

* Information Retrieval => Activity of finding material (usually documents) of an unstructured nature (usually text) that satisfies an information need (query) from within large digital collections.

* Data => Raw, unprocessed facts and figures.

* Information => Data that has been processed into meaningful context

* Knowledge => Combination of information, experience and understanding that provides a context for decision-making.

* Wisdom => Application of knowledge and experience to make informed decisions, often in complex and uncertain ~~situations~~ situations.

* DB V/S IR =>

A DBS is a software application that is designed to store and manage large amounts of structured data. It provides an organized and systematic way of storing data and allows for efficient retrieval through methods such as SQL queries.

IR is concerned with storage, retrieval and ranking of unstructured data such as documents, images and videos. It involves the development of algorithms and methods that allow users to quickly search and retrieve information from large collection of data.

* Data Mining V/S Text Mining =>

Data Mining is the process of discovering hidden patterns, relationships and insights from large datasets. It involves the use of statistical techniques, machine learning algorithms and data visualization to extract meaningful information from large datasets. The goal of data mining is to find correlations in data that can help make informed decisions.

Text mining is a specialized area of data mining that is focused on extracting informations from large collection of unstructured text data. It involves the use of Natural Language Processing Techniques (NLP), ML algorithms and IR methods to extract meaningful information from unstructured text data.

The goal is to discover hidden relationships, sentiments, opinions and topics from large collections of text data.

* IR VS Text Mining =>

IR is concerned with finding, organizing and retrieving information from large collections of data. IR systems use algorithms and methods such as keyword-base search, Boolean search and probabilistic retrieval to match user queries with relevant information from the collections. The goal of IR is to provide users with the most relevant information in response to their queries as efficiently and effectively as possible.

* TEXT MINING EXPLAINED EARLIER *

* IR VS Information Extraction =>

IE is the process of automatically extracting structured information from unstructured or semi-structured data. IE involves use of NLP, ML and pattern recognition techniques to identify and extract relevant information such as entities, relationships and events. The goal is to transform the data into structured data that can be easily analyzed and utilized for decision making and other applications.

* IR EXPLAINED EARLIER *

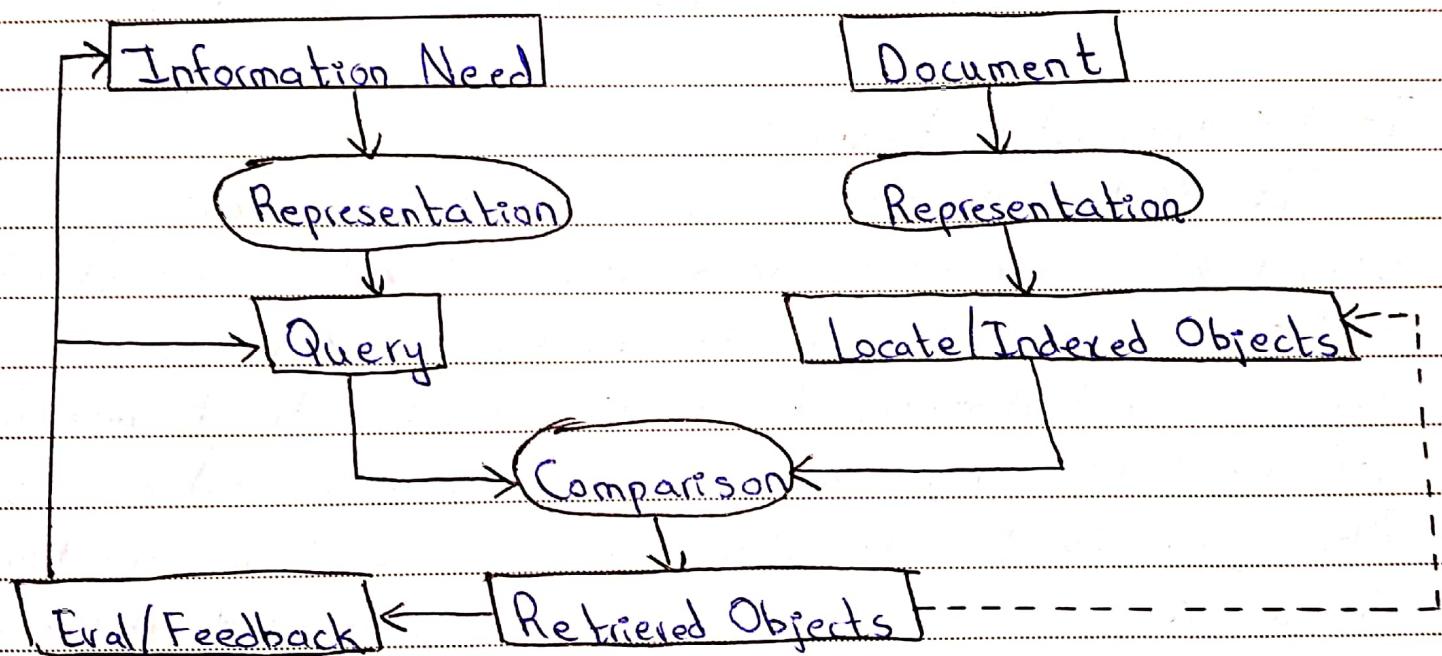
* Adhoc Retrieval Systems => IR systems designed to handle unstructured or semi-structured data. They are based on keyword-based search or probabilistic algorithms and do not have a predefined schema. They are commonly used in web search engines, intranet search, e-discovery, digital libraries. They are known for their ease of use and flexibility but can also be less precise than structured DBS.

* Routines VS Adhoc queries =>

Routines are predefined re-usable queries that are used to perform specific tasks stored in a DB.

Adhoc queries are one time, on-the-fly user generated queries that are executed as needed. They allow users to quickly and easily retrieve information from DB without the need for programming or routines.

* Basic IR process =>



* An IR system returns 8 relevant docs, 10 irrelevant.

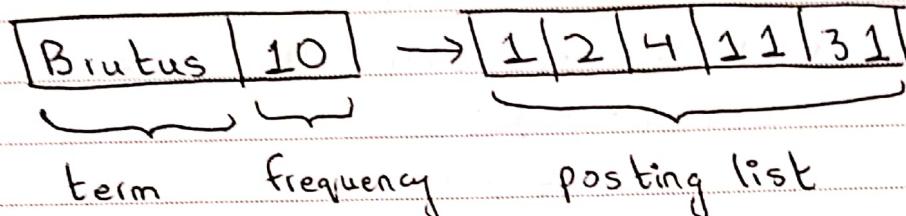
Total = 20 relevant docs in collection.

$$\text{Precision} = 8/18 = 0.44$$

$$\text{Recall} = 8/20 = 0.4$$

BACK TO PRESENT TIMELINE

- * Inverted Index => Actual index of the vocab that occurs in a collection. It consists of two parts:
 - 1) Lexicon / Dictionary / Vocabulary
 - 2) Posting / Posting List



↳ How are they made?

Consider two documents

Doc1 => My name is Wahaj Javed Alam.

Doc2 => I am Wahaj.

Step 1 => List the terms and their doc ID.

My	1	Wahaj	1	I	2
name	1	Javed	1	am	2
is	1	Alam	1	Wahaj	2

Step 2 => Sort the terms

Alam	1	is	1	name	1
Ilm	2	Javed	1	Wahaj	1
I	2	My	1	Wahaj	2

Step 3 => TADAA! Posting List time

Alam	1	→	1
am	1	→	1
I	1	→	1
is	1	→	1
Javed	1	→	1
My	1	→	1
name	1	→	1
Wahaj	2	→	1 2

Terms & Connectors \rightarrow words & features

Dated: _____

Transpose = Doc Term

* Term Document Matrix [possibly 2D Array] [DS]

	Documents		
	Play 1	Play 2	Play 3
Term 1	1	0	1
Term 2	0	0	0
Term 3	1	1	1

Since, we need Term1 and Term2 we AND them
and not needed Term3 so !AND Term3.

Brutus & Caesar & !Calpurnia.

\hookrightarrow can be really huge for memory \hookrightarrow we cannot machine read each doco

* Problems \Rightarrow Sparse matrix \hookrightarrow not noticing the order [only focus is features]
↳ might make Shakespeare has 44,000 words in each doco
very simple or would become a $44K \times 56$ matrix assuming there
very complex queries. \hookrightarrow exact matching {forms of words} {morph}

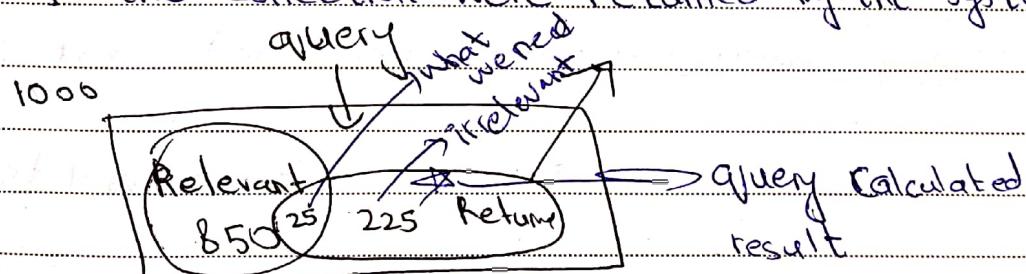
~~IR EVAL ASSUMPTIONS~~ Complex query formulation {what is needed VS isn't}
 \hookrightarrow user knows the features \rightarrow we can machine
 \hookrightarrow user can formulate logical queries. read each doc

* Effectiveness \Rightarrow Evaluation matrix

\rightarrow Effectiveness is the quality of its search results.

\rightarrow Precision \Rightarrow What fraction of returned results are
{Doc based} relevant to the need. {Term based}

\rightarrow Recall \Rightarrow What fraction of relevant document in
{static data} the collection were returned by the system.



Problems \Rightarrow can determine the two without knowing total docs.

\rightarrow Precision & Recall can be manipulated as
needed.

* Processing Boolean Queries =>

INTERSECT (p_1, p_2): [Posting List Algorithm]

- 1 answer $\leftarrow ()$
- 2 while $p_1 \& \& p_2$:
 - 3 if $\text{docID}(p_1) = \text{docID}(p_2)$
 - 4 add (answer, $\text{docID}(p_1)$)
 - 5 $p_1 \leftarrow \text{next}(p_1)$
 - 6 $p_2 \leftarrow \text{next}(p_2)$
 - 7 else if $\text{docID}(p_1) < \text{docID}(p_2)$
 - 8 $p_1 \leftarrow \text{next}(p_1)$
 - 9 else $p_2 \leftarrow \text{next}(p_2)$
- 10 return answer

↳ How intersection works?

- 1) Find terms in dictionary
- 2) Retrieve posting lists
- 3) Intersect

* Query Optimization => Selecting how to organize the work of answering a query so that the least total amount of work needs to be done by the system.

INTERSECT (t_1, \dots, t_n):

- 1 terms $\leftarrow \text{sortByFrequency}(t_1, \dots, t_n)$
- 2 result $\leftarrow \text{postings}(\text{first}(\text{terms}))$
- 3 terms $\leftarrow \text{rest}(\text{terms})$
- 4 while terms AND result:
result $\leftarrow \text{INTERSECT}(\text{result}, \text{postings}(\text{first}(\text{terms})))$
- 5 terms $\leftarrow \text{rest}(\text{terms})$
- 6 return result

$$\Rightarrow q_{v1} = 1000, q_{v2} = 500, q_{v3} = 4, q_{v4} = 4000$$

$$\text{optimized query} = q_{v3} \wedge q_{v2} \wedge q_{v1} \wedge q_{v4}$$

* Boolean Model (contd) =>

↳ Advantages =>

↳ Simplicity => easy to use and understand.

↳ Flexibility => allows user to search multiple terms and use operators to define criteria for search, thus conducting complex queries.

↳ Precision => only returns documents that exactly match the search criteria.

↳ Relevance => users can specify what they are looking for and exclude irrelevant results.

↳ Disadvantages =>

↳ Lack of Context => searches do not consider context of terms.

↳ Limited Ranking => no ranking of results, so it is difficult to determine relevant documents.

↳ Complex Queries => May become too complex and difficult to manage on extremely restrictive queries or extremely simple.

↳ Inflexibility => no partial matches or synonyms. This can lead to a large number of false negatives.

As a consequence, it might return too few or too many documents.

* Westlaw => Online legal research platform. It is a DB-driven IR system.

↳ AND => simple AND [contract AND dispute]

↳ OR/NOT ↳ citations => 's parallel citation

↳ Proximity Operators => contract w/5 dispute

[words should appear within 5 words]

↳ Wildcard Searches => allows to search terms with a specific pattern.

[contra*] → contract, contractory