

Unit - I

Information Retrieval

Introduction to Information retrieval systems

- Definition of IRS
- Objectives of IRS
- Functional overview

- Relationship to database management systems

- Digital libraries & data warehouses

Information retrieval system capabilities

- Search capabilities
- Browse capabilities
- Miscellaneous capabilities.

Definition of TRS

An Information Retrieval System (IRS) is a system designed to store, manage, search and retrieve relevant information for large collections of data. The information can be in the form of text, images, audio, video or multimedia content.

Ex:- Google search, library databases

1. Main purpose
 - helps users find relevant information quickly
 - stores large amounts of data and retrieves it when needed

2. How TRS works
 - stores information in a structured way

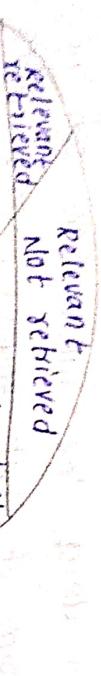
- uses indexing to organize data for fast retrieval
 - uses search algo to find the best results for user queries.

Key difference :

- Data Retrieval gives exact results based on structured queries.
- Information Retrieval finds relevant information from unstructured data using ranking algorithms.

Chapter 4. Information Retrieval Systems (IRS)

The primary objective of an IRS is to help users efficiently locate and retrieve relevant information while minimizing the effort required.



$$\text{Recall} = \frac{\text{Number of Relevant documents retrieved}}{\text{Total No. of Relevant Documents in the Database}}$$

Retrieved

Minimizing user overhead

- Reducing the time and effort spent on query generation, execution, scanning results, and reviewing non-relevant items.
- Ensuring users can quickly access the most useful information

2. Maximizing Precision and Recall

a. precision

Precision measures how many of the retrieved documents are actually relevant to the user's query

$$\text{Precision} = \frac{\text{Number of Relevant documents Retrieved}}{\text{Total No. of Retrieved documents}}$$

- Suppose a search query retrieves 50 documents out of which 30 are relevant.

$$\text{precision} = \frac{30}{50} = 0.6 \text{ (or } 60\%)$$

This means 40% of the retrieved documents were irrelevant.

b. Recall

recall measures how many of the total relevant documents in the database were retrieved.

$$\text{Recall} = \frac{\text{Number of Relevant documents retrieved}}{\text{Total No. of Relevant Documents in the Database}}$$

Ex:-

- Suppose there are 100 relevant documents in total, but only 30 were retrieved by the system.

- Recall = $\frac{30}{100} = 0.3$ (or 30%).
- This means 70% of relevant documents were missed, reducing the completeness of search results.

Trade-off :-

- High precision, low recall:

Retrieves fewer but more irrelevant results (e.g. strict keyword matching)

- High recall, low precision:

Retrieves more results, including many irrelevant ones (e.g., broad searches)

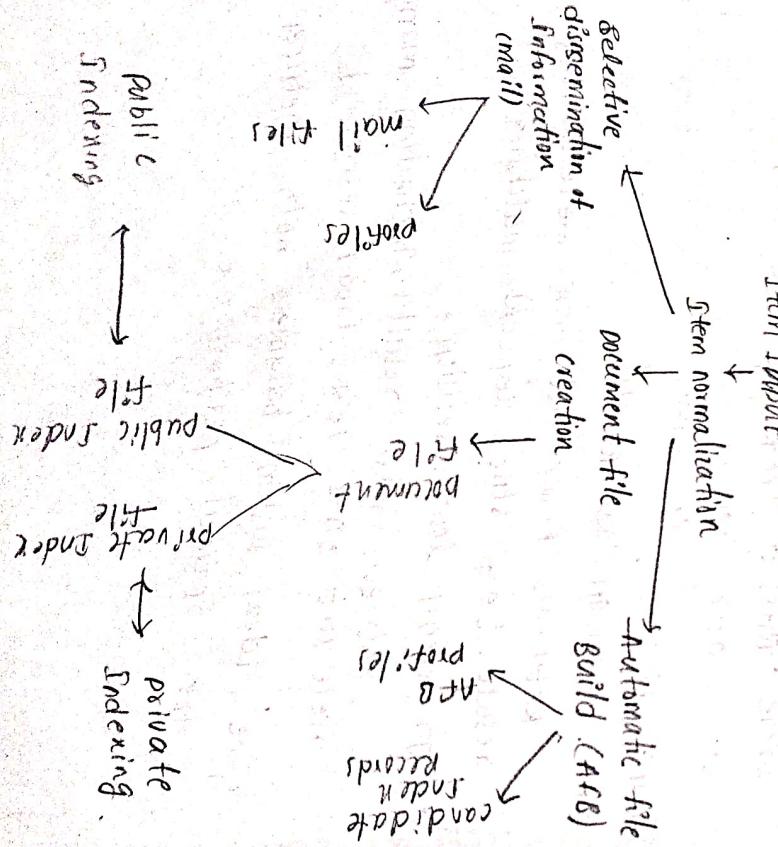
- The ideal system balances both metrics for optimal performance

Functional overview

An Information Retrieval System (IRS) is designed to help users find relevant information efficiently from a large collection of data. The system processes queries, searches the database to retrieves results based on certain criterion.

It comprises four major functional processes

1. Item Normalization
 2. Selective Dissemination of Information
 3. Document Database Search
 4. Index Data base search (including AFIS project)



Item Normalization is the first step in an integrated IRS. It involves converting incoming items into a standardized format to ensure consistency and ease of processing.

This process involves:
a standardization of Input formats:-

- Translating various external formats into a single, consistent data structure
- logical restructuring:- Dividing items into logical subdivisions through the process called zoning

This helps in improving search precision and optimizing display.

g. Tokenization:

- Identifying processing tokens from the text.

- Applying stemming (e.g. removing word endings) to normalize tokens for each purposes.
- Using stop lists / algo to eliminate low-value tokens

d. Characterization of Tokens

- Assigning characteristics to tokens to aid in disambiguation during searches.

e. Multimedia Normalization

- Synchronization b/w text and multimedia is crucial for precise retrieval.

f. Selective Dissemination of Information

- It is also known as mail process. The main purpose of this process is to perform comparison b/w recently received items and the existing item statements which are required by user key components

g. Index Database Search (IDBS)

- The Index Database search process allows users to create and search indexes for efficient retrieval of items key aspects:-

1. Index files

- Users can create private index files or use public index files.

2. Automatic file Build

- A process that automatically generates candidate index records for incoming items.

3. Document Database Search

- The Document Database search process allows user to perform retrospective searches against all items stored in the system.

Key features

- Document Database :-
 - Contains all items received and processed by the system.
 - Items are typically static and may be partitioned by time for archiving process.

Relationship to Database Management Systems

DBMS and DBMS serve different purposes but are increasingly being integrated for more efficient data management. The relationship stems from the need to handle both structured data and unstructured information.

DBMS vs. Information Retrieval Systems

Structured (e.g., tables, records, attributes)

Unstructured or semi-structured (e.g., document, articles, abstracts)

Well-defined schema with pre-defined relationships

Free from text with minimal structure.

Precise queries using SQL

Keyword-based, full-text searches, relevance ranking.

Exact matches, tabular format

Ranked results, relevance-based

Retrieves specific data refining queries over time.

Instantly

Integration of DBMS & PRS

Modern DBMS platforms incorporate

PR features to handle both structured & unstructured data.

examples of Integration:

- Inquire DBMS:

One of the earliest systems to combine structured and unstructured data retrieval.

- Oracle DBMS Connectors:

embeds an information retrieval system that uses a thesaurus based search

to improve accuracy.

- Informix DBMS (Cognitiveware):

links structured databases with powerful search tools for text based data.

Importance of Integration

- enables efficient data retrieval across structured and unstructured data sources.
- enhances search capabilities with ranking, relevance feedback, and thesaurus-based suggestions.
- supports business intelligence and decision making by providing a holistic view of data.

Digital Libraries and Data warehouses

Digital libraries and data warehouses (or Data Mart) are two major systems that overlap with PRS, but they serve distinct purposes.

while all three systems act as repository of information, they differ in their focus.

Structure and retrieval mechanisms.

1. Digital libraries

This originate from traditional libraries which have always been concerned with storing and retrieving information

As digital formats expanded, libraries

evolved from electronic libraries (1991-1993) to digital libraries (post-1995)

challenges & opportunities in digital libraries

Indexing & cataloging

Search intermediaries

Copyright & Intellectual property

Digital preservation

Digital Libraries vs IRS

Similarities:

Both store and retrieve textual information

Difference:

- IRS mainly focuses on search & retrieval
- Digital libraries address broader concerns like copyright, long-term accessibility & source authenticity

a accessibility & source authenticity

2. Data warehouses

Originating in the commercial sector, data warehouses help organizations control, organize & retrieve structured digital data.

Unlike digital libraries which focus on textual data, data warehouses handle structured, quantitative data for decision-making.

Components:-

1. Data storage
2. Information directory

3. Data processing & movement

4. Search & retrieval tools

5. Data export mechanism

Data warehouses vs IRS

Similarities:

Both allow users to search & retrieve info

Differences

- Data warehouses deal with structured data for decision support & analytical

- IRS focuses on textual, unstructured data with iterative searching and relevance ranking.

Conclusion:- Integration of these systems while digital libraries, data warehouses, and information retrieval systems have distinct roles, they are increasingly integrated to provide comprehensive data management.

- Digital libraries enhance access to unstructured textual data.
- Data warehouses facilitate structured data analytics and decision making.
- IRS enables search, ranking and retrieval of relevant textual information.

Information Retrieval System Capabilities

An Information Retrieval system (IRS) provides various capabilities that allow users to search, browse and interact with information effectively.

These capabilities are generally grouped into 3 categories.

1. Search capabilities

Search capabilities define how users query an IRS to find relevant information. A robust system offers multiple search options to accommodate different user needs and query types.

Boolean Logic

It has values True or False

Operator :- NOT, AND, OR

| | | |
|------------|--------------------------|------------------------|
| example :- | 1. Mobile OR Laptop ✓ | 2. mobile AND laptop |
| | doc 1 (mobile, laptop) ✓ | doc 1 (mobile, laptop) |
| | doc 2 (mobile) ✗ | doc 2 (mobile) ✗ |
| | doc 3 (laptop) ✓ | doc 3 (laptop) ✗ |

Doc 4 ✗

- Proximity search

finds documents where 2 or more terms appear close together, within a specified distance.

e.g:- "data" Neatly "mining" → finds "data" and "mining" within 3 words of each other.

- Contiguous word phrases when 2 or more words that are treated as a single semantic unit.

Searches for exact sequences of words in the same order.
e.g:- United state of America → only retrieves document containing that exact phrase.

- fuzzy logic search

- Allows approximate matching rather than exact matches
- Handles misspelling, typos or variations in text.

e.g:- Searching for "analyze" may also match "analyse", uses special characters to match partial terms or unknown characters.

e.g. - search matches computer, computing, computation, text, matches text, test, test.

- numerical and date range search.

- finds documents that contain numbers or dates within a specified range.

- especially useful in finance or academic database.

e.g:- 2016 - 2012

price [100 to 500]

- Thesaurus expansion (concept-based search)

- expands the search to include synonyms or related terms using a thesaurus or ontology.

e.g:- search term "car" may also results containing "automobile", "vehicle",

- Natural language queries:

- allows users to types queries in everyday human

language.

e.g:- "what are the applications of artificial intelligence?"

- Multimedia Queries

- enables search using non-textual input like images, audio, video or sketches.

e.g:- upload an image to search for visually similar image (e.g. Google Images).

- search for audio clips that sound similar.

Browse capabilities

- It allows users to navigate through retrieved documents or databases without specifying precise search term.

- They are essential for exploratory search where user may not know exactly what they're looking for.

User → query → PRS → Results → Browse capabilities. sorting results.

• Navigate • visualization • organize.

Key :-

↳ Relevance Ranking

- each retrieved item is assigned a relevance score

usually between 0.0 and 1.0

- results are presented in ranked order, showing the most relevant items first.

e.g:- A document with score 0.92 appears above one with 0.60.

User → query → PRS → Result Doc 1, Doc 2
Doc 3, Doc 4
↓ relevance ranking

Doc 1
Doc 3
Doc 4

2. Zoning :-

Only relevant parts (zones) of documents are shown such as title, abstract, or specific matching passage.

- helps users quickly assess relevance without reading the full document.

3. highlighting

This capability significantly enhances the browsing experience by making it easier for users to identify relevant information quickly.

- Matched terms are visually highlighted using bold text, color or underlining.
- Helps users immediately spot why a document was retrieved.
- Especially useful in long documents.

Browse capabilities enhance the user experience by helping users navigate, evaluate and define search results more easily - especially when they're uncertain about their query or exploring a new topic.

Miscellaneous Capabilities:

These are additional features in an IRS that enhance, usability, flexibility and efficiency during search and retrieval processes. They support both new and experienced users in defining and refining the queries.

i. Search History Log

It records all the queries submitted by the user during a session.

- Allows users to revisit previous searches

- Helps in refining or modifying past queries without retyping.

- useful for long research tasks

and comparative analysis.

e.g:- You search "Machine learning", later refine to "machine learning in healthcare" - the system keeps both entries in the history.

ii. Bookmarking

A bookmark saves the location or reference to a retrieved item, so the user can revisit it later.

- used when users plan to return to the result after reviewing others.

e.g:- In a digital library, you bookmark a journal article that you want to read later.

3. Canned Query / stored query

A canned query (also called a stored query) is a pre defined search query saved by the system or the user for future reuse.

purpose:- saves time for repetitive or common searches

- can include placeholders or variables for runtime inputs.

e.g:- A librarian creates a stored query for

"title contains 'AI' AND year > 2020".

each time it's used, the user just changes the keyword or year.

Miscellaneous capabilities improve the overall effectiveness and reusability of an IRS by enabling better interaction, reuse and system compatibility.

Unit-II

Cataloguing and Indexing

- History and objectives of Indexing

- Indexing process
- Automatic indexing

Information extraction Data Structure

- Introduction to Data structure
- Stemming Algo
- Inverted file structure
- N-gram Data structure
- PAT Data structure
- Signature file structure
- Hypertext and XML Data structure
- Hidden Markov model.

1. History of Indexing

- Ancient cataloguing (before the 19th century)

- Babylonian Libraries (3rd millennium BCE):

early forms of indexing appeared in cuneiform-tablets, arranged by subject.

- medieval libraries:-

Basic cataloguing methods developed, but indexing remained mostly manual and localized.

19th & 20th century Developments

• late 1800s:

- The Dewey decimal system introduced hierarchical subject indexing.
- Libraries started organizing materials systematically.

~~19th & 20th century Developments~~

◦ Labels - MARC Machine Readable cataloguing

- Developed by the Library of Congress (1966-69) for computerized cataloguing

- Standardized bibliographic records for better sharing among libraries

◦ 1970's - DIALOG system

- first commercial indexing system, originally developed by Lockheed for NASA
- Became a public service in 1978, due to expanding access to scientific & technical

Modern Developments (1990s - present)

- 1990s
- full-text indexing became possible
- with fast digital storage and computing advances.
- indexing shifted from manual cataloguing to automated systems using algo.
- 2000 - AI & ML
- + semantic indexing and concept-based search started replacing keyword search in NLP improved search accuracy.
- Objectives of Indexing
- Indexing help users find relevant information efficiently. Its main objectives are:
 - * Enhance searchability & retrieval
 - * Improve precision & recall
 - * Support user needs
 - * Enable ranking & clustering
 - * Adapt to evolving information delivery

Data Structure

- each data structure has a set of capabilities
- Ability to represent the concept and their relationships.
- Two major data structures
 - stores and manages the received items in their normalized form (Document manager)
 - contains the preprocessing tokens & associated data to support search (document search manager)
- When searching is completed
 - the result is passed to the document for review
 - the result is passed to Document normalization
 - the result is passed to Document file creation
- focus:- DS supports search in Document manager
- stemming :-
 - It is transformation of often applied to data before placing it in searchable ds.
- DS :-
 - 1. stemming
 - Porter stemming algo
 - Dictionary lookup stemmers
 - successors stemmers.
- It represents in canonical & morphological forms.

244 → , 761

Plated - Plate

— 61 —

happy → happy

'Able' → nothing

Adjustable → adjust

Dictionary lookup stemming

- Tries to avoid collapsing words with different meaning to same root
 - The original word or stemmed version is locked up in a dictionary and replaced by the best stem.
 - This technique has been implemented in Inquiry and Remarque system.
 - Inquiry system technique called kstem.
 - kstem is a morphological analyzer that conflates words variants to a root form.
 - It uses 6 major data files to control & limit the stemming process.
 - Successor stemmers.
 - It is based on the length of prefixes
 - The smallest unit of speech that distinguishes one word from another

These process uses successor varieties for each node.

word

Here once the algorithm takes substrings of a term will decrease as more characters are added until a segment boundary is reached.

- The stemming method based on this work

uses letters in place of phonemically

we can use 4 methods to successfully

variety of words

卷之六

A cutoff value is selected to define the

stem length.

segment when successor variety > threshold

EM. - READABLE
CORPUS - ABLE, APE, REATABLE, FIXABLE, READ,

READABLE, READING, READS, RED, ROPE,
ROPE

prefix successor letter
epo

R
P
C
A
D

REPRINTED FROM THE JOURNAL OF POLYMER SCIENCE

READAB
READABL
L
E

Peak and plateau

A segment break is made after a character whose successor variety exceeds that of the character.

- break at the character whose successor variety is greater than both its preceding and following characters.

variety is greater than both its preceding and following characters.

variety is greater than both its preceding and following characters.

complete word method

Breaks on boundaries of complete words

Breaks made if the segment is a complete word in the corpus (READ)

entropy methods-

uses the distribution method of successor variety letters

Inverted-file structure

Linear Kumar

- It is a data structure, it allows efficient full-text search in the database.

- It stores a mapping of words to their locations in the database table or document.

- Inverted file based on methodology of storing an inversion of document

3 files are created or accessed

1. Inversion / posting list / posting file

2. dictionary / index file

3. Inverted file / index file

- For each word a list of documents in which the word is found is stored.

- Each document is given a unique the numerical identifier that is stored in

inversion list which is stored list in the system

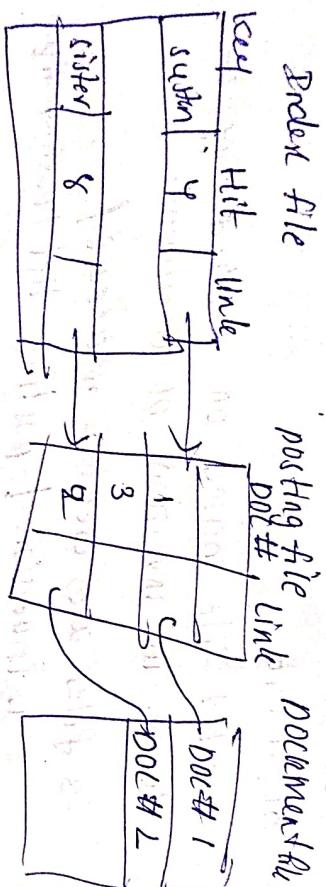
and a pointer to the location of the its inversion list.

- Using zoning to improve

- Inversion list consists of document identifier for each document in which the word is found.

- Inversion list are used to store concept and their relationship.

here using $\#$ to represent the entwining symbol which any one of a set of



- n grams can be viewed as special technique for conflation and as a unique data structure: in DS
 - These are fixed length of consecutive series of ' n ' characters
 - In stemming we care about semantic but n -gram we don't care.
 - Here we are observing unigrams - how often a word occur without looking at previous word
 - bigram - only the previous word to predict the current word
 - Trigram - previous two words are

count(wi-1)

hitting press — no word

2 years - 1 year from present - 2 years - 3 years - 2 years after

Hello! Where are you now?

600

Chuio 6

- each of the n-grams created becomes a separate processing token and searchable.
 - It is possible that the same n-gram can be created multiple times from a

Indexing, processing and cataloguing

Indexing:

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

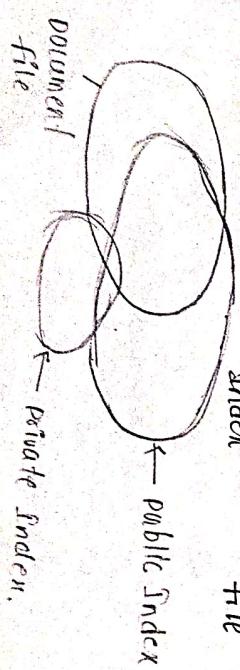
- Indexing process can be manual or automatic.
- Creating a direct search in document database or indirect search through index files.
- Once the searchable data structure has been created, techniques must be defined to correlate the user entered query statement to the set of items in DB.
- Indexing originally called cataloguing.

Objectives:

- Represent the concepts within an item to facilitate the users finding relevant information.
- The full text searchable data structures for items in document file provides a new class of indexing called total document indexing.
- Other objectives are indexing, ranking, item clustering.

Index processing

Indexer → public / private → document file



When an organization with multiple indexers decide to create a public or private index.

- Scope of the indexing - define what level of detail the subject index will contain.
- Second decide to link index term together in a single index for a particular concept.

Scope of indexing

- When performing the indexing manually, problems arise from a sources the author and the indexes.
- Vocabulary domain may be different the author and the indexes.
- The indexer must determine when to stop the indexing process.

- Two factors involved in deciding
 - o exhaustivity - refers to how many concepts from a document are indexed.
 - o specificity - How precise or general the index terms are.

Precordination and linkages

- Precordination is a process of linking related index related index terms during the indexing phase of document processing.

- Post coordination - in which index terms are not linked during indexing, which index terms are not linked during indexing, and the relationships are established during search and the relationships are established during search time, it's called post coordination.

- Use of linkages:
 - o linkages connect attributes of concept (e.g., actor, action, location).
 - o help disambiguate similar terms

Text processing

1. Document Parsing:-

documents come in all sorts

of the same document may

contain multiple languages or

formats.

- Document parsing deals with the recognition and "breaking down" of the document structure into individual components.

2. Lexical analysis

Document parsing after lexical analysis tokenizes a document, seen as input

stream into words.

3. Stop-word removal:-

A subsequent step optionally applied to the results of lexical analysis is stop-word removal.

- The removal of high-frequency words

4. Phrase detection

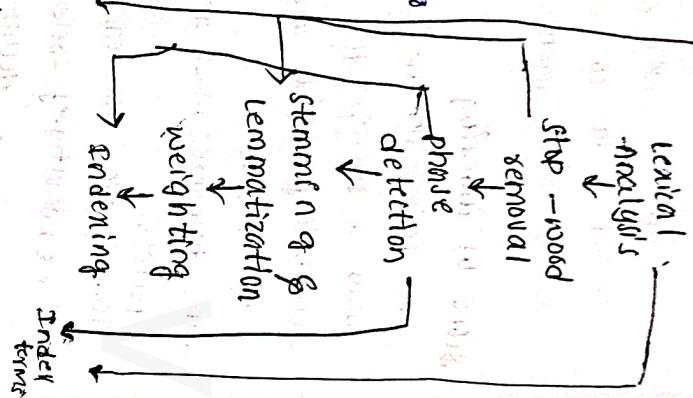
This step captures text meaning beyond word sense.

bag-of-word approaches, thanks to the identification of noun groups and other groups.

5. stemming and lemmatization:

stemming & lemmatization aim at stripping down word suffixes in order to normalize the word.

- Stemming is a heuristic process that "chops off" the ends of words in the hope of achieving the goal correctly most of the time.



Automatic Indexing :-

It is the process by which a computer system automatically selects and assigns index terms to a document, without human intervention. It is faster, more scalable, and less biased compared to manual indexing.

- Automatic indexing having advanced technologies such as natural language processing (NLP) and machine learning.

- It finds relevant terms or keywords from the documents automatically and create an index.

- First system scans the text and identify important words to create an index.

How it works?

Advantages:-

1. Document parsing
2. Text preprocessing
3. Term Selection.
4. Term weighting
5. Index Creation.

Challenges:-

- May miss contextual meaning
- Ambiguity & polysemy
- Harder to handle domain-specific language without proper training.

Inverted file structure

It is mostly widely used data structure in Information retrieval systems. It is used to map terms (words) to the list of documents in which they appear. This is the core of fast full-text searching.

Why Inverted?

Because instead of listing all words per document, we store all documents per word - which is the inverse of the document structure.

Components

1. Document file : stores the actual document
2. Dictionary : stores all unique terms in the document collection in sorted order.

3. Inversion list / posting list :-

for each term in the dictionary, stores a list of document IDs where the term appears.

How it works:-

1. each document is assigned a unique numeric document ID.

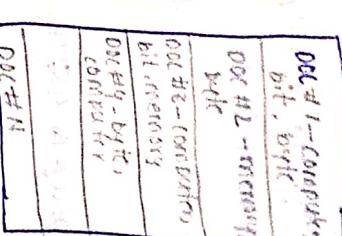
2. for every term an inversion list (posting list) stores:

- The document IDs where the term occurs.
- Optionally : positions, frequencies or weights.

exo- document file

Dictionary file

Inversion file list



Additional features

- Zoning: limits search to certain zones to increase precision
- weights
- ranking support.

Disadvantages

- Post search
- Scalable
- Flexible
- supports ranking.

The inverted file structure is the backbone of modern information retrieval. It allows fast, flexible and precise searching by mapping terms to documents, enabling efficient query processing even at massive scale.

N-Gram data structure

- N-Gram is a sequence of n consecutive characters extracted from text. It is used to analyze, search and match sets of words or phrases without considering the meaning (semantics).

sky Al-gramm

- N-Gem's break tent into overlapping fixed-length tokens, helping match terms even if misspelled

- 13 -

1. Convert each word / phrase into overlapping N-grams.
 2. Store them in the Index.
 3. During search, the query converted into N-grams

and matched.

$$p(w_i | w_{i-2}, w_{i-1}) \rightarrow \frac{\text{count}(w_{i-2}, w_{i-1}, w_i)}{\text{count}(w_{i-2}, w_{i-1})}$$

Statement

DOD like

corpus database

I am Amit
like computer

Do emit like computers

Do you like computers?

I do like am I.

$$\text{if } n_1 = 1 \\ w_{i-1} = d_0 \rightarrow \frac{\text{count}(d_0, T)}{\text{count}(d_0)} = \frac{2}{4}$$

$$w_i = \alpha m \rightarrow \frac{\text{count}(a_i)}{\text{count}} = \frac{0}{1}$$

3. $w_i^* = \text{elke const. do, stuk}$ \Rightarrow const. (do) $= \frac{q}{q}$

$w_{i-1} = \text{do}$ const. (do)

4. $w_i^* = \text{elke const. (stuk, do, enz)}$ \Rightarrow const. (do) $= \frac{1}{q}$

$w_{i-1} = \text{do}$ const. (do)

⇒ Do I like, in this way colour

Pat Data Structure

PAT Data Structure
practical Algorithm to retrieve information coded in
A Iophonumeric

Alophanemic

A PAT tree (or PAT array) is a compressed binary tree used to efficiently store and search substring

of text. It is especially useful in string matching, prefix searching and text indexing.

key character:
-Designed for continuous text input

– Designed for extracting substrings called sibling

-works with **unsorted** data
-**inserts** new elements into the **structure** that supports logical sorting

—A binary search is better than *inverted filter* for range searches better than *inverted filter*.

— each position in the input becomes an anchor point.

for a sis ting

A **substring** is a substring starting from a specific character position to the end of the input text.

How it works:
Every position in the input string becomes a starting

2. AVL strings are sorted as leaf nodes in a binary point for a sibling

3. Intermediate nodes store bit positions where strings differ.

- a. The comparison is done bitwise for efficiency.
- b. The tree is compressed to reduce size by skipping redundant bits.

e.g. cat → c A T = 3 bytes

car → c A R = 3 + 6 bytes.

Regular tree compact pat tree.



Binary representation

CAT = 110 100 001
CAR = 110 100 101
CAR = 110 100 R
CAT = 110 100 001
CAT = 110 100 001
CAT = 110 100 001

Signature file structure.
A signature file is a filtering-based data structure used in information retrieval to quickly eliminate non-matching documents before doing a full-text check. It uses bitwise codes (signatures) to represent the presence of terms in documents.

key points

- each word is converted into a fixed-length binary code (word signature) using a hash function.
- the bit positions set to 1 in the signature depend on the hash output.
- all word signatures in a document are OR-ed together to form a block / document signature.

110 100 001

110 100 100

110 100 101

[110 100]

#

110 100 001

110 100 100
110 100 101

110 100 100
110 100 101

R
T

LA T

- during search, a query signature is generated and matched against document signature using bitwise comparison.

Advantages:-

1. Efficient prefix search
2. Compact representation
3. Logical sorting
4. Better than inverted file

creation of signatures

3

Hash function is used for creating signatures.

that happens input data at

Hashing is a process that converts any size string into a fixed-size value or key, typically a string of characters or a numerical value.

```

graph LR
    Input["Input: 'This is going' (document)"] --> HashFunction[hash function]
    HashFunction --> Output["Output: 0001000 (pseudored sign)"]

```

The diagram illustrates a process flow. On the left, the text "Input: 'This is going' (document)" is shown above a downward-pointing arrow. This arrow points to a box labeled "hash function". From the "hash function" box, another downward-pointing arrow leads to the text "Output: 0001000 (pseudored sign)".

Eng- "computer science graduate student study"

word signature

| | | | | |
|------------------------------|---|--------|-------|------|
| computer → 00010100000001010 | } | origin | D 110 | |
| science → 00000000000010010 | | } | 0100 | |
| graduate → 10001000000000100 | | | } | 0100 |
| students → 00000111100001000 | | | | 0100 |

0110
0100
0100
0100

signature → 100101111001
→ 000001110011001000
↓
oring

user → query computer: 1001011100010
 001 010 0000 010 → 1001 011100010
 searching
 o' clock

~~00001 0110 0000~~ 0110 → match found

Advantages:-

- compact structure , —fast filtering , —parallel processing

- works with worm media
 - low-frequency terms.

- False positives: Because unrelated documents may match.

- No ranking: unlike inverted indexes, it doesn't easily support term weighting or ranking.

Collections

Hidden markov model

Hidden Markov Model (HMM) is a statistical model used to represent systems that are probabilistic and sequential where the system being modeled is not directly observable (hidden), but can be inferred through

Observable outputs

HMMs have been successfully applied in:

- speech recognition
 - named entity recognition (NER)
 - optical character recognition (OCR)
 - topic identification.
 - information retrieval.

concept :-
A markov model assumes :
The system is in one of a finite no. of states
which depends only on current state.

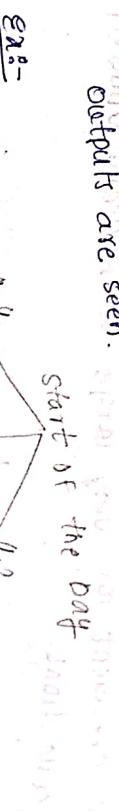
(not full history) - called a first order Markov property.

How it works:-

- The model starts in an initial state.
- It transitions to the next state using transition probabilities (α).

- At each state, it emits an output symbol using output probabilities (β).

- The actual sequence of states is hidden, only outputs are seen.



$$\begin{aligned} \text{event 1} &= \text{Cloudy | sunny} \cdot \text{Cloudy | sunny} \cdot \text{Cloudy | sunny} \\ \text{event 2} &= \text{Cloudy | sunny} \cdot \text{Cloudy | sunny} \cdot \text{Cloudy | sunny} \\ \text{event 3} &= \text{Cloudy | sunny} \cdot \text{Cloudy | sunny} \cdot \text{Cloudy | sunny} \\ \text{event 4} &= \text{Cloudy | sunny} \cdot \text{Cloudy | sunny} \cdot \text{Cloudy | sunny} \\ \text{P}_{\text{event 1}} &= 0.2 \cdot 0.4 \cdot 0.1 = 0.008 \\ \text{P}_{\text{event 2}} &= 0.4 \cdot 0.4 \cdot 0.2 = 0.032 \\ \text{P}_{\text{event 3}} &= 0.2 \cdot 0.4 \cdot 0.2 = 0.016 \\ \text{P}_{\text{event 4}} &= 0.4 \cdot 0.4 \cdot 0.2 = 0.064 \end{aligned}$$

$$\begin{aligned} \text{1. } \text{Cloudy} &\rightarrow \text{sunny} \rightarrow \text{Rainy} \rightarrow \text{Cloudy} = P_1 \\ \text{2. } \text{Cloudy} &\rightarrow \text{Rainy} \rightarrow \text{Cloudy} \rightarrow \text{sunny} = P_2 \\ \text{3. } \text{Sunny} &\rightarrow \text{Cloudy} \rightarrow \text{Rainy} \rightarrow \text{Cloudy} = P_3 \\ \text{4. } \text{Cloudy} &\rightarrow \text{Rainy} \rightarrow \text{Cloudy} \rightarrow \text{Rainy} = P_4 \end{aligned}$$

Applications:-

- Ranking documents
- Sequence modeling
- Language modeling.

Advantages :-

- Handles uncertainty & sequential patterns

- Can model temporal or structured data

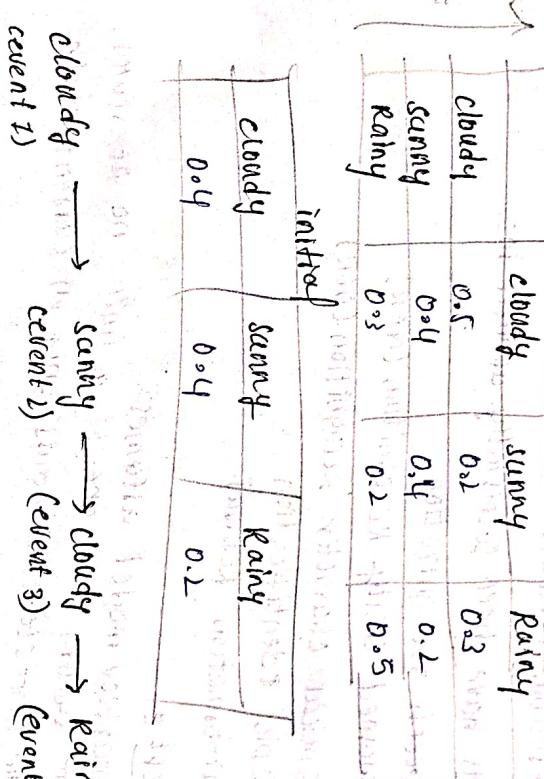
- Useful for ranking, prediction and pattern detection

Limitations :-

- Computationally intensive for large state spaces

- Requires training data

- Assumes Markov property.



Hypertext and XML Data Structures

Information Retrieval

What is Hypertext?

- It is a form of text that enable includes links (hyperlinks) connecting it to other text.

HTML documents, or multimedia resources.

- This structure enables non-linear navigation and is the foundation of the world wide web.

Key: - supports inter-document and intra-document linking

- Helps user browse and navigate through information intuitively.

- Typically represented using HTML.

Information Retrieval (IR)

- Hypertext adds context and connectivity to documents.
- Used in web search engines, digital libraries, and hypermedia systems.

- Links help in ranking.

What is XML? Extensible Markup Language

- XML is a structured document format that defines rules for encoding documents in a way that is both human-readable and machine-readable.

Structure: - XML documents are made of

- Elements (tags like <title>, <author>)
- Attributes
- Hierarchical tree structure.

Why use XML in IR?

- Precise tagging and structuring of information
- Using semantic relationships b/w parts of a document
- Efficient querying using XML-aware search engines (e.g. XPath, XQuery).

Conclusion

- Hypertext helps organize documents non-linearly, enabling web-style navigation.
- XML provides structure, semantics and metadata to documents, making them easier to index, search and retrieve.
- Together, they enrich document representation of retrieval in modern IR systems.

Unit - III

Automatic Indexing

- classes of automatic indexing
 - Statistical indexing
 - Natural language
 - concept indexing
 - Hypertext linkages
- Document and term clustering
 - Introduction to clustering
 - thesaurus generation
 - item clustering
 - hierarchy of clusters.

Automatic indexing

- A method of indexing in which an algorithm is applied by a computer to the title/text of a work to identify and extract words and phrases representing subjects based on entries of index.
- Automatic indexing is the process of assigning keywords or index terms to documents without human intervention, using computational techniques. It helps in organizing and retrieving information efficiently.

Types

| Statistical Indexing | Natural language indexing | concept indexing | Hypertext linkage |
|----------------------|---------------------------|------------------|-------------------|
|----------------------|---------------------------|------------------|-------------------|

- It requires few seconds based on the processor & complexity of algo to generate Indexes.
- It is the process of analyzing an item to extract the information to be permanently kept in index.
- Here we can associate with generation of the searchable data structure associated with an item.

→ Statistical indexing fall into two classes

- weighted
- unweighted

Statistical Indexing

Statistical indexing is a technique used in information retrieval to ranks and retrieve documents based on statistical properties of terms in collection.

- Instead of relying on predefined keywords it analyzer term frequency and distribution patterns to improve search accuracy.
- (a) In statistical indexing is rank items based on their relevance.
- Because of their user are easy to find the most relevant documents when searching.
- The documents are found by normal boolean search and then statistical calculation are performed on the hit file.
- Ranking the output
- Statistical indexing based upon different models. Those are
 - i. probability weighting:-

Assigns weights based on the probability of a term appearing in relevant vs non relevant documents.

4. vector weighting:-

Represents documents as vectors and ranks them using cosine similarity.

3. Simple Term frequency Algo:-

Weights terms based on how often they appear in a document.

a. Frequency of occurrence

To calculate a relevance value for each item. The frequency is used to determine the importance of an item. in relation

b. Document frequency:-

It refers to the no. of documents in a corpus that contain a particular term.

c. Term frequency

It refers to the no. of times a particular term appears in a document.

4. Inverse Document frequency Algo

Gives higher weight to rare terms

5. Signal weighting:-

User signal processing techniques to detect important terms.

6. Discrimination weighting

Assigns higher weight to terms that differentiate relevant from non-relevant documents.

7. Bayesian model:-

Uses Bayes Theorem to calculate the probability of a document being relevant

These weighting techniques are used to improve information retrieval, search engine ranking, document classification.

Natural language

The goal of natural language processing is to use the semantic information in addition to the statistical information to enhance the indexing of the item.

- It improves the precision of searches reducing the no. of false hits or user reviews.
- The information is extracted as a result of processing the language rather than treating each word as independent entity.
- After the processing, the output is generation of phrases that become index to an item.
- More complex analysis generates thematic representation - word phrases generated by natural language processing algorithms. enden to specification and provide another level of disambiguation.

In Natural language

1. Index phrase Generation
2. Natural Language processing.

Index phrase Generation:-

The goal of indexing is to represent the semantic concepts of an item in the information.

Natural language processing

It provides higher-level semantic information identifying relationships b/w concepts.

- The processing tokens in the document are mapped to subject codes.



- Assignment of terms to components, classifying discourse level areas within an item.



- Tent structure, which attempts to identify general assignments of terms to components, classifying the intent of the terms in the tent.



- Identifies interrelationships b/w the concepts.



- Finally assign final weight to the established relationship.

challenges:-

Advantages:-

- Better precision

and complex

- phrase normalization

- requires large knowledge bases of training data

- context-aware

- computationally intensive compared to keyword matching.

system to support finding relevant information.

- single words have conceptual context, these are very easy to find normal info.

- term phrases allow additional specification and focusing of concept to provide better result except non-relevant items.

Concept Indexing

In previous indexing i.e Natural language processing starts with a basis of the terms and extends the information to phrases.

- higher level concepts builds relationship b/w concepts
- concept indexing takes the abstraction a level further.
- it starts with a no. of unlabeled concepts:

e.g: automobile - term.

Concepts are - vehicle, transportation, environment.

fuel, mechanical devices.

vehicle - 0.65

transportation - 0.60

environment - 0.25

fuel - 0.33

Mechanical - 0.15.

concept indexing means indexing a document based on the meaning of words, not just the exact words.

- It focuses on what the document is talking about, not just which words are used.

why do we need it

sometimes, different words mean the same thing

for e.g

• car = automobile

• heart attack" = myocardial infarction

• mobile = "cell phone"

If your search for "automobile", but the document says "car", a basic system may miss it.

But concept indexing will understand both mean the same and show the result.

How it works

- Concept indexing uses tools like
thesaures (like a smart dictionary of related words)
- Ontologies (big knowledge maps of terms & meanings)
- These tools help group words with similar meanings into the same concept.

Advantages :-

- finds similar meaning
- handles synonyms
- Better search accuracy
- Good for smart systems

Drawbacks

- Needs smart tools
- Slower than simple search
- Errors possible.

Hyper linkages (Hypertext linkages)

Hypertext linkages are clickable links that connect one document to another.

In IRS, hypertext linkages help in

- connecting related documents
- providing extra information
- Add potentially improving search results by including linked content.

Why they are important

- Traditional documents are 2D - they only contain text.
- With hypertext, documents become 3D - they contain text connections to other documents.

How can hyperlinks linkages be used in Indexing?

Basic idea :-

Treat the content in linked documents as part of the main document - but with less weight.

That means : If Document A links to Document B then some keywords from B should be added to A's index , but with reduced importance.

Weight formula:

- Direct text = high weight

- Linked content = lower weight

- The closer the link is to relevant text, the stronger its influence.

There are 3 main ways hyperlink is currently used

1. Manual Indexes:

e.g : YAHOO

- People organize websites into categories.

2. Automatic Indexes:-

e.g :- LYCOS, Altavista

- Bots read pages and index words automatically

3. Web crawlers

e.g :- Google, webcrawler

- Bots follow hyperlinks across websites to build index.

Advantages :-

- connects related info

- better search results

- helps navigation

Challenges :-

- Link passing errors

- synonym confusion

- overhead

- link spam.

Introduction to clustering

clustering is the process of grouping similar items or terms together into categories called clusters (or classes)

In Information Retrieval (IR), clustering helps:-

o Group documents discussing the same topic

o Organize terms that have related meanings (thesaurus)

o Improve search effectiveness by finding related results.

clustering useful for organizes data, improves recall, enhances searching.

2. Main uses of clustering

1. Term clustering (thesaurus generation)

- groups similar terms, like synonyms

- improves query expansion.

2. Document clustering

- groups similar documents, even if exact query terms

don't match.

- improves retrieval of similar content.

steps in clustering process

1. Define the domain

- what are we clustering (e.g medical terms, articles, emails)

2. Determine attributes

- choose what part of the document to consider

(title, abstract, body)

3. measure Relationship

- use co-occurrence of words or similarity functions

like cosine similarity.

Document and Term clustering

4. Apply clustering algorithm

- used a method like k-Means or hierarchical clustering to form groups

Challenges in Clustering:

- Ambiguity of language
 - too much recall can reduce precision
 - homographs (same word, different meaning)

Thesaurus Generation

- A thesaurus in information (IR) is collection of related terms (like synonyms, hierarchical or associated terms) used to improve search precision and recall by expanding queries.
- A book or electronic resource that lists words in groups of synonyms and selected concepts.

Types of thesaurus generation

1. Manual clustering

- Human experts manually group words into conceptually similar clusters
- Relies on knowledge of:
 - Domain (e.g. medical, legal)
 - Concordances (word frequencies + occurrence)
 - KWOC (keyword out of context)
 - KWIC (keyword in context)
 - KWAC (keyword and context)
- pros:- High precision
- cons:- Time consuming, hard to scale.

Item: An item could be a document, a web page, an image, a video, or any other piece of information that can be indexed and searched for.

Term: A "term" refers to a significant word

2. Automatic term clustering

Likely idea:- words that appear together in many documents likely belong to the same concept.

Techniques for automatic thesaurus creation:-

1. Complete term relation method / Term Clustering
 - compares every pair of words using similarity measures (like co-occurrence).
 - creates a term - term matrix (word vs word)
 - thresholding: If similarity > threshold, they're grouped.

Clustering algorithms used:-

• C linkage: all terms are similar to each other (high precision)

◦ single link: if one term matches, add it (high recall, low precision)

◦ star: one loose term, add all related terms

◦ string: like a chain of connected terms.

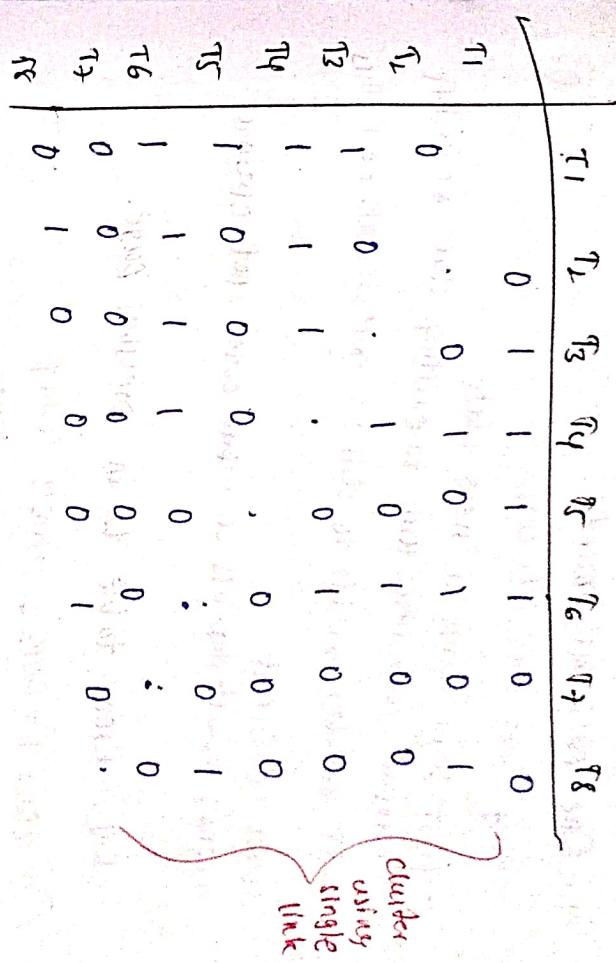
Term clustering in information retrieval is the process of grouping similar terms together based on their occurrence and context within the dataset.

e.g.: Given stems / documents & their extracted terms:

| Item | Term 1 | Term 2 | Term 3 | Term 4 | Term 5 | Term 6 | Term 7 | Term 8 |
|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| stem 1 | 0 | 4 | 0 | 0 | 0 | 2 | 1 | 3 |
| stem 2 | 3 | 1 | 4 | 3 | 1 | 2 | 0 | 1 |
| stem 3 | 3 | 0 | 0 | 0 | 3 | 0 | 3 | 0 |
| stem 4 | 0 | 1 | 0 | 3 | 0 | 0 | 2 | 0 |
| stem 5 | 2 | 2 | 2 | 3 | 1 | 4 | 0 | 2 |

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

| Term | Term 1 | Term 2 | Term 3 | Term 4 | Term 5 | Term 6 | Term 7 | Term 8 |
|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| Term 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Term 2 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Term 3 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| Term 4 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| Term 5 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| Term 6 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| Term 7 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| Term 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |



i.e., Clustering using existing clusters

- start with predefined clusters
- compute centroids coverage position in vector space
- compare centroids coverage position in vector space
- repeat until stable iterative refinement.

The next step is to select a threshold that determines if 2 terms are associated consider

similar enough to each other to be in the same class. The threshold value of 10 is used.

Ref take term as T

- pros :- fewer calculations
- cons :- Need to predefine no. of clusters.

iii. One pass assignments.

- very fast method (O(n))
- first term forms a new class
- compare each new term to existing class centroids
- If similar → adds to class; if not → create new class
- Pros: efficient
- Cons: Result depends on input order, not optimal.

e.g :- from table t_2 in previous page.

```
class 1 = term 1, term 3, term 4
class 2 = term 2, term 6, term 5
class 3 = term 5
class 4 = term 4
```

Note the centroid values used during the one-pass process:-

```
class1(term2, term3) = 0, 2/2, 3/2, 0, 4/2
class2(term1, term3, term4) = 0, 10/3, 3/3, 3/3, 3/3
class3(term1, term6) = 6/2, 3/2, 0/2, 1/2, 6/2.
```

How it works:-

Item clustering uses similarity functions - just like term clustering - but instead of comparing terms across documents, we now compare entire documents using their term vectors.

The similarity b/w item calculated using a vector model, and various clustering algorithms are applied such as:

- clique method
- single link method
- star method
- string method.

examples:-

Advantages

- Improved recall
- Better precision
- query expansion
- semantic search
- supports NLP & AI

Drawbacks

- Time consuming (manual)
- limited domain
- generalization
- Ambiguity & homographs
- overgeneralization

Item clustering (Document clustering)

Item clustering is the grouping of documents/items based on their content similarity. The main goal is to cluster together items that talk about similar topics so that users can:

- find related documents easily.
- expand or refine search results

Applications

- * PES
- * Knowledge Representation
- * Search engines
- * Digital libraries.

Step 1 :- convert document to term vectors

| Term | Doc 1 | Doc 2 | Doc 3 | Doc 4 | Doc 5 |
|-------------|-------|-------|-------|-------|-------|
| Cats | 1 | 0 | 1 | 0 | 0 |
| Dogs | 0 | 1 | 0 | 0 | 0 |
| Python | 0 | 0 | 0 | 1 | 1 |
| Java | 0 | 0 | 0 | 0 | 1 |
| Programming | 0 | 0 | 0 | 1 | 0 |

Step 2 :- Measure Similarity
We compare documents based on shared terms.

- Doc 1 and Doc 3 share "cats"
- Doc 2 and Doc 3 share "dogs"
- Doc 4 and Doc 5 share "python"

Step 3 :- cluster them using a simple single-link clustering. we get documents.

cluster

- cluster A - Doc 1, Doc 2, Doc 3 → all about pets/hanmey
- cluster B - Doc 4, Doc 5 → both about programming

Results :-

- Now if a user searches for "cats", the system may also show Doc 3 (because it's in the same cluster) even if "cats" isn't the top keyword → this improves recall
- If the user searches "Java", they'll see Doc 5 and also Doc 4 because they are clustered

together, helping the user discover related info.

Benefits :-

- enhance navigation in large datasets
- improves search relevance
- enables query expansion

Drawbacks :-

- precision loss
- topic mixing
- computational cost.

Hierarchical clustering

- Hierarchical clustering is a method used to group documents or terms into a tree-like structure of clusters. It shows which documents are closely related, which are broader and how they are nested within each other.

Applications:-

- used in systems like Scatter Gathers
- helps in browsing, filtering & query expansion
- used in yahoo style - category systems.

Challenges of drawbacks :-

- Automatic hierarchy creation especially for terms is error-prone due to language ambiguity
- LL poorly done, can reduce recall or precision
- manual hierarchies are better for users but require effort.

1. Explain briefly about hierarchy of clusters

- It builds hierarchy of clusters.

- This algo starts with all the data points assigned to a cluster of their own.

The two nearest clusters are merged into the same cluster.

A technique to create a tree like structures

Hierarchy of clusters from documents or terms based on similarity.

Two main types:-

1. Agglomerative (Bottom-up); start with individual items merge similar ones.

2. Divisive (Top down); start with one big cluster, split into smaller ones.

* Hierarchical agglomerative clustering (HAC) is widely used in TRS:

Goals or objectives of hierarchical clustering

1. Improve search efficiency
2. enhancing user experience
3. enhancing data organization
4. Reduce information overload
5. Improving relevance of search results

key concepts:-

- Dendograms: tree-like diagrams that show the merging or splitting of clusters.
- centroids: represent the average vector of items in a cluster.
- Ward's Method: uses minimum variance and Euclidean distance to form compact clusters.
- $$d_{ij}^2 = \frac{C(\min_{ik})}{(C(\min_{ik}))} d_{ik}^2$$
$$d_{ik}^2 = \sum_{k=1}^n (\alpha_{ik} - \alpha_{jk})^2$$
 where α_{ik}
- Lance-Williams formula:
General formula for calculating dissimilarity between clusters.
- $$D(C_i, C_k) = \alpha_i D(C_i, C_k) + \alpha_j D(C_j, C_k) + \beta D(C_i, C_j)$$
$$\text{where } \alpha_i, \alpha_j, \beta \text{ are } |\{c_i\}|, |\{c_j\}|, |\{c_i \cup c_j\}|.$$
- ### cluster representation
- Monolithic clusters: focused on a single topic.
- Polythetic clusters: represented using top keywords from multiple topics.
- Hierarchical clustering provides a structured and scalable approach to organizing large document collections, supporting both efficient search and intuitive browsing.

Unit - IV

User search Techniques

- Search statements and Binding
- Similarity measures and Ranking
- Relevance feedback
- Selective Dissemination of Information
- Weighted searches of Boolean systems
- Searching the Internet and Hypertext.

Information Visualization

- Introduction to Information Visualization
- cognition and perception.
- Information visualization technologies.

Search statement are the statements to

of an information need generated by user to specify the concepts they are trying to locate in items.

Original user search techniques are

entity to search statements

Technique 2. Binding Query Binding

Information retrieval system

entity to query -
query to database -
query to query -

Ex:-
Basic search statement:-
"Climate change and Arctic wildlife".

Advanced:- "Climate change AND Arctic Wildlife NOT polar bear"

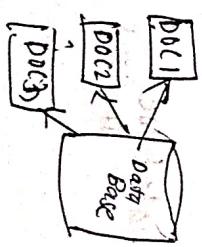
Binding or Search statement Binding (b)

Query Binding

Binding refers to the process of matching the search statement or query with the relevant data within the database.

Search statement binding

Search statement / Query Binding



Search statements and Binding

Search statement (Query)

A search statement is essentially the query that a user inputs into an information retrieval system.

This can include keywords, phrases or even more complex queries involving Boolean operators (like AND, OR, NOT).

The primary purpose of search statement is to define the user's information need as precisely as possible. Since users want to get the most relevant results, the search statement must be well-defined and specific.

Basic search statement:-

- binding is to the vocabulary and past experience of the user.

- The search statement is the users attempt to specify the conditions needed to logically subset the total items space to clarify a subset of items that contains that information needed by the user.

Binding (first level) - How the users own knowledge and vocabulary shape their search query (User specific) - When user creates a search statement they are using their understanding and past experience to form a query that captures what they need from database.

Binding (second level) - The next binding level happens when system specific

- the search system processes the users query
- it translates the query into a standard format that the system can understand (Query processing)

Binding (third/final level) - Finally, the search query is applied to a specific database (database specific)

- adjusting the query based on the specific characteristics of the database like the frequency of certain terms.

Example :-

Input

- i) "find me information on the impact of the oil spills the using vocabulary Alaska on the price of oil" (not used).

Impact oil (petroleum), a statistical system which spills accidents), Alaska binding extracts price (cost, value) to weight assigned to search tokens in processing of tokens, petroleum (.65), soils (.12) terms based upon level accidents (.23), Alaska (.45) inverse document price (.16) cost (.25), value (.10) frequency algo & database

Similarity Measures

- Similarity measures is a function which is used to measure the similarity b/w user query and documents.
- It is possible to derive documents in the order of presumed importance.
- Their function calculates the degree of similarity b/w a pair of text objects.

Similarity measure methods

1. cosine similarity method
2. Jaccards method
3. Dice method.

Using similarity method in information retrieval

It is a measure used in text mining to determine how similar two documents are in content.

documents, for a query of documents based on their content.

$$\text{Similarity } (A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

$\cos \theta$ is the cosine of the angle θ below the horizontal.

-two vectors.
 $A \cdot B = \text{Dot product of } A \text{ and } B$

$$\frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n -a_i b_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

ex-^{er} doc 1 : The quick brown fox

vector 2: The quick brown dog

$$D_1 = [1, 1, 1, 1, 0],$$

$$D_2 = [1, 1, 1, 1, 0],$$

$$A \cdot B = D_1 \cdot D_2 = (1 \cdot 1) + (1 \cdot 1) + (1 \cdot 1) + (1 \cdot 0) + (0 \cdot 1) \\ = 1 + 1 + 1 + 0 + 0 = 3$$

$$|\overrightarrow{OP_1}| = \sqrt{1^2 + 1^2 + 1^2 + 1^2 + 0} = \sqrt{4} = 2$$

$$MP_1 = \sqrt{1^2 + 1^2 + 1^2 + 0^2 + 1^2} = \sqrt{4} = 2$$

Q. Jaccard similarity
It compares the intersection and union of two

self. It is useful when we want to search the similarity b/w 2 documents or an query and document.

formula for Jaccard similarity = $\frac{\text{common terms}}{\text{total terms}}$

q "data", "machine" $p \rightarrow \text{size} = 2$
 q "data", "mining", "machine", "learning"
 "scientist", "intelligence" $p = 6$

3. Die Regelmässigkeit der Verteilung

Dice coefficient is another measure of similarity, but gives more weight to the intersection.

$$\text{formula} = \frac{a \times |A \cap B|}{|A| + |B|}$$

E.g.— using the same sets

$$\bullet |A \cap B| = 2$$

$$\text{Dice} = \frac{2 \times 2}{4 + 4} = \frac{4}{8} = 0.5$$

Role in IRs

- Helps to rank documents from most to least relevant.

- Useful in models like vector space, probabilistic

\Rightarrow Neural IR

- Supports fuzzy searching and partial matching

- Comparison of terms based on TF-IDF

Ranking

\Rightarrow Ranking is the process of sorting document

that match a query in order to display the most relevant results to the user first.

Query - Document matching

Ranking approach

Dataset, Query, Algo

Information system

Information retrieval system

Information retrieval system

Information retrieval system

Information

Information retrieval system

Information retrieval system

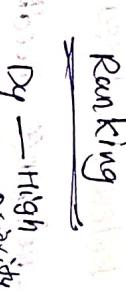
Information retrieval system

Evaluation of ranking

\Rightarrow retrieval results

Document collection

calculating similarity:-



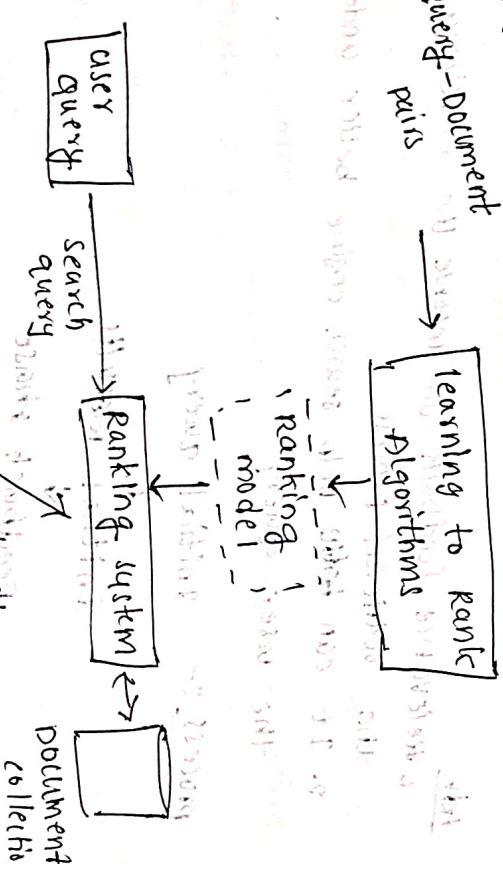
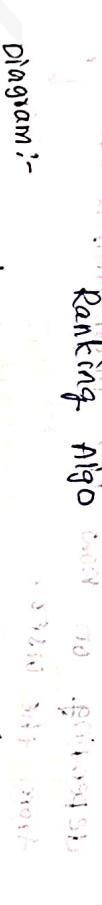
steps :-

1. Query understanding

2. Document matching

3. Relevance scoring

4. Ranking algo



Importance:-

- Improve precision and recall

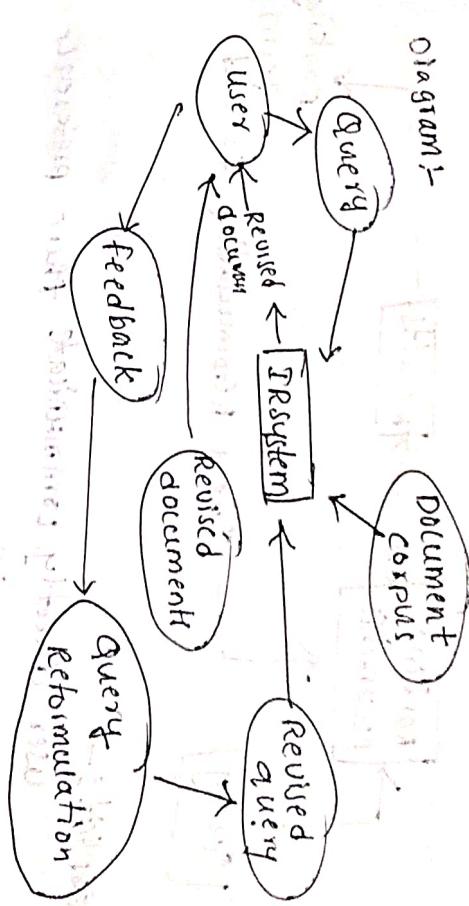
- Enhance user satisfaction

- Helps in relevance-based filtering

- Used in web search engines, digital libraries

Relevance Feedback

- Relevance feedback is a feature of IRS that improves the effectiveness of search queries.
 - It is a mechanism that allows users to provide feedback on the relevance of returned results to their search query. This feedback is then used to refine and improve future search results.
 - Relevance feedback can be explicit or implicit



Rele
t
dback

- Implicit feedback
 - Explicit feedback
 - Pseudo feedback
 - Users are often relevant to provide relevant

25

- Judgement is gathered by observing user actions and behavior during their interaction.

w

- Time spent on a document with subsystems.

x3

- | No. of click. | filling |
|---------------|---------|
| 1 | 1 |
| 2 | 1 |
| 3 | 1 |
| 4 | 1 |
| 5 | 1 |
| 6 | 1 |
| 7 | 1 |
| 8 | 1 |
| 9 | 1 |
| 10 | 1 |
| 11 | 1 |
| 12 | 1 |
| 13 | 1 |
| 14 | 1 |
| 15 | 1 |
| 16 | 1 |
| 17 | 1 |
| 18 | 1 |
| 19 | 1 |
| 20 | 1 |
| 21 | 1 |
| 22 | 1 |
| 23 | 1 |
| 24 | 1 |
| 25 | 1 |
| 26 | 1 |
| 27 | 1 |
| 28 | 1 |
| 29 | 1 |
| 30 | 1 |
| 31 | 1 |
| 32 | 1 |
| 33 | 1 |
| 34 | 1 |
| 35 | 1 |
| 36 | 1 |
| 37 | 1 |
| 38 | 1 |
| 39 | 1 |
| 40 | 1 |
| 41 | 1 |
| 42 | 1 |
| 43 | 1 |
| 44 | 1 |
| 45 | 1 |
| 46 | 1 |
| 47 | 1 |
| 48 | 1 |
| 49 | 1 |
| 50 | 1 |
| 51 | 1 |
| 52 | 1 |
| 53 | 1 |
| 54 | 1 |
| 55 | 1 |
| 56 | 1 |
| 57 | 1 |
| 58 | 1 |
| 59 | 1 |
| 60 | 1 |
| 61 | 1 |
| 62 | 1 |
| 63 | 1 |
| 64 | 1 |
| 65 | 1 |
| 66 | 1 |
| 67 | 1 |
| 68 | 1 |
| 69 | 1 |
| 70 | 1 |
| 71 | 1 |
| 72 | 1 |
| 73 | 1 |
| 74 | 1 |
| 75 | 1 |
| 76 | 1 |
| 77 | 1 |
| 78 | 1 |
| 79 | 1 |
| 80 | 1 |
| 81 | 1 |
| 82 | 1 |
| 83 | 1 |
| 84 | 1 |
| 85 | 1 |
| 86 | 1 |
| 87 | 1 |
| 88 | 1 |
| 89 | 1 |
| 90 | 1 |
| 91 | 1 |
| 92 | 1 |
| 93 | 1 |
| 94 | 1 |
| 95 | 1 |
| 96 | 1 |
| 97 | 1 |
| 98 | 1 |
| 99 | 1 |
| 100 | 1 |
| 101 | 1 |
| 102 | 1 |
| 103 | 1 |
| 104 | 1 |
| 105 | 1 |
| 106 | 1 |
| 107 | 1 |
| 108 | 1 |
| 109 | 1 |
| 110 | 1 |
| 111 | 1 |
| 112 | 1 |
| 113 | 1 |
| 114 | 1 |
| 115 | 1 |
| 116 | 1 |
| 117 | 1 |
| 118 | 1 |
| 119 | 1 |
| 120 | 1 |
| 121 | 1 |
| 122 | 1 |
| 123 | 1 |
| 124 | 1 |
| 125 | 1 |
| 126 | 1 |
| 127 | 1 |
| 128 | 1 |
| 129 | 1 |
| 130 | 1 |
| 131 | 1 |
| 132 | 1 |
| 133 | 1 |
| 134 | 1 |
| 135 | 1 |
| 136 | 1 |
| 137 | 1 |
| 138 | 1 |
| 139 | 1 |
| 140 | 1 |
| 141 | 1 |
| 142 | 1 |
| 143 | 1 |
| 144 | 1 |
| 145 | 1 |
| 146 | 1 |
| 147 | 1 |
| 148 | 1 |
| 149 | 1 |
| 150 | 1 |
| 151 | 1 |
| 152 | 1 |
| 153 | 1 |
| 154 | 1 |
| 155 | 1 |
| 156 | 1 |
| 157 | 1 |
| 158 | 1 |
| 159 | 1 |
| 160 | 1 |
| 161 | 1 |
| 162 | 1 |
| 163 | 1 |
| 164 | 1 |
| 165 | 1 |
| 166 | 1 |
| 167 | 1 |
| 168 | 1 |
| 169 | 1 |
| 170 | 1 |
| 171 | 1 |
| 172 | 1 |
| 173 | 1 |
| 174 | 1 |
| 175 | 1 |
| 176 | 1 |
| 177 | 1 |
| 178 | 1 |
| 179 | 1 |
| 180 | 1 |
| 181 | 1 |
| 182 | 1 |
| 183 | 1 |
| 184 | 1 |
| 185 | 1 |
| 186 | 1 |
| 187 | 1 |
| 188 | 1 |
| 189 | 1 |
| 190 | 1 |
| 191 | 1 |
| 192 | 1 |
| 193 | 1 |
| 194 | 1 |
| 195 | 1 |
| 196 | 1 |
| 197 | 1 |
| 198 | 1 |
| 199 | 1 |
| 200 | 1 |
| 201 | 1 |
| 202 | 1 |
| 203 | 1 |
| 204 | 1 |
| 205 | 1 |
| 206 | 1 |
| 207 | 1 |
| 208 | 1 |
| 209 | 1 |
| 210 | 1 |
| 211 | 1 |
| 212 | 1 |
| 213 | 1 |
| 214 | 1 |
| 215 | 1 |
| 216 | 1 |
| 217 | 1 |
| 218 | 1 |
| 219 | 1 |
| 220 | 1 |
| 221 | 1 |
| 222 | 1 |
| 223 | 1 |
| 224 | 1 |
| 225 | 1 |
| 226 | 1 |
| 227 | 1 |
| 228 | 1 |
| 229 | 1 |
| 230 | 1 |
| 231 | 1 |
| 232 | 1 |
| 233 | 1 |
| 234 | 1 |
| 235 | 1 |
| 236 | 1 |
| 237 | 1 |
| 238 | 1 |
| 239 | 1 |
| 240 | 1 |
| 241 | 1 |
| 242 | 1 |
| 243 | 1 |
| 244 | 1 |
| 245 | 1 |
| 246 | 1 |
| 247 | 1 |
| 248 | 1 |
| 249 | 1 |
| 250 | 1 |
| 251 | 1 |
| 252 | 1 |
| 253 | 1 |
| 254 | 1 |
| 255 | 1 |
| 256 | 1 |
| 257 | 1 |
| 258 | 1 |
| 259 | 1 |
| 260 | 1 |
| 261 | 1 |
| 262 | 1 |
| 263 | 1 |
| 264 | 1 |
| 265 | 1 |
| 266 | 1 |
| 267 | 1 |
| 268 | 1 |
| 269 | 1 |
| 270 | 1 |
| 271 | 1 |
| 272 | 1 |
| 273 | 1 |
| 274 | 1 |
| 275 | 1 |
| 276 | 1 |
| 277 | 1 |
| 278 | 1 |
| 279 | 1 |
| 280 | 1 |
| 281 | 1 |
| 282 | 1 |
| 283 | 1 |
| 284 | 1 |
| 285 | 1 |
| 286 | 1 |
| 287 | 1 |
| 288 | 1 |
| 289 | 1 |
| 290 | 1 |
| 291 | 1 |
| 292 | 1 |
| 293 | 1 |
| 294 | 1 |
| 295 | 1 |
| 296 | 1 |
| 297 | 1 |
| 298 | 1 |
| 299 | 1 |
| 300 | 1 |
| 301 | 1 |
| 302 | 1 |
| 303 | 1 |
| 304 | 1 |
| 305 | 1 |
| 306 | 1 |
| 307 | 1 |
| 308 | 1 |
| 309 | 1 |
| 310 | 1 |
| 311 | 1 |
| 312 | 1 |
| 313 | 1 |
| 314 | 1 |
| 315 | 1 |
| 316 | 1 |
| 317 | 1 |
| 318 | 1 |
| 319 | 1 |
| 320 | 1 |
| 321 | 1 |
| 322 | 1 |
| 323 | 1 |
| 324 | 1 |
| 325 | 1 |
| 326 | 1 |
| 327 | 1 |
| 328 | 1 |
| 329 | 1 |
| 330 | 1 |
| 331 | 1 |
| 332 | 1 |
| 333 | 1 |
| 334 | 1 |
| 335 | 1 |
| 336 | 1 |
| 337 | 1 |
| 338 | 1 |
| 339 | 1 |
| 340 | 1 |
| 341 | 1 |
| 342 | 1 |
| 343 | 1 |
| 344 | 1 |
| 345 | 1 |
| 346 | 1 |
| 347 | 1 |
| 348 | 1 |
| 349 | 1 |
| 350 | 1 |
| 351 | 1 |
| 352 | 1 |
| 353 | 1 |
| 354 | 1 |
| 355 | 1 |
| 356 | 1 |
| 357 | 1 |
| 358 | 1 |
| 359 | 1 |
| 360 | 1 |
| 361 | 1 |
| 362 | 1 |
| 363 | 1 |
| 364 | 1 |
| 365 | 1 |
| 366 | 1 |
| 367 | 1 |
| 368 | 1 |
| 369 | 1 |
| 370 | 1 |
| 371 | 1 |
| 372 | 1 |
| 373 | 1 |
| 374 | 1 |
| 375 | 1 |
| 376 | 1 |
| 377 | 1 |
| 378 | 1 |
| 379 | 1 |
| 380 | 1 |
| 381 | 1 |
| 382 | 1 |
| 383 | 1 |
| 384 | 1 |
| 385 | 1 |
| 386 | 1 |
| 387 | 1 |
| 388 | 1 |
| 389 | 1 |
| 390 | 1 |
| 391 | 1 |
| 392 | 1 |
| 393 | 1 |
| 394 | 1 |
| 395 | 1 |
| 396 | 1 |
| 397 | 1 |
| 398 | 1 |
| 399 | 1 |
| 400 | 1 |
| 401 | 1 |
| 402 | 1 |
| 403 | 1 |
| 404 | 1 |
| 405 | 1 |
| 406 | 1 |
| 407 | 1 |
| 408 | 1 |
| 409 | 1 |
| 410 | 1 |
| 411 | 1 |
| 412 | 1 |
| 413 | 1 |
| 414 | 1 |
| 415 | 1 |
| 416 | 1 |
| 417 | 1 |
| 418 | 1 |
| 419 | 1 |
| 420 | 1 |
| 421 | 1 |
| 422 | 1 |
| 423 | 1 |
| 424 | 1 |
| 425 | 1 |
| 426 | 1 |
| 427 | 1 |
| 428 | 1 |
| 429 | 1 |
| 430 | 1 |
| 431 | 1 |
| 432 | 1 |
| 433 | 1 |
| 434 | 1 |
| 435 | 1 |
| 436 | 1 |
| 437 | 1 |
| 438 | 1 |
| 439 | 1 |
| 440 | 1 |
| 441 | 1 |
| 442 | 1 |
| 443 | 1 |
| 444 | 1 |
| 445 | 1 |
| 446 | 1 |
| 447 | 1 |
| 448 | 1 |
| 449 | 1 |
| 450 | 1 |
| 451 | 1 |
| 452 | 1 |
| 453 | 1 |
| 454 | 1 |
| 455 | 1 |
| 456 | 1 |
| 457 | 1 |
| 458 | 1 |
| 459 | 1 |
| 460 | 1 |
| 461 | 1 |
| 462 | 1 |
| 463 | 1 |
| 464 | 1 |
| 465 | 1 |
| 466 | 1 |
| 467 | 1 |
| 468 | 1 |
| 469 | 1 |
| 470 | 1 |
| 471 | 1 |
| 472 | 1 |
| 473 | 1 |
| 474 | 1 |
| 475 | 1 |
| 476 | 1 |
| 477 | 1 |
| 478 | 1 |
| 479 | 1 |
| 480 | 1 |
| 481 | 1 |
| 482 | 1 |
| 483 | 1 |
| 484 | 1 |
| 485 | 1 |
| 486 | 1 |
| 487 | 1 |
| 488 | 1 |
| 489 | 1 |
| 490 | 1 |
| 491 | 1 |
| 492 | 1 |
| 493 | 1 |
| 494 | 1 |
| 495 | 1 |
| 496 | 1 |
| 497 | 1 |
| 498 | 1 |
| 499 | 1 |
| 500 | 1 |
| 501 | 1 |
| 502 | 1 |
| 503 | 1 |
| 504 | 1 |
| 505 | 1 |
| 506 | 1 |
| 507 | 1 |
| 508 | 1 |
| 509 | 1 |
| 510 | 1 |
| 511 | 1 |
| 512 | 1 |
| 513 | 1 |
| 514 | 1 |
| 515 | 1 |
| 516 | 1 |
| 517 | 1 |
| 518 | 1 |
| 519 | 1 |
| 520 | 1 |
| 521 | 1 |
| 522 | 1 |
| 523 | 1 |
| 524 | 1 |
| 525 | 1 |
| 526 | 1 |
| 527 | 1 |
| 528 | 1 |
| 529 | 1 |
| 530 | 1 |
| 531 | 1 |
| 532 | 1 |
| 533 | 1 |
| 534 | 1 |
| 535 | 1 |
| 536 | 1 |
| 537 | 1 |
| 538 | 1 |
| 539 | 1 |
| 540 | 1 |
| 541 | 1 |
| 542 | 1 |
| 543 | 1 |
| 544 | 1 |
| 545 | 1 |
| 546 | 1 |
| 547 | 1 |
| 548 | 1 |
| 549 | 1 |
| 550 | 1 |
| 551 | 1 |
| 552 | 1 |
| 553 | 1 |
| 554 | 1 |
| 555 | 1 |
| 556 | 1 |
| 557 | 1 |
| 558 | 1 |
| 559 | 1 |
| 560 | 1 |
| 561 | 1 |
| 562 | 1 |
| 563 | 1 |
| 564 | 1 |
| 565 | 1 |
| 566 | 1 |
| 567 | 1 |
| 568 | 1 |
| 569 | 1 |
| 570 | 1 |
| 571 | 1 |
| 572 | 1 |
| 573 | 1 |
| 574 | 1 |
| 575 | 1 |
| 576 | 1 |
| 577 | 1 |
| 578 | 1 |
| 579 | 1 |
| 580 | 1 |
| 581 | 1 |
| 582 | 1 |
| 583 | 1 |
| 584 | 1 |
| 585 | 1 |
| 586 | 1 |
| 587 | 1 |
| 588 | 1 |
| 589 | 1 |
| 590 | 1 |
| 591 | 1 |
| 592 | 1 |
| 593 | 1 |
| 594 | 1 |
| 595 | 1 |
| 596 | 1 |
| 597 | 1 |
| 598 | 1 |
| 599 | 1 |
| 600 | 1 |
| 601 | 1 |
| 602 | 1 |
| 603 | 1 |
| 604 | 1 |
| 605 | 1 |
| 606 | 1 |
| 607 | 1 |
| 608 | 1 |
| 609 | 1 |
| 610 | 1 |
| 611 | 1 |
| 612 | 1 |
| 613 | 1 |
| 614 | 1 |
| 615 | 1 |
| 616 | 1 |
| 617 | 1 |
| 618 | 1 |
| 619 | 1 |
| 620 | 1 |
| 621 | 1 |
| 622 | 1 |
| 623 | 1 |
| 624 | 1 |
| 625 | 1 |
| 626 | 1 |
| 627 | 1 |
| 628 | 1 |
| 629 | 1 |
| 630 | 1 |
| 631 | 1 |
| 632 | 1 |
| 633 | 1 |
| 634 | 1 |
| 635 | 1 |
| 636 | 1 |
| 637 | 1 |
| 638 | 1 |
| 639 | 1 |
| 640 | 1 |
| 641 | 1 |
| 642 | 1 |
| 643 | 1 |
| 644 | 1 |
| 645 | 1 |
| 646 | 1 |
| 647 | 1 |
| 648 | 1 |
| 649 | 1 |
| 650 | 1 |
| 651 | 1 |
| 652 | 1 |
| 653 | 1 |
| 654 | 1 |
| 655 | 1 |
| 656 | 1 |
| 657 | 1 |
| 658 | 1 |
| 659 | 1 |
| 660 | 1 |
| 661 | 1 |
| 662 | 1 |
| 663 | 1 |
| 664 | 1 |
| 665 | 1 |
| 666 | 1 |
| 667 | 1 |
| 668 | 1 |
| 669 | 1 |
| 670 | 1 |
| 671 | 1 |
| 672 | 1 |
| 673 | 1 |
| 674 | 1 |
| 675 | 1 |
| 676 | 1 |
| 677 | 1 |
| 678 | 1 |
| 679 | 1 |
| 680 | |

It
en

- is a method of providing feedback about the relevance of search results by directly asking

use

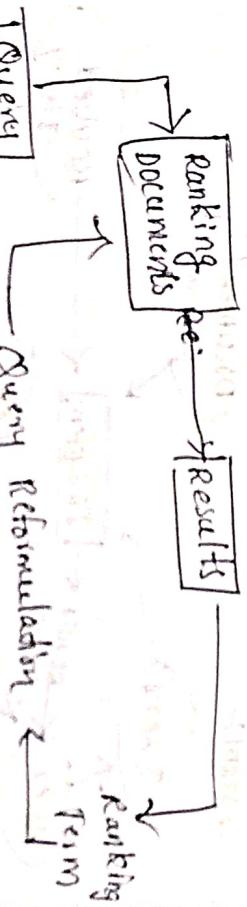
- Ran the original review
 - ranking the document

end

- ranking the document
 - Assume top documents pseudo relevant

- Construct new query representations
 - Run it, compare rankings.

Selective Dissemination of Information (SDI)

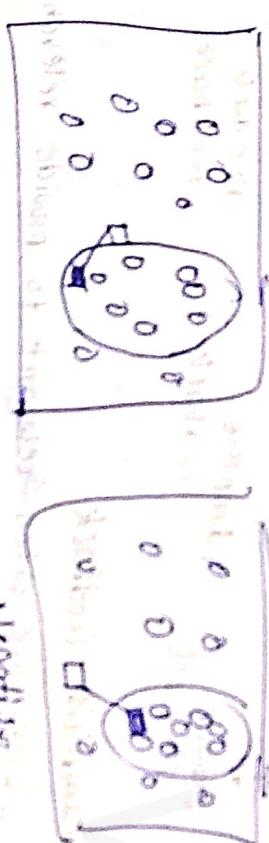


It is a service that automatically delivers information to users based on their interests and preferences.

Explicit:-

User directly communicate their preference

- Positive
- Negative



Implicit positive

'The items that are selected'

'The oval represent the items that are selected from the query.'

'The solid box is logically where the query initially.'

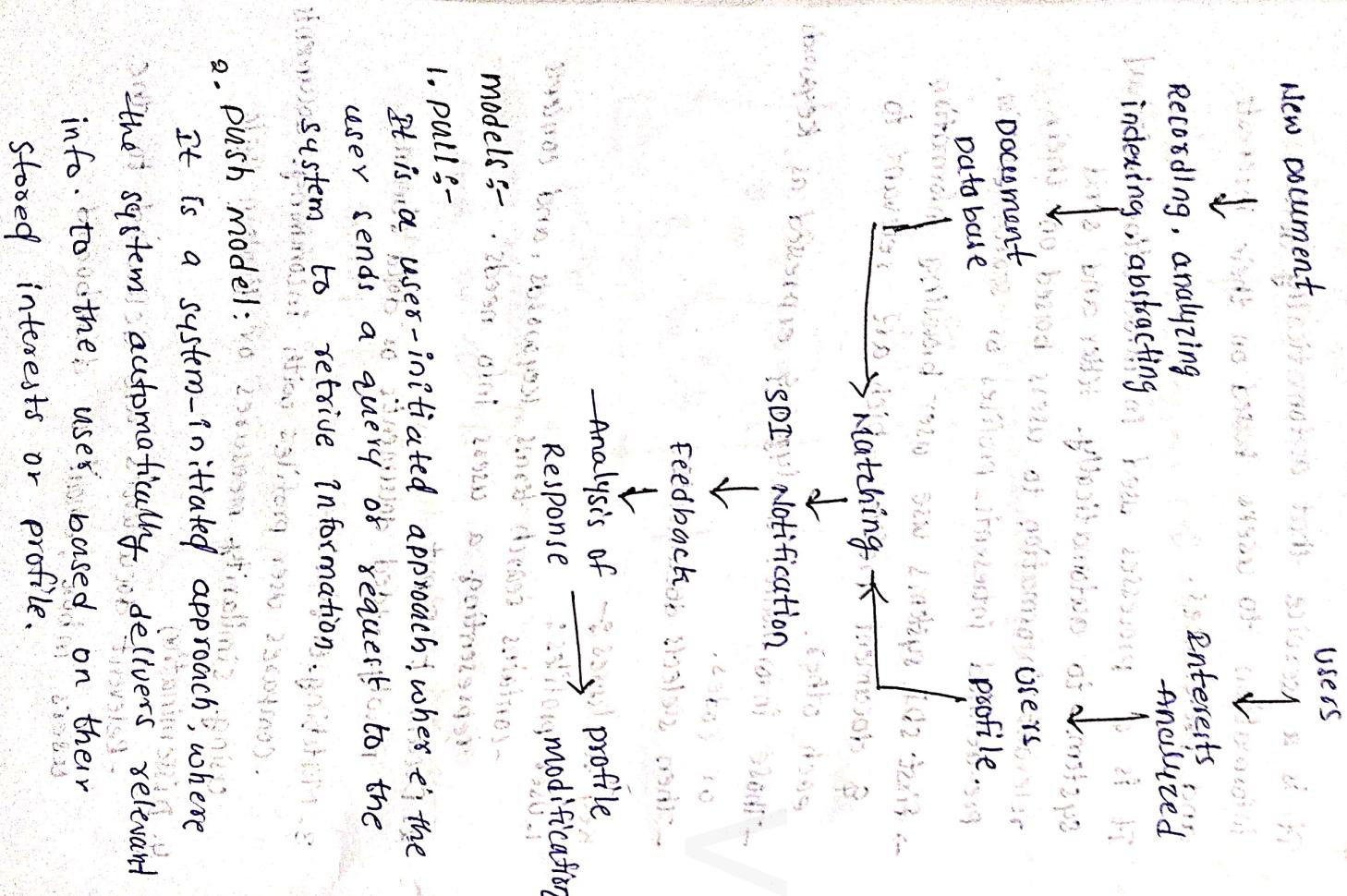
'The hollow box is the query modified for the relevance feedback.'

'(modified relevance feedback)'

Key features:-

1. User profiles:
 - contains search terms, keywords, and concepts representing a user's info needs.
2. Document Input:
 - newly added documents or data streams.
3. Matching Engine:
 - compares user profiles with incoming document using similarity measures or Boolean logic
4. Dissemination:
 - relevant documents are delivered to the user's inbox, email or dashboard.

SDE Diagram



Benefits of SDE

- Increase user satisfaction & engagement
- Reduce information overload via filtering
- Improve decision-making by providing timely info
- Limitations**

 - Needs accurate and updated user profiles
 - Risk of missing relevant data if profiles are too narrow
 - Complex to implement in multimedia or large scale systems.

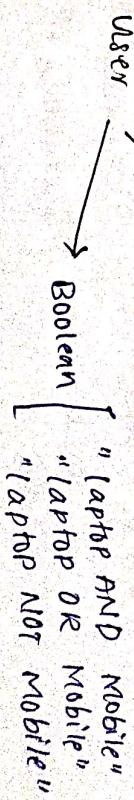
Weighted Searches of Boolean Systems

Boolean systems are classical information retrieval (CR) models that use Boolean logic operators **AND**, **OR**, and **NOT** to retrieve relevant documents based on exact term matches.

- Pulls** -
It is a user-initiated approach, where the user sends a query or request to the system to retrieve information.
 - Push model**:
It is a system-initiated approach, where the system automatically delivers relevant info. to the user based on their stored interests or profile.
- A weighted Boolean search assigns importance weights to search terms in a Boolean query to influence the ranking of retrieved documents.

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

The system automatically delivers relevant info. to the user based on their stored interests or profile.



Query \leftarrow AND operation. results must contain

LAPTOP AND MOBILE \rightarrow IRS \rightarrow Laptop and Mobile

Only few results will contain both, cause very

LAPTOP OR MOBILE \rightarrow IRS \rightarrow results will contain either laptop or mobile, causes too many irrelevant results.

LAPTOP AND MOBILE \rightarrow Laptop, AND Mobile

LAPTOP AND MOBILE \rightarrow IRS \rightarrow either laptop or mobile, causes too many irrelevant results.

Ex:- query computer program AND sale

D_1 contains all sales related to computer

D_2 contains all sales related to program

D_3 contains all sales related to sale

D_4 contains all sales related to computer program

D_5 contains all sales related to computer and sale

D_6 contains all sales related to program and sale

D_7 contains all sales related to computer and program

D_8 contains all sales related to computer, program and sale

step 1 - query

$Q_{1\text{ optional}} = \text{computer} \cup \text{program}$

$Q_{2\text{ optional}} = \text{cost}_{0.75} \text{ AND sale}_{1.0}$

using AND operation, result is

selected via step 1

Q. Step 2 :- strict Interpretation = IRS selection

- $Q_1 = \text{computer} \text{ OR } \text{program}$
- $Q_2 = \text{cost}_{0.75} \text{ AND sale}$

$Q_1 \text{ optional} = (D_1, D_2, D_3, D_4, D_5, D_6, D_7)$

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

$Q_1 \text{ strict interpretation} = (D_1, D_2, D_3, D_4, D_5, D_6, D_7)$

$Q_2 \text{ optional} = (D_1, D_2, D_3, D_4, D_5, D_6) \rightarrow [1.0, 99] = 2 \text{ doc}$

$Q_1 \text{ invariant} = (D_8) \rightarrow [1.0, 5] = 1 \text{ doc}$

Step 3 :- Invariant set (use documents)

- $Q_1 \text{ optional} = (D_1, D_2, D_3, D_4, D_5, D_6) \rightarrow [1.0, 99] = 2 \text{ doc}$
- $Q_1 \text{ optional} = (D_1, D_2, D_3, D_4, D_5, D_6) \rightarrow [1.0, 99] = 2 \text{ doc}$
- $Q_2 \text{ optional} = (D_1, D_2) \rightarrow [1.0, 2.5] = 2 \text{ doc}$
- $Q_2 \text{ optional} = (D_1, D_2) \rightarrow [1.0, 2.5] = 2 \text{ doc}$
- $Q_1 \text{ invariant} = (D_8) \rightarrow [1.0, 5] = 1 \text{ doc}$
- $Q_1 \text{ invariant} = (D_8) \rightarrow [1.0, 5] = 1 \text{ doc}$

Step 5 :- Calculate centroid w.r.t to invariant set.

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

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

$$1/3 (4+0+0, 0+6+4, 2+4+6, 4+6+4)$$

$$= 1/3 (4, 10, 12, 14)$$

step 6 :- calculate similarity between centroid

$$(Q_1 \text{ optional} = (D_1, D_2, D_3, D_4, D_5, D_6)) = 2 \text{ doc}$$

$$Q_2 \text{ optional} = (D_1, D_2) = 1 \text{ doc}$$

$$\text{Centroid } (Q_1) = (D_1) = (C_1, 2, 0, D_2)$$

$$\text{Centroid } (Q_2) = (C_2, D_4, D_5) = \frac{1}{3}(4, 10, 12, 14)$$

$$\text{similarity } (\text{Centroid } Q_1, D_1) = (0+4+0+0) = 4$$

$$\text{similarity } (\text{Centroid } Q_1, D_2) = (0+4+0+0) = 4$$

$$\text{similarity } (\text{Centroid } Q_1, D_3) = (16+0+0+0) = 16$$

$$\text{similarity } (\text{Centroid } Q_1, D_4) = (0+12+0+0) = 12$$

$$\text{similarity } (\text{Centroid } Q_1, D_5) = (0+8+0+0) = 8$$

$$\text{similarity } (\text{Centroid } Q_1, D_6) = (24+0+0+24) = 24$$

$$\text{similarity } (\text{Centroid } Q_2, D_1) = \frac{1}{3}(0+4+0+0) = \frac{4}{3}$$

$$\text{similarity } (\text{Centroid } Q_2, D_2) = \frac{1}{3}(16+10+0+0) = \frac{26}{3}$$

$$\text{similarity } (\text{Centroid } Q_2, D_3) = \frac{1}{3}(0+12+0+0) = \frac{12}{3} = 4$$

$$\text{similarity } (\text{Centroid } Q_2, D_4) = \frac{1}{3}(0+8+0+0) = \frac{8}{3}$$

$$\text{similarity } (\text{Centroid } Q_2, D_5) = \frac{1}{3}(24+0+0+24) = \frac{48}{3} = 16$$

$$\text{similarity } (\text{Centroid } Q_2, D_6) = \frac{1}{3}(24+0+0+0) = \frac{24}{3} = 8$$

$$\text{similarity } (\text{Centroid } Q_1, Q_2) = \frac{1}{3}(12+16+0+0) = \frac{28}{3}$$

$$\text{similarity } (\text{Centroid } Q_1, Q_3) = \frac{1}{3}(12+16+0+0) = \frac{28}{3}$$

$$\text{similarity } (\text{Centroid } Q_1, Q_4) = \frac{1}{3}(12+16+0+0) = \frac{28}{3}$$

$$\text{similarity } (\text{Centroid } Q_1, Q_5) = \frac{1}{3}(12+16+0+0) = \frac{28}{3}$$

$$\text{similarity } (\text{Centroid } Q_1, Q_6) = \frac{1}{3}(12+16+0+0) = \frac{28}{3}$$

$$\text{similarity } (\text{Centroid } Q_2, Q_3) = \frac{1}{3}(12+16+0+0) = \frac{28}{3}$$

$$\text{similarity } (\text{Centroid } Q_2, Q_4) = \frac{1}{3}(12+16+0+0) = \frac{28}{3}$$

$$\text{similarity } (\text{Centroid } Q_2, Q_5) = \frac{1}{3}(12+16+0+0) = \frac{28}{3}$$

$$\text{similarity } (\text{Centroid } Q_2, Q_6) = \frac{1}{3}(12+16+0+0) = \frac{28}{3}$$

Introduction to Information Visualization

Information visualization is the practice of representing data and information in graphical or visual format, enabling users to see, understand, and interact with large and complex data sets more effectively.

- Benefits: - Reduces cognitive load by visualizing complex data sets and making them easier to understand and interact with.
- Reduces time to interpret large result sets.
- Improves document-term relationships.
- Helps users refine queries.
- Makes information retrieval, interactive & engaging.
- Improves efficiency.

Cognition and Perception

Cognition: process of thinking, understanding, learning

Perception: process of thinking, understanding, learning

It refers to the mental process involved in

acquiring and understanding knowledge, including attention, memory, problem-solving etc.

In visualization, cognition is concerned with how users:

- Interpret visual structures
- Retain information
- Make decisions based on visual cues.

In the context of information visualization,

cognition and perception play a central role

in designing visual interfaces that help

users process, understand, and analyse

information effectively.

These concepts come from cognitive psychology

and human-computer interaction (HCI) and

guide how visual representation of data

are created in IRS.

Perception :-

how we see

It is ability to see, hear or become aware

of something through the senses, especially

sight. It affects how user notice, group

or prioritize visual information.

Perception allows users to:

- Recognize shapes, colors, patterns

- Detect differences in size, orientation & motion

- Understand groupings in visual data.

Aspects of visualization process

ex:- para A → with text

para B → text with key underline

para C → the key points coloured.

and also by using WIMP

WIMP stands for windows, icons, menus, pointers.

4. Aspects of the visualization process

a. Preattentive processing

- a low-level, unconscious visual process

- detects borders, shapes & color changes almost

- useful for grouping in document clustering or highlighting instantly

b. Object rotation & symmetry

- rotated objects are harder to recognise

- vertically aligned items are easier to process

c. Color usage

key for classification & emphasis

attribute: hue, saturation, lightness.

Information Visualization Technologies

These are systems and techniques designed

to help users understand large volume of

search results by visually representing

information.

Main Goals :-

i. Document clustering

- visually group documents based on content

- e.g. showing documents as clusters in a scatter plot or 3D map.

2. Search statement analysis:

- help users understand why documents were retrieved especially when complex ranking or query expansion is used.

e.g.: - Visualizing term contributions using bar charts or heatmaps.

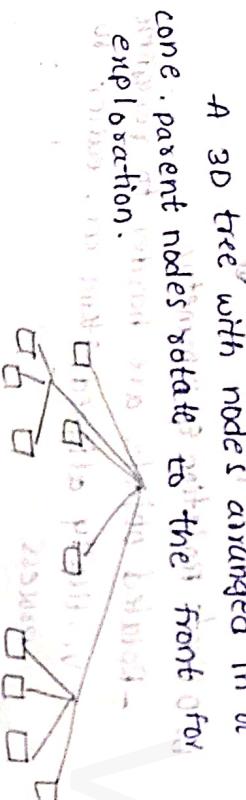
Bar charts show individual term contributions.

Heatmaps represent document similarity based on context.

Visualization Techniques & Tools

A. Hierarchical & 3D structures:

- cone tree



- A 3D tree with nodes arranged in a cone. Parent nodes rotate to the front of exploration.

- Perspective walls
- shows the focus area in the center with sides faded into background.
- Tree Maps: space-filling technique using nested rectangles based on data hierarchy



B. Clustering & spatial layouts:

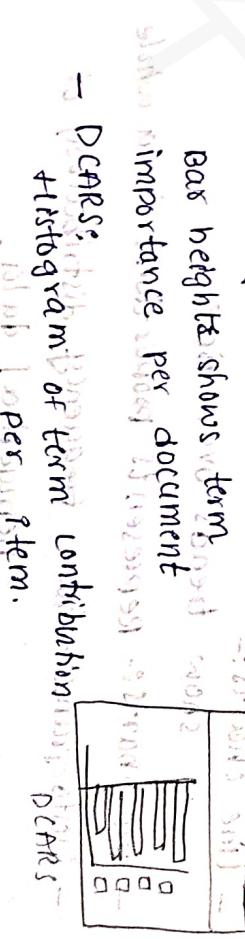
- Scatterplots plot documents in 2D or 3D based on context similarity or term scatterplots

- Semantic landscapes & hills represent topic elevation, valleys & hills represent topic density.

c. Term relationship & query effect.

- envision system: displays document relevance in scatterplots.

- Veerasamy - Belkin Bars

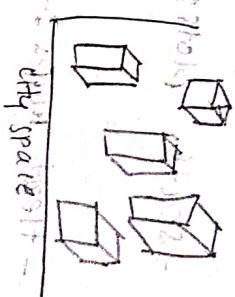


- DCArs histogram of term contribution per document.

D. Cityspace & Geometric views:

- city's pale: uses buildings

- skyscrapers: builds concept importance in 3D.



- Webbook: mimics a virtual book navigation library books where documents are borrowed like flipping a real book pages.

- allows to help refine a query by adding to it.

E. Statistical visualization.

- **counts** :- plots term frequencies over time.
- **Cross field matrix** :- 2D matrix showing intersections of two fields distributions.
- **Title bags** :-
 - shows term distribution within a document.
- **Basic visualization techniques** :-
 - **Bar charts** :- compares quantity across categories.
 - **Line charts** :-
 - show trends over time.
 - **Pie charts** :- represents proportions in a whole.
 - **Histograms** :- shows frequency distribution of numerical data.
 - **Scatter plots** :- shows relationships / correlations between two variables.
 - **Heat maps** :- show intensity or density over a map.
 - **Tree maps** :- show hierarchical data using nested rectangles.
 - **Dendograms** :- tree diagrams for showing clustering or taxonomy.

Unit - V

Text search Algorithms:

- Introduction to Text search Techniques
- Software Text search Algorithms
- Hardware Text search systems.

Multimedia Information Retrieval

- Spoken language Audio retrieval
- Non-speech Audio retrieval
- Graph retrieval
- Imagery retrieval
- Video retrieval.

Introduction to Text Search Techniques

They are used in Information Retrieval systems to efficiently locate relevant documents based on user's query. These techniques aim to match patterns (query terms) within a large collection of text and retrieve items and that satisfy the user's search conditions.

- Quickly identify relevant items
- eliminate irrelevant ones
- Improve search performance and user satisfaction.

In a textual database three classical text retrieval techniques have been defined.

Those techniques are:

1. Streaming (full text scanning)
 - scans the entire database to find matches
 - Useful for exact matches, stop words, or highlighting terms in results.

2. Inverted Index (Inverted Index)

- Builds a word-to-document index

- enables fast look up by skipping non-matching documents.

3. Multi Attribute Retrieval

- searches across multiple fields like title, author and body.
- provides flexible search capabilities.

- First we use Index mechanism for searching.
Text streaming of text was frequently found in the system as search technique.

- Generally query entered by one user. In fact scanning system one or more users enter queries.

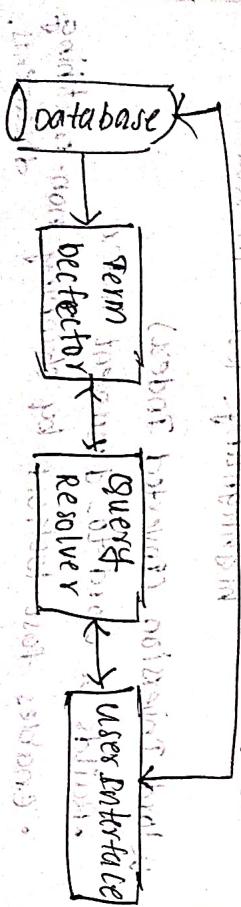
- Earlier systems used hardware based text streamers due to CPU limitations. But modern systems mostly rely on software based search due to increased computing power.

Drawbacks of early techniques

- manual - slow, error-prone
- No partial matching or ranking
- hardware = costly, outdated

Example of limits of index searches are

- search for stop words
- search for exact matches when stemming is performed
- search for terms that contain both leading and trailing "don't care's"



- Term detector is a specially contains all of the terms being searched. It also contains logic to the items. It will give output to the query resolver.

- It process finally. It accept search statements from the user to extract the logic to search terms.

- It accepts results from the detector & determines which queries are satisfied by the item & possibility.

Finally it pass the info to the user interface that will continually updating search status.

These are two approaches to the data beam.

1. Complete database is being sent to the detector.
2. Random retrieved items are being passed to the detectors.

Many of these, like text searches use finite automata as a basis for their algorithms.
I - set of IP symbols from alphabet
S - set of states
P - set of productions
So - initial state
Sr - One or more final state.

Software Text Search Algorithms

- Software Text search algorithms are software based techniques used to locate specific patterns (search terms) within some body of text.
- These are typically used when full text is available (in memory) and no hardware support is required (by function).
- These are four major techniques associated with software text search.
 - 1. KMP
 - 2. Rabin-Karp
 - 3. Aho-Corasick
 - 4. Boyer-Moore
- KMP is a powerful pattern matching algorithm that reduces time complexity from $O(nxm)$ to $O(n+m)$ using the LPS table.
- It's efficient for searching long patterns in large text.
- Avoids rechecking matched characters using a partial match table (prefix function).
- Pre processing takes $O(m)$, matching takes $O(n)$.
- Time Complexity: $O(n+m)$
- Good for predictable and long patterns searches.

a b c d a b c
pattern
abc

suffix :- a, ab, abc, abcd

puttin

string :- a b a b c a b c a b a b d
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15
lps :- a b c d e c i a b c d a b c
0 0 0 0 1 2 0 1 2 0 1 2 0 1 2 0

π or lps

P_1 : a b c d a b e a b b
0 0 0 0 1 2 0 0 1 2 0
 P_2 : a b c d e c i a b c
0 0 0 0 1 2 0 1 2 0 1 2 0 1 2 0

steps :-
1. Build LPS table (longest prefix suffix)

2. Perform search.

- i. Start with two pointers i for text and j for pattern.
- ii. Compare $text[i:j]$ with $pattern[i:j]$. If they match, move both $i+1$ & $j+1$ to next test.
- iii. If j reaches end of pattern
 - pattern found at position $i-j$
 - Reset $j = lps[i-j]$ to look for more matches
 - If mismatch,

- If no backtracking using LPS $\rightarrow j = lps[i-j]$
- If $j=0$ move $i+1$ only (start next search)

a. Brute-force algorithm:-

- compares pattern with the text character by character.

- If mismatch occurs shift by one and repeat

- Time complexity : worst case $\rightarrow O(n \times m)$

- simple but inefficient.

- It is a simplest way to search for pattern in text.

Ex:- Text = ababcaabcab

Pattern = abc

| | | | | | | | | | | | | | | | |
|----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | a | b | c | a | b | c | a | b | c | a | b | c | a | b | c |
| 1 | | | | | | | | | | | | | | | |
| 2 | | | | | | | | | | | | | | | |
| 3 | | | | | | | | | | | | | | | |
| 4 | | | | | | | | | | | | | | | |
| 5 | | | | | | | | | | | | | | | |
| 6 | | | | | | | | | | | | | | | |
| 7 | | | | | | | | | | | | | | | |
| 8 | | | | | | | | | | | | | | | |
| 9 | | | | | | | | | | | | | | | |
| 10 | | | | | | | | | | | | | | | |
| 11 | | | | | | | | | | | | | | | |
| 12 | | | | | | | | | | | | | | | |
| 13 | | | | | | | | | | | | | | | |
| 14 | | | | | | | | | | | | | | | |

→ welcome to LAB mug

3. Boyer-Moore Algorithm

It is an efficient pattern matching algorithm

that searches from right to left and uses

smart skipping to avoid unnecessary comparisons

steps:-

1. first construct 'Bad Match Table'

2. compare right most character of pattern with given string based on value of bad match table.

3. If mismatch then shift the pattern to the right position corresponding to the value.

- extremely efficient in practice

- Bestcase : sub linear time, worst case

worst case $O(nm)$

- most effective for long patterns, and large texts.

- faster than KMP

→ ~~Worst case O(nm) due to hash collisions~~

Ex:- Text = WELCOME TO LAB MUG

Bad match table

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| L | A | B | M | U | G | T | O | L | E | W | C | F | E | N |

length - index - 1
 $L = 6 - 0 - 1 = 5$

→ W E L C O F E T O L A B M U G

* not matched
 * value of m=2
 → W E L C O F E T O L A B M U G

* again shift
 2 steps right

| | | | | | |
|---|---|---|---|---|---|
| L | A | B | M | U | G |
| L | A | B | M | U | G |

→ Then it is matched. ok now we just have to move the pattern to the right position.

4. Rabin-Karp Algorithm

- uses hashing for pattern and substring comparison.

- good for multiple pattern search.

- expected time: $O(n)$, but worst case can degrade to $O(n \times m)$ (due to hash collisions)

- often used in plagiarism detection or searching multiple patterns.

Hardware Text search systems

While software based text search algorithms are widely used, they face performance

- Limitations when dealing with multiple documents
- Processing many search terms simultaneously
- Dealing with large-scale text databases
- Or facing I/O speed bottlenecks.

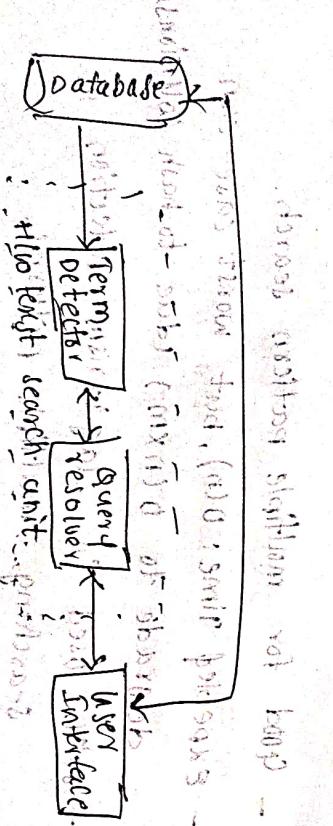
To overcome these, hardware text search systems were developed to

- Offload search processing from the CPU.
- Improve speed and scalability
- Eliminate the need for index structures.

→ It's a specialized machine that performs searches on text data.

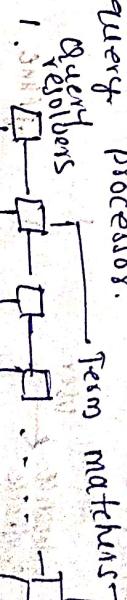
- first machine to perform the searches & pass the results to the main computer which supported the user interface & retrieved hits.

Architecture



- One of the earliest I/Os text searching search unit was the rapid search. It developed by General electric.
- The machine consisted of a special purpose search unit where a single query was passed against a magnetic tape containing the documents.
- The next developed associative file processor (Afp)
- DST following a different approach developed by high speed text search (HST)
- GE redesigned their rapid search machine into the GESCAN unit
- Here uses Text Array Processor (TAP)
- An array of four to 128 query processor
- The text is loaded into the cache & searched by the query processors.
- Each query processor is independent and can be loaded at any time.
- A complete query is handled by each query processor.

query resolution step by step -



N



A query processor works two operations in parallel

1. \rightarrow matching query terms to input text
2. \rightarrow Boolean logic resolution

Term matching is performed by a series of characters cells each containing one character.

Multimedia Information Retrieval (MMIR) also refers to searching multimedia information retrieval system.

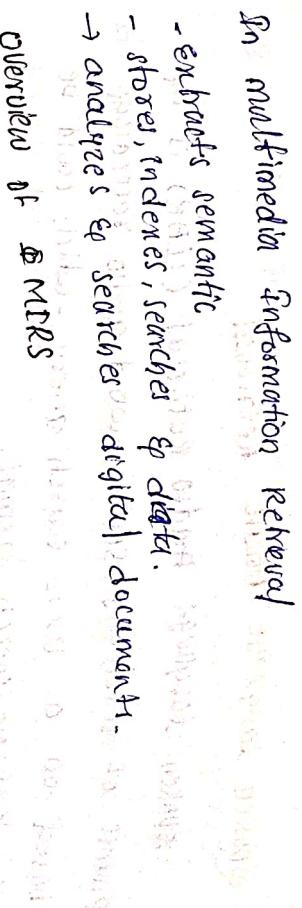
and retrieving relevant multimedia content like audio, image, graphs, video, animations from large databases, using content-based or metadata-based technique.

It aims at extracting semantic information from multimedia data sources. Data sources like audio, video, images depends on the architecture of multimedia database. It characterizes of the multimedia databases online, offline, user interface.

Query specification

Search engine

Database



Component of MMIR system

1. Data collection
2. Feature extraction
3. Indexing
4. Query Interface
5. Matching & Ranking
6. Retrieval output.

MRIS = Multimedia + Relevance search

- involves analyzing and retrieving non-textual content
- uses content-based features of metadata
- supports multi-modal queries

Feedback

Visualization

Annotations

Libraries

Auto/manual

Spoken language audio retrieval

spoken language audio retrieval (SLNR) is the process of retrieving spoken-word audio recordings based on a user's search query - which could be in text or speech format.

e.g:-

"Give me all audio clips where someone says

"climate change impact".

The system must identify all audio files containing that spoken phrase, even though the data is stored as audio - not text.

Why is this important :-

Spoken content is produced in large amounts in areas like news, lectures & interviews. Manually reviewing it takes too much time, so smart systems are needed to quickly find important parts.

How it works :-

Step 1 :- Speech Recognition (ASR)

Automatic speech recognition (ASR) is used to convert speech into text (called transcription). This acts as a bridge between the spoken data and traditional text retrieval.

audio → [ASR] → text → (Indexing & Retrieval)

Step 2 :- Indexing the text.

The text is organized using eco-friendly methods

- A user types or says a query, which is matched against the indexed text.

Step 3 : Get Audio :-

The system finds and returns audio clips and timestamps where the query is spoken.

Advantages :-

- Access spoken content without listening to hours of data.
- Supports search by speech and search-by-text.
- Helpful in archiving, journalism, legal, education.
- enables fast retrieval of relevant spoken content.

Challenges :- ASR errors, semantic gap, multiple speaker noisy audio, language dialects

* Spoken audio → ASR → text → search

↓
Retrieves audio.

Non-speech audio retrieval

Unlike spoken language retrieval which relies on converting speech to text, non-speech audio retrieval is about searching sounds that are not speech - like rain, alarm, animal calls, or environmental sounds, traffic etc.

How it work:

1. first you extract audio features from the sound
- pitch, rhythm, timbre, energy, spectral features.
2. these features are then stored in an index.
3. when the user submits a query - for example by giving an example sound or describing it - the system tries to match those features to sounds in its database.
4. It then returns relevant results : e.g. "find me all files with sirens".
- Techniques :- - query-by-humming / singing / lip-synching - audio fingerprinting
- Applications - music retrieval, audio surveillance, wildlife monitoring

How does graph retrieval work:-

1. First a graph database or index is built that stores nodes, edges and their attributes.
2. when a user gives a query (often a subgraph pattern) the system:
 - matches the pattern
 - retrieves graphs (or subgraphs) that match

- "find all subgraphs similar to the triangle shape"
- "find all subgraphs similar to the triangle shape"

common techniques :-
subgraph isomorphism, graph similarity matching,
graph embeddings, graph indexing

Applications:-

- drug discovery, social network analysis
- fraud detection
- web structure mining

Challenges:-

- audio signals are more variable and harder to classify than text
- background noise can interfere
- lack of standard "words" to represent many sounds.

Image Retrieval

Image retrieval is about retrieving relevant images from a large collection, using either text-label (CBIR) or visual features (CBCE), with growing interest in combining the two or for better results.

Graph Retrieval

It is about searching and retrieving information from data that is organized as graph nodes also from collections of nodes (vertices) and edges (links) between them.

Types of Image Retrieval

There are 2 major approaches:

1. Text-Based Image Retrieval (CBIR)

- Images are described by manually assigned keywords, captions, tags, or textual labels
- Search uses these textual labels
- Limitation: subjective, inconsistent, labor intensive.

2. Content-Based Image Retrieval (CBIR)

- Images are indexed using their actual visual content.
- (e.g., color, texture, shapes, spatial layout).
- No manual tagging needed.

- Example: query by example - you give an example image, and system retrieves visually similar images.

How does CBIR work?

- feature extraction
- the system analyzes images to extract low-level features.

• Color histograms, shapes, edges.

- texture patterns, colors, shapes, edges.
- indexing
- features are organized into an efficient data structure (like Inverted file structure) for fast searching.

3. Query Processing

- A user provides a query (keywords, or even an example image).

4. Similarity matching

- . The system compares the query features with stored image features using distance / similarity measures.

5. Ranking of retrieval

- Images are ranked by similarity to the query.

Challenges:-

- Semantic gap
- Variation → lighting, viewpoint, quality
- Subjectivity → "similar" is subjective.

Applications:-

- Digital photo albums
- medical image retrieval - museum and heritage
- e-commerce - surveillance.

Video Retrieval

Video retrieval is the process of searching for and retrieving relevant video segments from a large collection based on a user query.

Unlike text or image, video is more complex because it has:-

- visual frames, scenes,
- digital speech, music, sound effects
- temporal information (sequence, motion)

How it works:-

1. Video segmentation

- Divides the video into meaningful shots or scenes
- A shot is an unbroken sequence of frames from one camera.

- A scene is group of selected shots.

2. Feature extraction:-

- Extracts visual features (color, shape, motion patterns)
- Extracts audio features (speech recognition, music detection)
- Sometimes includes text (closed captions, superimposed text)

3. Indexing

- Stores features, metadata, and scene boundaries in an efficient data structure
- to allow fast searching.

4. Query formulation

- Users can search by keywords (e.g. "football goals")

- example video query - by - sample
- sketch usketch based retrieval

- And no clues (hamming a song) for visual

5. Similarity matching

- compares features of the query to features of the stored video
- Ranks videos by how well they match.

6. Retrieval and display

- Retrieves, ranked results.
- Sometimes displays key frames so user can quickly browse them

Challenges :-

- Large data volumes
- Semantic gap
- Temporal structure
- Subjectivity

Applications

- News archives, sports highlights, video-on-demand
- Movie scene searches
- Education / lecture retrieval systems.

Video retrieval is about intelligently searching large video collections using visual, audio and sometimes textual cues to help users find the right content quickly.