

* Introduction

- An IRS stores, retrieves and maintains information.
- This information can exist in multiple forms such as :- text format
 - audio
 - video
 - images
 - multimedia objects.
- An IRS refers to a system which is capable of representing, organizing, storing & accessing information.
- It is a proper way of representing & organizing information.
- And it provides an easy access to it users.
- The main objective of an IRS is to reduce user's overhead while searching for desired information.

→ Overhead is the difference b/w :

- time required to obtain information
- & time required to read the extracted information

$$\text{Overhead} = (\text{Time required to obtain info} - \text{Actual reading time required to read data})$$

- An IRS is a S/w.
- It consists of all the features & functions needed to manipulate information.
- Retrieval System is a system used to storing & extracting information.
- IRS is a S/w program, that helps to users in obtaining desired information.

* Objectives of IRS

- The main objective of an IRS is to reduce user's overhead while searching for desired information.
- Overhead means the difference b/w the time required to obtain the information & time required to read the data.
- Overhead = $\left\{ \begin{array}{l} \text{time required to obtain the desired info} \\ \text{time required to read the data} \end{array} \right\}$
- To provide right information to the users at the right time.
- To retrieve or provide relevant information to users in time.

→ Two measures of IRS are: 1. Precision
2. Recall

1. Precision

Precision = $\frac{\text{No. of relevant items retrieved}}{\text{Total no. of retrieved items}}$

2. Recall

Recall = $\frac{\text{No. of relevant items retrieved}}{\text{Possible no. of relevant items}}$
 \downarrow
all relevant items available in DB



Ideal Precision/Recall Graph.

Goal: Maximize recall to ~~get~~ most relevant items

Ex: If a search engine has high recall of 80%

- It means, it retrieves 80% of relevant items
- but misses 20% items.

High Precision: Retrieve only relevant documents/items.

High Recall: Retrieve all relevant documents/items

If a search engine has recall of 100%.

- If a search engine has recall of 100%.
- It means it retrieves all relevant document,
- missing none

* Functional Overview

→ IRS consists of following four functional processes, major

1. Item Normalization.

2. SDI (Selective Dissemination of Information)

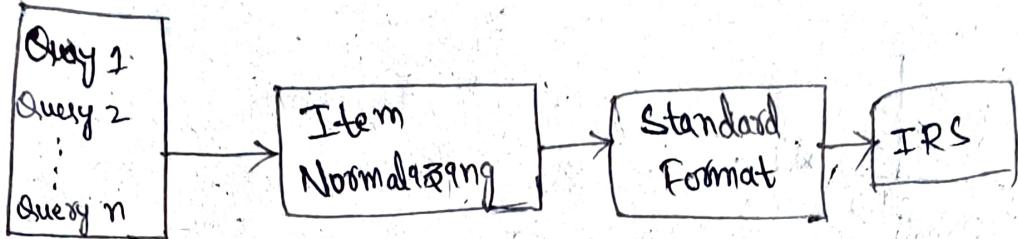
3. Document database Search

4. Index database Search

1. Item Normalization.

→ It is a process of decomposing the incoming item into a standard form.

→ It normalizes the incoming data in to a Standard format.



→ To normalize the data into standard form, it includes :

- Token identification
- Token categorization
- Apply stop lists
- Apply stemming Algorithms
- ~~Zeros~~, Zoning

• Apply Stop lists

→ It removes the common words like 'the', 'is', 'what'

→ It removes the common words like 'the', 'is', 'what'

Ex: what is the current cricket match score now

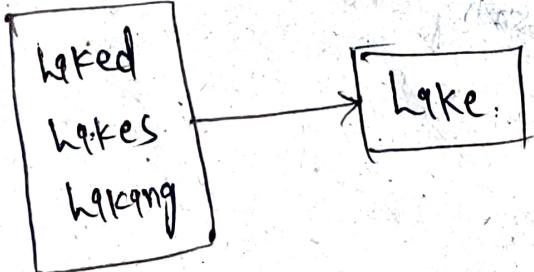
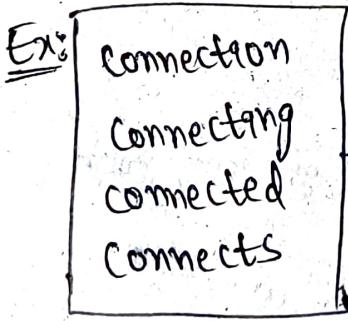
After applying Stop lists: Cricket score.

• Applying Stemming Algorithms

→ It reduces words to their base form.

→ It improves recall & reduce index size

→ Enhance query matching.



~~Z~~ Zoning

- It parses the data into logical subparts.
- This parsing is user-transparent.
- It is done, to improve the precision.
- the term ^{zone} considered as field.

Example: (Normalization)

~~Ex~~ John

Name	age	DOB	address
John	20	02/05/1992	123, Main Street, 2nd town, USA
Smith	30	05/02/1990	321, 2nd street, Town, USA

⇒ It is converted into:

John	20	02/05/1992	123, Main Street, town, USA
Name	age	DD/MM/YYYY	(street, city, state, pin code)

2. Selective Dissemination of Information

- It is commonly referred as SDI.
- SDI is a system that sends relevant information to users based on their interests.
- It is a process of categorizing data.
- SDI is also known as Mail Process.
- It is a process where users receive new updates based on their interests or profiles.

Ex: 1. A farmer subscribes to a weather update service.

- And gets notification about rainfall or temperature changes.

2. A business person subscribes to stock market alerts.

- So, he gets all updates about their industries.

3. A teacher subscribes to an education portal.

- So, they get all new updates about teaching of their subject.

→ SDI system push information to user based on their profile interests.

Pull System: A System where users pull information by searching.

Push System: A System where information is pushed to users based on their interests.

Difference:

- Pull System requires users to search actively.
- Push System that sends information to users based on their interests.

3. Document database Search Process

User enters a query:

→ A user submits a search query to the system

Search Process:

→ The system processes the query & search ^{on} the document database

↓ it contains only documents

Document database: the system searches all items in the document database which contains all received, processed & stored items.

Ex: User enters a query

- A user searches for "AI applications"
- the system processes the query and searches the document database for documents containing the known word "AI application"
- the system searches all documents on ~~the DB~~ document DB

Ex: If it contains 10,000 documents including research papers, articles & reports

Search Result

list of

→ The system returns a ^{list of} 20 documents that match the search query.

H. Index & Database Search

- When you find something interesting you want to save at for later in an information system
— this is called Indexing
- there are two classes of index files:
 1. Public Index file
 2. Private Index file

1. Public Index file :

- Maintained by library professionals
- Index every item in the Document DB
- Smaller no. of files
- Accessible by any one with proper privileges
↓
Permissions

Ex: A public index file called "All documents" that contains entries for every documents in the DB

2. Private Index file

- Created by individual user
- Each user can have multiple private index file
- Limited access list (only owners & specified users can access)

* Relationship to DBMS

DBMS

- 1) It is an organizational structure and handles structured data.
- 2) It provides precise semantics.
- 3) It provides Data Modeling Facility.
- 4) Every language is artificial.
- 5) Query Specification is complete.
- 6) In DBMS, there exists structured data format:
 - ↳ Structured data
 - ↳ Query matching is exact match
 - ↳ Inference is deduction.
 - ↳ Uses deterministic model.
 - ↳ We get exact results.
 - ↳ Knowledge Representation is transparent.

IRS

- 1) It is a system which stores, organizes, retrieves info.
- 2) It provides imprecise semantics.
- 3) It does not provide DMF.
- 4) Query language is almost similar to natural language.
- 5) Query specification is incomplete.
- 6) In IRS, there exists unstructured data format:
 - ↳ Unstructured data
 - ↳ Query matching is partial match.
 - ↳ Inference is Induction.
 - ↳ Uses Probabilistic model.
 - ↳ Sometimes relevant, often not.
 - ↳ Knowledge Representation is complex process.

* Digital Libraries

- Digital library acts as a repository.
- It stores data.
- It is a system, used to extract relevant information.
- At first, it is known as Electronic Libraries.
- ~~bcz~~, they utilize electronic devices to perform functions.
- Later US government renamed electronic libraries as Digital Libraries.
- As Internet technology is growing exponentially, the digital libraries conception is also emerging.
- It is responsible for extracting relevant information.

* Data Warehouses

- A Data warehouse is a special type of database.
- And it is used by the organizations for storing large amount of business data.
- It is a centralized repository.
- It stores data from various sources.
- It is designed to handle large volume of data.
- It makes easier to access, analyze data.
- It supports decision-making & analytics.

Example:

~~A. student warehouse~~

1. A school can store all the data about student.

- Such as grades, attendance, test results from different years in one place.
- We can see student's performance has improved or not.
- So, based on this information, we can make decisions.

2. A hospital stores Patient records, treatment history, appointment data.

3. Banks stores transaction data.

4. Library.

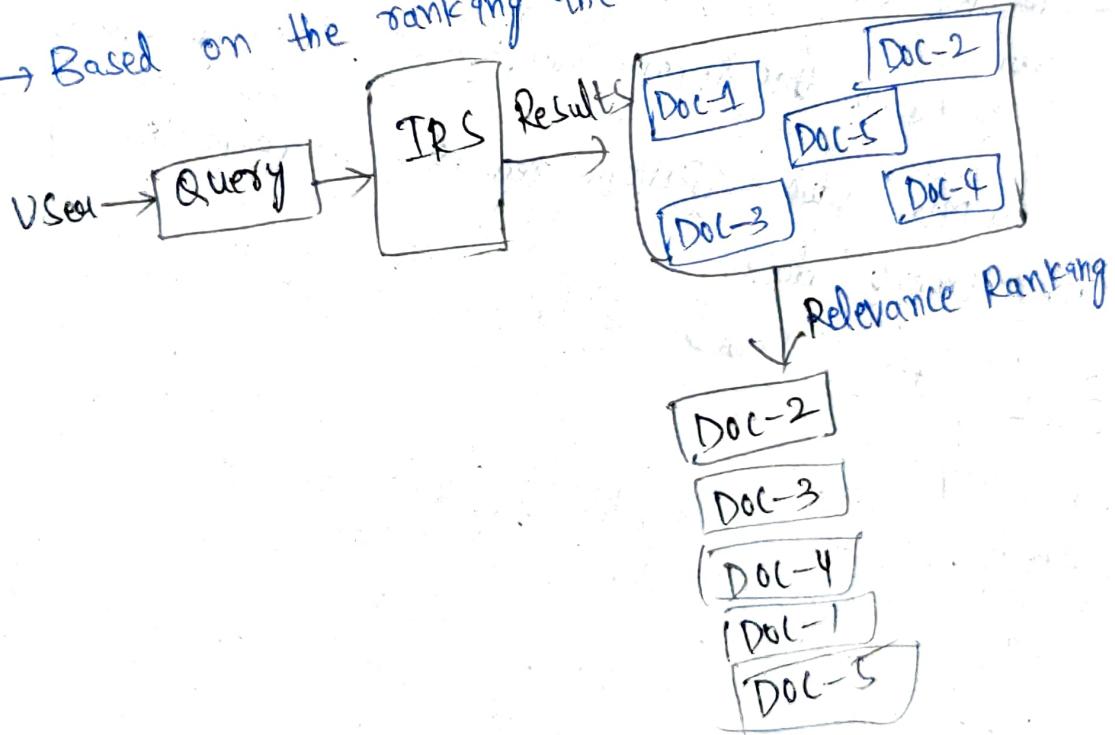
5. Airlines : Stores flight data, schedules, & customer experience.

* Browse Capabilities

- After the completion of the search, browse capabilities allows user to identify and display the required items.
- It allows users to explore data from categories or topics.
- Data Summary can be displayed ~~in~~ in two ways:
 1. Line the status of items
 2. Data Visualization.
- Browse capabilities includes:-
 - Relevance Ranking
 - Zoning
 - Highlighting.

1. Relevance Ranking :

- Based on the ranking the results will be displayed.

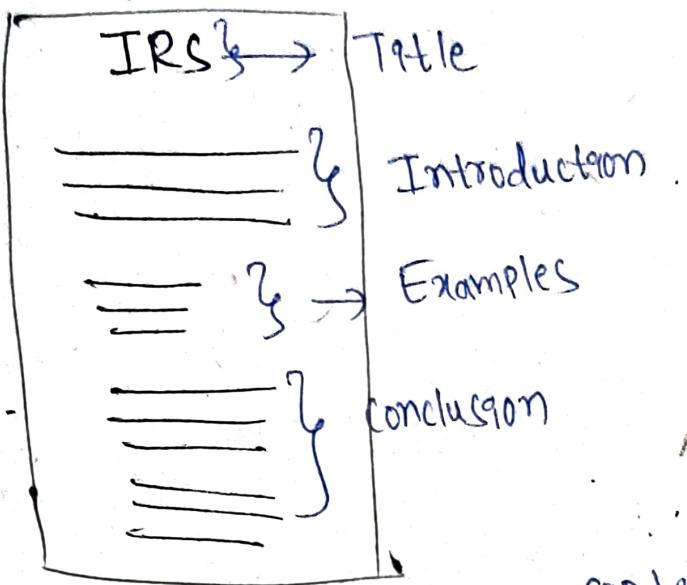


- Users can find relevant information faster.

2. Zoning

→ It plays a significant role in browsing.

Ex:



- Zoning enriches browsing capabilities.
- Ex: when user looking for main concept in book.
 - user can browse through conclusion.

3. Highlighting

→ It makes easier for users to identify relevant information quickly.

Ex: If a user search for IRS topic in google.

→ In results, main points are highlighted.

Give a brief note on ~~capabilities~~ capabilities.



Capabilities

search capabilities

- It creates mapping between user's need and items in the database.

User's need mapping items in Database (result)
Source Destination
- User's need (or) Search item consists of natural language and boolean logic.

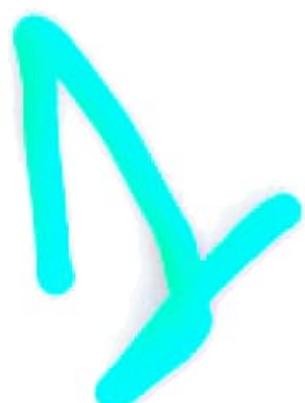
Ex: AI and IRS
- Every computer contains some algorithms.
- Different algorithms contain different functions.
- Algorithms are used for processing the search statement.
- The functions define the relationship between the terms in the search statement.

Functions

- Functions read search items and creates mapping.
- Boolean logic
- Proximity
- Contiguous word phrases
- Fuzzy Search
- Term Masking
- Numeric and date ranges
- Natural language queries

1. Boolean logic

- It is used in between query terms.
- It is a way to combine multiple search terms using logical operators (AND, OR, NOT)



Boolean Operators :

AND (Intersection)

→ It returns the document containing both terms.

Ex: AI and IRS

OR (Union)

→ It returns document either term

Ex: DataScience OR Machine learning

NOT (Difference)

→ excludes document containing term

Ex: AI NOT IRS

Q. Natural language Query

→ It allows to specify the importance of each search term.
using a value b/w 0.0 and 1.0

Ex: AI (0.8) and IRS (0.6) or DS(0.4)

→ It means AI is most important
- IRS is moderate in important
- DS is less important

3. Proximity

→ It specifies distance b/w two items in search terms.

→ It is used to increase the precision of the search

→ The typical form of Proximity is TERM1 within n units of

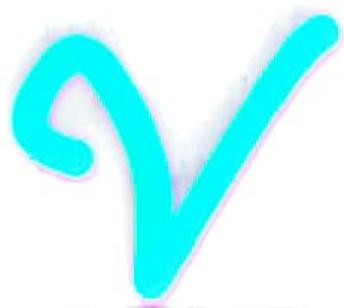
Ex: Data within 5 words of machine TERM 2

3. 4. Fuzzy Search

→ It search for similar words.

Ex: Fuzzy search on "Computer"

- computer (exact match)
- Comptor (minor spelling error)
- Computer (common misspelling)
- Compute (another common misspelling)



5. Contiguous word Phrases:

→ It is a sequence of two or more words that are treated as a single unit.

Ex: United States of America

→ It represents a single concept and can be used in search queries.

→ And combine a contiguous word phrase with Search operators like "AND", "OR" etc....

6. Term Masking:

→ In the context of term masking, a wild card character is a special symbol.

→ It is used to represent one or more characters in a search query.

→ The most common wild card characters are:

* - for multiple characters

? - for single character

\$ - for a special character

Ex: - *puter → Computer

- C?mputer → Computer

- multi\$national → multi-national