Class 1 = Item 1, 2, 5

Class 2 = Item 3, 2

Class 3 = Item 4, 2, 5.

clustering by starting with existing clusters can be performed in a manner manner similar to term model.

Let's start with item 1 & 3 in class 1, and item 2 & 4 in class 2. Then the centroids are:-

$$\text{class 1} = \{ \frac{3}{2}, \frac{4}{2}, \frac{0}{2}, \frac{0}{2}, \frac{3}{2}, \frac{2}{2}, \frac{4}{2}, \frac{3}{2} \}$$
$$\text{class 2} = \{ \frac{3}{2}, \frac{2}{2}, \frac{4}{2}, \frac{6}{2}, \frac{1}{2}, \frac{2}{2}, \frac{2}{2}, \frac{1}{2} \}$$

The results of recalculating the similarities of each item to each centroid and reassigning turns is as shown:-

	Class 1	Class 2	Assign
I ₁	($\frac{3}{2}$)	($\frac{1}{2}$)	Class 1
I ₂	($\frac{2}{2}$)	($\frac{5}{2}$)	Class 2
I ₃	($\frac{0}{2}$)	($\frac{1}{2}$)	$\{ \text{Class 1} \} \rightarrow \text{Class 2}$
I ₄	($\frac{1}{2}$)	($\frac{2}{2}$)	Class 2
I ₅	($\frac{3}{2}$)	($\frac{4}{2}$)	Class 2

After calculating the centroid of class 2, item 3 will also take place in class 2.

III. Hierarchy of clusters:-

Hierarchical clustering in IRS focuses on the area of hierarchical agglomerative clustering methods.

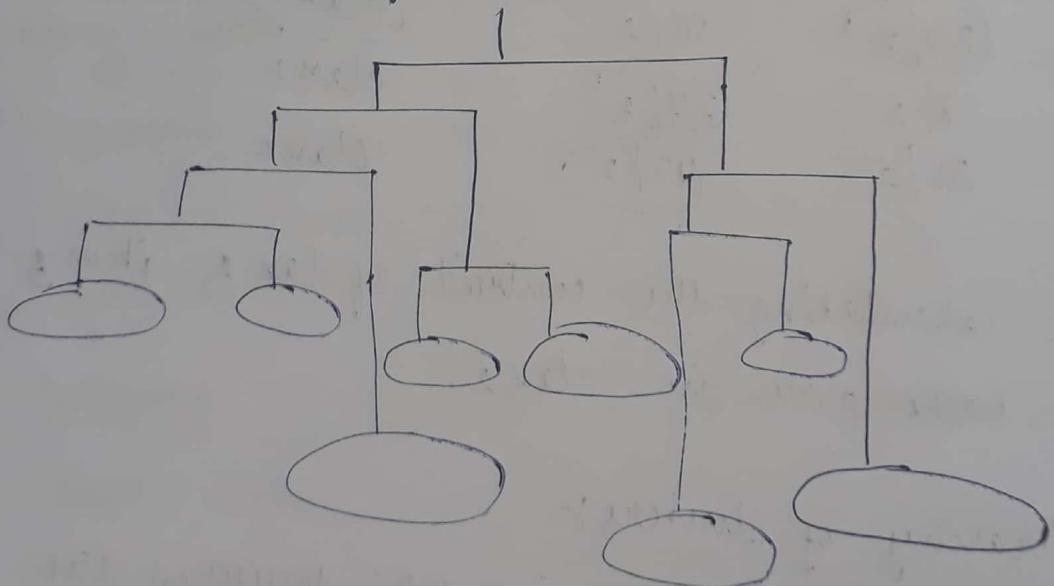
The term agglomerative means the clustering process starts with unclustered items and performs pairwise similarity measures to determine the clusters.

Divisive is the term applied to starting with a cluster and breaking it down into smaller clusters.

- Objectives:-
- * Reduce the overhead of search.
 - * Provide for a visual rep of info upon
 - * Expand the retrieval of relevant items

Search overhead is reduced by performing top-down searches of the outside of the cluster in the hierarchy and trimming those branches that are not relevant.

The use of dendograms along with visual cues on the size of clusters and strengths of linkage b/w clusters allows a user to determine alternate paths of browsing database.



Even without the visual display of the hierarchy, a user can use the logical hierarchy to browse items of interest.

A user, once having identified an item of interest, can request to see other items in the cluster.

Sanderson & Agt proposed the following methodology to build a concept hierarchy. This had also considered the ways of extracting terms from the doc. to rep. the hierarchy.

The terms had the foll characteristics:-

- * Terms had to best reflect the topic.
- * A parent term would refer to a more general concept than its child.
- * A child would cover a related subtopic of the parent.
- * A directed acyclic graph would represent relationships vs a pure hierarchy.
- * Ambiguous terms would have separate entries in the hierarchy for each meaning.

USER SEARCH TECHNIQUES

To understand the search process, it is first necessary to look at the different binding levels of the search str entered by the user to the db being searched. The selection and ranking of items is accomplished via similarity measures that calculate the similarity b/w the user search str and the wtd stored rep of the semantic of an item. Searching centroid can reduce search computation, but there is an associated risk of missing relevant items because of the averaging nature of centroid.

Hyperlinked items introduce new concepts in search originating from the dynamic nature of

(1) Search statements & Binding:

Search st are the sts of an info need generated by user to specify the concepts they are trying to locate in items. The search st uses traditional boolean logic and/or Natural Lang.

Binding is when a more abstract form is modified into a more specific form.

The user search statement is the user's attempt to specify the conditions needed to subset logically the total item space to that cluster of items that contain the info needed by the user.

The next level of binding comes when the search st is paired from use by a specific search sys.

The search st sys translates the query to its own meta language. This process is similar to the indexing of item process.

The final level of binding comes at the search sys. It is applied to a specific database. This binding is based upon the statistics of the processing tokens in the database and the semantics used in the database. This is especially true in statistical & concept indexing sys.

Frequently, in concept indexing sys, the concepts that are used as the basis for indexing are determined by applying a statistical alg against a rep sample of database.

Input

"find me info on the impact of the oilspill in Alaska on the price of oil"

impact, oil (petroleum), spills (accidents), Alaska, price (cost, value)

impact (0.308), oil (0.606),

Binding

User search \rightarrow using vocabulary of user.

statistical binding
extracts processing tokens.

Weights assigned to
search terms based
upon inverse doc freq
alg and database.

The above mentioned are the 3 potential levels of binding.

Parenthesis are used in the second binding step to indicate expansion by a thesaurus.

(*) The length of search \rightarrow directly affect the ability of IRS to find relevant items.
The longer the search query, the easier it is for the sys to find items.

* Profiles ~~are~~ used as search \rightarrow for selective dissemination of Info sys. are usually very long, typically 75 to 100 terms.

(2) Similarity Measures & Ranking :-

↳ Similarity Measures
HMM Techniques
Ranking Alg.

Searching in general is concerned with calculating similarity b/w a user search & items in the database.

The similarity may be applied to the total item or constrained to logical passage in the item.

→ Limiting the size of a passage to a fixed length size, locality based searching and similarity allows variable length passages based upon similarity of content.

This then leads to the ability to define locality based searching and retrieval of the precise locations of info that satisfy the query.

* Restricting the similarity measure to passage gain significant precision with minimal impact on recall.

The lack of a large no of terms makes it harder to find shorter passage that contain the search terms expanded from the shorter query.

⇒ Once items are identified as possibly relevant to the user query, it is best to present the most likely relevant items first. This process is called "Ranking".

(i) Similarity Measure:-

A characteristic of a similarity formula is that the results of the formula increase as the items become more similar.

The value is zero if the items are totally dissimilar.

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

If Item_j is replaced with Query_j , then the same formula generates the similarity b/w every item and Query_j .

The big problem with this simple measure is in the normalization needed to account for variance in the length of items.

→ Robertson & Spark Jones suggest that knowledge of terms in relevant items retrieved from a query should adjust the weights of those terms in the weighting process.

Taking into account the freq of occurrence of terms within an item producing the foll similarity formula, developed by Croft.

$$\text{SIM}(\text{DOC}_i, \text{Query}_j) = \sum_{i=1}^q ((C + \text{IDF}_i) * f_{i,j})$$

where C - constant used in tuning.

IDF_i - inverse doc freq for term i .

$$f_{i,j} = k + (k-1) \text{TF}_{i,j} / \text{maxfreq}_j$$

where k - tuning constant.

$\text{TF}_{i,j}$ - freq of term i in item j

maxfreq_j - max freq of any term in item j .

The best value of k ranges b/w 0.3 & 0.5.

To determine the "weight", an item has weight. The search st, the cosine formula is used to calculate the distance b/w the vector for the item and the vector for the query :-

$$\text{SIM}(\text{Doc}_i, \text{Query}_j) = \frac{\sum_{k=1}^n (\text{Doc}_{i,k} * \text{QTerm}_{j,k})}{\sqrt{\sum_{k=1}^n (\text{Doc}_{i,k})^2} * \sqrt{\sum_{k=1}^n (\text{QTerm}_{j,k})^2}}$$

where $\text{Doc}_{i,k}$ - k^{th} term in the weighted vector
for Item "i".

$\text{QTerm}_{j,k}$ - k^{th} term in the query "j"

The cosine formula calculates the cosine of
the angle b/w the two vectors.

If the two vectors are totally unrelated, then
they will be orthogonal and the value of
cosine is "0" otherwise "1".

Salton & Buckley improved the above formula
as :-

$$\text{QTerm}_{i,k} = (0.5 + (0.5 \cdot \text{TF}_{i,k} / \text{max freq}_k)) * \text{IDF}_i$$

where $\text{TF}_{i,k}$ is the freq of term "i" in the query
'k'
 max freq_k is the max freq of any term in
query 'k' and IDF_i is the Inverse Doc Freq
for term "i".

→ Commonly used other measures are the
Jaccard & the Dice similarity measures.
Hence, there is a change in the normalizing factor in the denominator to
account for diff characteristics of the
data.

The denominator in the cosine formula is invariant to the no of terms in common and produces very small nos when the vectors are large and the no of common elements is small.

* The Jaccard formula is as follows:-

$$\text{SIM}(\text{DOC}_i, \text{Q}_j) = \frac{\sum_{k=1}^n (\text{DOC}_{i,k} * \text{QTerm}_{j,k})}{\sum_{k=1}^n \text{DOC}_{i,k} + \sum_{k=1}^n \text{QTerm}_{j,k} - \sum_{k=1}^n (\text{DOC}_{i,k} * \text{QTerm}_{j,k})}$$

As the common elements increase, the similarity value decrease, but it always in the range -1 to +1.

* The Dice formula is :-

$$\text{SIM}(\text{DOC}_i, \text{Query}_j) = \frac{2 * \sum_{k=1}^n (\text{DOC}_{i,k} * \text{QTerm}_{j,k})}{\sum_{k=1}^n \text{DOC}_{i,k} + \sum_{k=1}^n \text{QTerm}_{j,k}}$$

Eg:- Query = (2, 2, 0, 0, 4)

DOC₁ = (0, 2, 6, 4, 0)

DOC₂ = (2, 6, 0, 0, 4)

	Cosine	Jaccard	Dice
DOC ₁	36.66	16	20
DOC ₂	36.66	-12	20

Note that: as long as the vector values are same independent of their order, the cosine and dice normalization factors do not change.

Many of the items have a similarity close or equal to zero.
For this reason, thresholds are usually associated with the search process.

The threshold defines the items in the resultant hit file from the query.

→ A document is always the case where the similarity is greater than zero.

The following illustrates the threshold process:-

Vector: American, geography, lake, Mexico, painting, oil, reserve, subject

DOC1 geography of Mexico suggests oil reserves are available vector (0, 1, 0, 2, 0, 3)

DOC2 American geography has lakes available everywhere.

vector (1, 3, 2, 0, 0, 0, 0, 0)

DOC3 paintings suggest Mexico takes as subject

vector (0, 0, 1, 3, 3, 0, 0, 2)

Query oil reserves in Mexico

vector (0, 0, 0, 1, 0, 1, 1, 0)

$$\text{SIM}(Q, \text{DOC1}) = 6$$

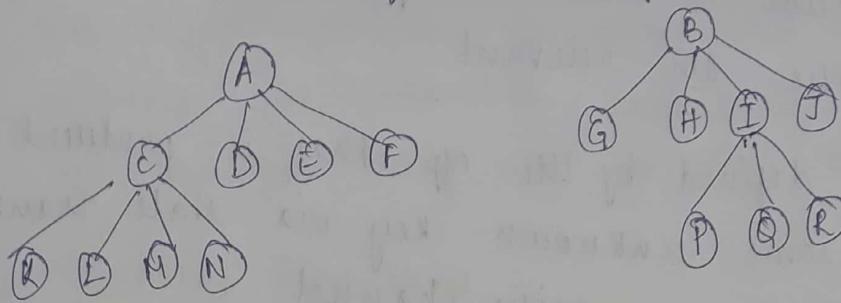
$$\text{SIM}(Q, \text{DOC2}) = 0$$

$$\text{SIM}(Q, \text{DOC3}) = 3$$

Here, the simple "sum of the products" minilab formula is used to calculate similarity b/w the query and each document.

If no threshold is specified, all three doce are considered true. If a threshold of 4 is selected, then only DOC1 is returned.

Another esp technique is to represent the clusters in the form of hierarchy.



The items are stored in clusters that are represented by the centroid for each cluster. In each letter at the leaf node represent an item in K, L, M... The letters at the higher nodes A, C, B, I represent the centroid of their immediate children nodes.

The hierarchy is used in search by performing a top-down process. The query is compared to the centroids "A" and "B".

If the results of the similarity measure are above the threshold, the query is then applied to the node's children. If not, that part of the tree is pruned and not searched. This continues until the actual leaf nodes that are not pruned are compared.

The use of centroids reduce the similarity computations but could cause a decrease in recall. It should have no effect on precision since that is based upon the similarity calculations at the leaf level.

(ii) Hidden Markov Model Techniques:

↳ used for searching textual corpora.
Here, the doc are considered as unknown statistical processes that can generate o/p that is equivalent to the set of queries that would consider the doc relevant.

HMM is defined by the o/p that is produced by passing some unknown key via state transitions through a noisy channel.

The observed o/p is the query, and the unknown keys are the relevant documents.

The noisy channel is the mismatch b/w the author's way of expressing ideas, and the user's ability to specify his query.

For each doc the probability that D was the relevant doc in the user's mind given that Q is the query produced is:

$$P(D \text{ is R} | Q) = P(Q|D \text{ is R}) * P(D \text{ is R}) / P(Q)$$

↳ This is application of Bayes rule to conditional probability.

Relevant doc will seem to be so sensitive to the specific queries, that trying to estimate $P(D|R)$ does not return any noticeable improvements in query resolution. Thus the probability that a doc is relevant given a specific query can be estimated by calculating the probability of the query being given the doc is relevant, i.e. $P(Q|R)$.

A HMM is defined by a set of states, a transition matrix defining the probability of moving between states, a set of OLP symbols and the probability of the OLP symbols given a particular state.

The biggest problem in using this approach is to estimate the transition probability matrix and the OLP for every doc in the corpus.

(iii) Ranking Algorithms :-

Ranking the OLP implies ordering the OLP from most likely item that satisfy the query to least likely item. This reduces the user overhead by allowing the user to display the most likely relevant items first.

→ With the inclusion of statistical similarity techniques into commercial systems and the large no of hits that originate from searching device corpora, such as the Internet, ranking has become a common feature of modern systems.

In commercial systems, heuristic rules are used to assist in the ranking of items.

Ranking
 └ Coarse grain ranking
 └ Fine grain ranking

The coarse grain ranking is based on the presence of query terms within items.

In the fine grain ranking, the exact rank of the item is calculated.

The coarse grain ranking is a weighted formula that can be adjusted based on completeness, contextual evidence or variety and semantic distance.

→ Completeness is the proportion of the no of query terms found in the item vs the no in the query.

→ Contextual evidence occurs when related words from the semantic w/w are also in the item.

→ Semantic distance evaluates how close the additional words are to the query term.

Synonyms add additional wt, antonyms decrease wt.

→ The coarse grain process provides an initial rank to the item based upon existence of words within the item.

→ Fine grain ranking consider the physical location of query terms and related words using factors of proximity in addition to the other three factors in coarse grain evaluation.

(3) Relevance Feedback :-

One of the major problems in finding relevant items lie in the difference in vocabulary b/w the authors and the user.

The user can use relevant items that have been found by the sys to improve future searches, which is the basis behind relevance feedback.

Relevant items are used to reweight the existing query terms and possibly expand the user search ~~st~~ with new items.

The relevance feedback concept was that the new query should be based on the old query modified to increase the wt of terms in relevant items and decrease the wt of terms that are in non-relevant items.

$$Q_n = Q_0 + \frac{1}{n} \sum_{i=1}^n DR_i - \frac{1}{nR} \sum_{j=1}^{nR} DNR_j$$

where

Q_n = the revised vector for the new query

Q_0 = original query

n = # of relevant items

DR_i = vector for the relevant item

nR = # of non relevant items

DNR_j = vector for the non-relevant items.

The revised version of the above formula is :-

$$Q_n = \alpha Q_0 + \beta \sum_{i=1}^n DR_i - \gamma \sum_{j=1}^{nR} DNR_j$$

where α, β, γ are the constants associated with each factor usually $1/m$ or $1/nR$ times a const.

$\beta \sum_{i=1}^n DR_i$ - the feedback

$\gamma \sum_{j=1}^{nR} DNR_j$ - -ve feedback

→ +ve feedback is more likely to move a query closer to user's info needs.

-ve feedback actually reduces the effectiveness of a query.

+ve feedback moves the query to retrieve items similar to the items retrieved and thus in the direction of more relevant items.

-ve feedback moves the query away from the non-relevant items retrieved.

	T_1	T_2	T_3	T_4	T_5
Q_0	3	0	0	2	0
DOC_1	2	4	0	0	2
DOC_2	1	3	0	0	0
DOC_{3n}	0	0	4	3	3
Q_n	$3\frac{3}{4}$	$1\frac{3}{4}$	0	$1\frac{1}{4}$	0

In the above table, Q_n is obtained as follows:-

Here there are 3 items - 2 relevant and 1 non-relevant.

Let us assume, $\alpha_s = 1$, $\beta = \frac{1}{4} (\frac{1}{2} \text{ times } b)$

$r = \frac{1}{4} (\frac{1}{4} \text{ times a constant } y_4)$

Now, sub these values in the formula for Q_n

$$\Rightarrow Q_n = (3, 0, 0, 2, 0) + \frac{1}{4} (2+1, 4+3, 0+0, 0+0, 0+0) \\ - \frac{1}{4} (0, 0, 4, 3, 3)$$

$$= (3\frac{3}{4}, 1\frac{3}{4}, 0, 8-13, 1\frac{1}{4}, 0) \\ \downarrow \text{choose}$$

Using the normalized similarity formula,

$$\text{SIM}(Q_K, \text{DOC}_i) = \sum_{i=1}^n \text{Term}_{K,i} * \text{Term}_{1,i}$$

Applying this results in :-

	DOC1	DOC2	DOC3
Q0	6	3	6
Qn	14.5	9.0	8.75

↳ Relevance feedback, (+ve feedback) has been proven to be of significant value in producing better queries.

↳ When the original query is modified based upon relevance feedback, the system ensure that the original query terms are in the modified query, even if -ve feedback would have eliminated them.

↳ Queries using relevance feedback produce significantly better results than those being manually enhanced.

↳ Pseudo-relevance feedback:- when user enter queries with a few ^{no} of terms, automatic relevance feedback based upon just the rank values of the items has been used.

This is also referred to as blind feedback or local context analysis.

(4) Selective Dissemination of Info Search :-

A dissemination sys is sometimes labeled as a "push" sys while a search sys is called a "pull" sys. The differences are that in a search sys the user proactively makes a query to the info sys to search.

In a dissemination sys, the user defines a profile and as new info is added to the sys it is automatically compared to the user profile.

A dissemination sys does not necessarily have a retrospective database associated with it.

One of the 1st commercial search techniques for dissemination was the Logicon Message Dissemination System (LMDS). It was designed for speed to support the search of thousand of profiles with items arriving every 20 sec.

This sys uses a least frequently occurring trigraph (3 characters) alg that quickly identifies which profiles are not satisfied by the item.

The potential hits profiles are analyzed in detail to confirm if the item is a hit.

→ An eg of a dissemination approach is the personal library s/w sys (PLS). This uses the approach of accumulating newly received items into the database and periodically running user profiles against the database. This makes max use of the retrospective search s/w but loses near real time delivery of items.

Eg: Retrieval Ware & Inloule s/w sys.

Retrieval ware uses a statistical alg but it does not include any corporate data.

INQUERY sys used against retrospective database, uses IDF info. It creates this info as it processes items, storing and modifying it for use as future items arrive.

The dissemination process is continuous, and the issue is the practicality of archiving all of the relevance judgements to be used in the relevance feedback process.

→ Another approach to dissemination uses a statistical classification approach and explicit error minimization to determine the decision criteria for selecting items for a particular profile.

In this case, the classification process is related to assignment for each item into one of two classes: * relevant to a user profile
* non-relevant.

↳ Eg. of alg used in linear classifiers that perform explicit error minimization are linear discriminant analysis, logistic regression and linear neural nw.

* ↳ Schütze et al. used two approaches to reduce dimensionality : * selecting a set of existing features to use
* creating a new much smaller set of features than the original features are mapped into.

→ A χ^2 measure was used to determine the most imp features.

The test was applied to a table that contained the no. of relevant - N_{R+} and non-relevant N_{R-} items in which a term occurs plus the no. of relevant and non-relevant items in which the term does not occurie $N_{R+} \rightarrow N_{R-}$ resp.

$$\chi^2 = \frac{N(N_{R+}N_{R-} - N_{R+}N_{R-})^2}{(N_{R+} + N_{R-})(N_{R+} + N_{R-})(N_{R+} + N_{R-})(N_{R+} + N_{R-})}$$

The chi-squared technique provides a more effective mechanism than freq of occurrence of terms. A high χ^2 score indicates a feature whose freq has a significant dependence on occurrence in a relevant or non-relevant item.

An alternative technique to identify the reduced feature (vector) set is to use a modified latent semantic index (LSI) technique to determine a new reduced set of concept vectors.

The use of the profile to define a local region is essential when working with large databases.

Linear discriminant analysis, logistic regression and neural nets are three possible techniques that were compared by Schütze et al.

Linear discrimination analysis uses the covariance class for each doc class to detect feature dependence.

Assuming a sample of data from two groups with n_1 and n_2 members, mean vectors \bar{x}_1 & \bar{x}_2 and covariance matrices C_1 and C_2 resp., the objective is to maximize the separation b/w the two groups.

This can be achieved by maximizing the distance b/w the vector means, scaling to reflect the structure in the pooled covariance matrix. Thus choose a such that :-

$$a^* = \arg_a \max \frac{a^T(\bar{x}_1 - \bar{x}_2)}{\sqrt{a^T C a}}$$

a^* is maximized where T is the transpose and

$$(n_1 + n_2 - 2)c = (n_1 - 1)c_1 + (n_2 - 1)c_2.$$

$\therefore c$ is the, the Cholesky decomposition of $C = RT$.

Let $b = Ra$, then the formula becomes:

$$a^* = \arg \max_b \frac{b^T R^{T-1} (\bar{x}_1 - \bar{x}_2)}{\sqrt{b^T b}}$$

which is maximized by choosing $b \propto R^{T-1} (\bar{x}_1 - \bar{x}_2)$

This means :-

$$a^* = R^{-1} b = C^T (\bar{x}_1 - \bar{x}_2)$$

The one dimensional space defined by $y = a^T x$ should cause the group means to be well separated.

Regularized Discriminant Analysis:- is to produce a non-linear classifier, a pair of shrinkage parameter is used to create a very general family of estimator for the group covariance matrix.

→ A 2nd approach is to use logistic regression proposed by Cox & Wermuth. It models a binary response variable by a linear combination of one or more predictor variables, using a logit link function :-

$$g(\pi) = \log(\pi / (1 - \pi))$$

This is achieved by modeling the dependent variable $\log(\pi/(1-\pi))$ as a linear combination of independent variables using a form $g(\pi) = \pi^\beta$. Here π - estimated response probability i.e. probability of relevance.

x_i - feature vector for document i

β - weight vector which is estimated from the matrix of feature vectors.

The optimal value of β can be calculated using the max likelihood and the Newton-Raphson method of numerical optimization.

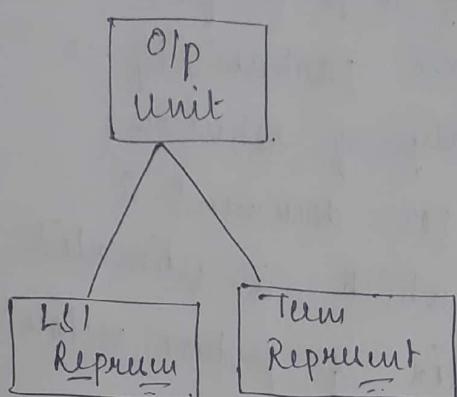
→ A 3rd technique is to use neural nw for the learning qn.

A neural nw is a nw of ilp and op cell based upon neuron fns in the brain, originating with the work of McCulloch and Pitts. Each ilp pattern is propagated forward through the nw.

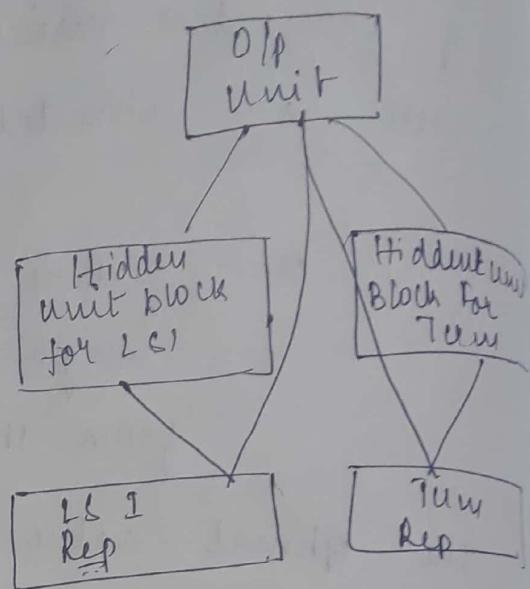
When an error is detected it is propagated backward adjusting the cell parameters to reduce the error, thus achieving learning. This technique is very flexible and can accommodate a wide range of distributions.

Whenever the error on the validation set increases, it indicates that overfitting is occurring and establishes the no of iterations on training that improve the parameter values while not halving

The linear & non-linear architectures for an implementation of neural n/w are as follows:



Linear neural n/w



Non linear Neural n/w

In the non-linear n/w, each of the hidden blocks consists of three hidden units.

A hidden unit can be interpreted as feature detector that estimate the probability of a feature being present in the i/p.

Propagating this to the o/p unit can improve the overall estimation of relevance in the o/p unit.

The n/w show the i/p of both terms and the LSI rep.

Relevance is computed by setting the activations of the i/p units to the doc rep and propagating the activation through the n/w to the o/p unit, then propagating the error back through the n/w using a gradient descent

alg as proposed by Rumelhart.

A sigmoid was chosen as:-

$$f(x) = \frac{e^x}{1+e^x}$$

This is an activation fn for the units of the n/w.

In this case backpropagation minimizes the same error as logistic regression.

The cross-entropy error is:-

$$L = -\sum (t_i \log \sigma_i + (1-t_i) \log (1-\sigma_i))$$

where t_i is the relevance for doc i and σ_i is the estimated relevance (activation of the o/p unit) for doc i.

The def of the sigmoid is:-

$$x = \log \left(\frac{f(x)}{1-f(x)} \right)$$

(5) Weighted Searches of Boolean Systems :-

The two major approaches to generating queries are Boolean and natural lang.

Natural lang queries are easily represented within statistical models and are usable by the similarity measure discussed.

Issues arise when Boolean queries are associated with weighted index sys.

Some of the issues are associated with how the logic operators fit with weighted values and how weights are associated with the query terms.

To integrate the Boolean and weighted system model, Fox and Sharat proposed a fuzzy set approach.

Fuzzy sets introduce the concept of degree of membership to a set.

The degree of membership for AND and OR operations are defined as:

$$DEG_{A \cap B} = \min(DEG_A, DEG_B)$$

$$DEG_{A \cup B} = \max(DEG_A, DEG_B)$$

where A and B are terms in an item.

DEG is the degree of membership.

The Mixed Min and Max (MMM) model considers the similarity b/w query and doc to be a linear combination of the min and max item wts.

Fox proposed the foll similarity formula:-

$$\text{SIM}(\text{QUERY}_{\text{OR}}, \text{DOC}) = C_{\text{OR}_1} * \max(\text{DOC}_1, \text{DOC}_2, \dots, \text{DOC}_n) + C_{\text{OR}_2} * \min(\text{DOC}_1, \text{DOC}_2, \dots, \text{DOC}_n)$$

$$\text{SIM}(\text{QUERY}_{\text{AND}}, \text{DOC}) = C_{\text{AND}_1} * \min(\text{DOC}_1, \text{DOC}_2, \dots, \text{DOC}_n) + C_{\text{AND}_2} * \max(\text{DOC}_1, \text{DOC}_2, \dots, \text{DOC}_n)$$

where C_{OR_1} and C_{OR_2} are weighting coefficients for the OR operation and

CAND1 and CAND2 - wt'ng coeff for the AND.

Lee and Fox found in their experiments that the best performance comes when CAND1 b/w 0.5 to 0.8 and COR1 > 0.2.

The MMM technique was expanded by Paice considering all item wts vs the max/min approach. The similarity measure is :-

$$SIM(\text{QUERY DOC}) = \frac{\sum_{i=1}^n \alpha^{i-1} d_i}{\sum_{i=1}^n \alpha^{i-1}}$$

where d_i 's are inspected in ascending order for AND queries and descending order for OR queries

α - wt'ng coeff.

Best values for α are 1.0 for AND queries and 0.7 for OR queries

→ An alternative approach is using the P-norm model which allows terms within the query to have wts in addition to the terms in the items.

Similar to the cosine similarity technique, it considers the membership values (d_{A1}, \dots, d_{An})

to be coordinates in an "n" dimensional space.

For an OR query, the origin is the worst possibility.

For an AND query the ideal pt is the unit vector where all the d_i values equal 1. Thus the best ranked docs will have max distance from the origin in an OR query and minimal distance from the unit vector. The generalized queries are:-

$$Q_{OR} = (A_1, a_1) \text{ OR } (A_2, a_2) \text{ OR } \dots \text{ OR } (A_n, a_n)$$

$$Q_{AND} = (A_1, a_1) \text{ AND } (A_2, a_2) \text{ AND } \dots \text{ AND } (A_n, a_n)$$

The operators will have a strictness value α that varies from 1 to ∞ where ' α ' is the strict def of Boolean operator.

The ' a_i ' values are the query term wts.

If we assign the strictness values to a parameter labeled "s" then the similarity formulas b/w queries and items are:-

$$\text{SIM}(Q_{OR}, \text{DOC}) = \sqrt{s(a_1^s d_{A1}^s + \dots + a_n^s d_{An}^s)} / (a_1^s + a_2^s + \dots + a_n^s)$$

$$\text{SIM}(Q_{AND}, \text{DOC}) = 1 - \sqrt{s(a_1^s (1-d_{A1})^s + \dots + a_n^s (1-d_{An})^s)} / (a_1^s + a_2^s + \dots + a_n^s)$$

$$\text{SIM}(Q_{not}, \text{DOC}) = 1 - \text{SIM}(Q, \text{DOC})$$

→ The objective of Salton's proposal is to perform the normal Boolean operations and then refine the results using weighting techniques.

The fol procedure is a modification to his approach for defining search results.

- Normal Boolean operations produce the fol results:
 - * "A OR B" retrieves those items that contain the term A or the term B or both.
 - * "A AND B" retrieves those items that contain both terms A and B.
 - * "A NOT B" retrieves those items that contain term A and not contain term B.

If wts are then assigned to the terms b/w the values 0.0 to 1.0, they may be interpreted as the significance that we are placing on each term.

The value 1.0 is assumed to be the strict interpretation of a Boolean query.

0.0 is interpreted ~~as~~ to mean that the user places little value on the term.

Under these assumptions, a term assigned a value of 0.0 should have no effect on the retrieved set.

Thus:

- * "A, or B₀" should return the set of items that contain A as a term.
- * "A, AND B₀" will also return the set of items that contain term A.
- * "A, NOT B₀" also return set A.

The alg follows the foll steps:-

1. Determine the items that are satisfied by applying strict interpretation.
2. Determine the items that are part of the set that is invariant.
3. Determine the centroid of the invariant set.
4. Determine the no. of items to be added or deleted by multiplying the item wt times the no. of items outside of the invariant set and round up to the nearest whole no.
5. Determine the similarity b/w items outside of the invariant set and the centroid.
6. Select the items to be included or removed from the final set.

Eg:- ~~Query ends up with~~

Weighted Boolean Query.

	Computer	Program	cost	Sale
D ₁	0	4	0	8
D ₂	0	2	0	0
D ₃	4	0	2	4
D ₄	0	6	4	6
D ₅	0	4	6	4
D ₆	6	0	4	0
D ₇	0	0	4	0
D ₈	4	2	0	2

$Q_1 = \text{QUERY}_1 = \text{Computer}_{1.0} \text{ OR Program}_{.833}$

$Q_2 = \text{QUERY}_2 = \text{Cost}_{0.75} \text{ AND Sale}_{1.0}$

Q_1 strict interpretation = $(D_1, D_2, D_3, D_4, D_5, D_6, D_8)$

Q_2 strict interpretation = (D_3, D_4, D_5)

Q_1 invariant = (D_8)

Q_2 invariant = (D_3, D_4, D_5)

Q_1 optional = $(D_1, D_2, D_3, D_4, D_5, D_6) \Rightarrow \lceil .333 \text{ times } 6 \text{ items} \rceil$

Q_2 optional = $(D_1, D_8) \Rightarrow \lceil (1 - 0.75) \text{ times } 2 \text{ items} \rceil = 2 \text{ items}$

Here, QUERY_1 ends up with a set containing all of the items that contain the term "Computer" and two items from the set "Computer" NOT "program". The symbol $\lceil \rceil$ stands for rounding up to the next integer.

In QUERY_2 , the final set contains all of set "cost" AND "sale" plus 0.25 of the set of "sale" NOT "cost".

Now, using the similarity measure

$$s(\text{Item}_i, \text{Item}_j) = \sum (T_{\text{Item}_i k}) T_{\text{Item}_j k}$$

Gives the following set of values based upon the centroids:

$$\text{Centroid } (Q_1) = (D_8) = (4, 2, 0, 2)$$

$$\text{Centroid } (Q_2) = (D_3, D_4, D_5)$$

$$= \frac{1}{3} (4+0+0, 0+6+4, 2+4+6, 4+6)$$

Now, apply similarity,

$$S(\text{Centroid } Q_1, D_1) = \cancel{S}(4 \times 2 + 2 \times 4 + 0 \times 0 + 2 \times 8) \\ = 24$$

$$S(CQ_1, D_2) = 0 + 4 + 0 + 0 = 4$$

$$S(CQ_1, D_3) = 24$$

$$S(CQ_1, D_4) = 24$$

$$S(CQ_1, D_5) = 16$$

$$S(CQ_1, D_6) = 24$$

$$S(CQ_2, D_1) = \frac{1}{3} (0+40+0+112) \\ = \frac{1}{3} (152)$$

$$S(CQ_2, D_8) = \frac{1}{3} (16+20+0+28) \\ = \frac{1}{3} (64)$$

Now, For Q_1 , two additional items are added to the invariant set $(D_8) \cup (D_1, D_3)$ by choosing the layout no items because of the tie at 24, giving the answer of (D_1, D_3, D_8) .

For Q_2 , one additional item is added to the invariant set $(D_3, D_4, D_5) \cup (D_1)$ giving the answer (D_1, D_3, D_4, D_5) .

(6) Searching the Internet and Hypertext :-

Some of the most commonly used nodes to search items on Internet are YAHOO, AltaVista and Lycos. The primary design decisions are on the level to which they retrieve data and their general philosophy on user access.

Lycos and Alta Vista automatically go out to other Internet sites and return the text at the site for automatic indexing.

Lycos returns home pages from each site for automatic indexing while Alta Vista indexes all of the text at a site.

The search process is directly available to the user via Intelligent Agents (IA).

IA provide the capability for a user to specify an info need which will be used by the IA as it independently moves b/w Internet sites locating info of interest.

There are 6 characteristics of IA:-

(i) Autonomy:- The search agent must be able to operate without interaction with a human agent. It must have ctrl over its own internal states and make independent decisions.

(ii) Communication Ability:- The agent must be able to communicate with the info sites as it traverses them. This implies a universal accepted lang defining the external interface.

(iii) Capacity for cooperation: This concept suggests that intelligent agents need to cooperate to perform mutually beneficial tasks.

(iv) Capacity for Reasoning:

- * Rule based - where user has defined a set of conditions and actions to be taken.
- * Knowledge-based - where the IA have stored previous conditions & actions taken which are used to deduce for future actions.
- * Artificial evolution based: where IA spawn new agents with higher logic capability to perform its objectives.

(v) Adaptive Behavior: closely tied to I and A, adaptive behavior permits the intelligent agent to assess its current state and make decisions on the actions it should take.

(vi) Trustworthiness - the user must trust that the intelligent agent will act on the user's behalf to locate info that the user has access to and is relevant to the user.

→ A Hyperlink is an embedded link to another item that can be instantiated by clicking on the item reference.