

## Unit-5:

### Syllabus

- Web IR system
- Search Engine
- Web Crawling
- Search Engine Optimization
- Web Analytics
- Web mining taxonomy
- Web mining framework
- Social web mining
- Text mining
- Opinion mining
- Recommendation system.

X

X

X

X

### \* Web IR system

Web can be considered as a large-scale document collection, for which classical text retrieval techniques can be applied.

WebIR examines the combination of evidence from both the textual content of documents and the structure of the web as well as search behavior of users and issues related to the evaluation of retrieval effectiveness in the web setting.

IR is dominant form of info access.

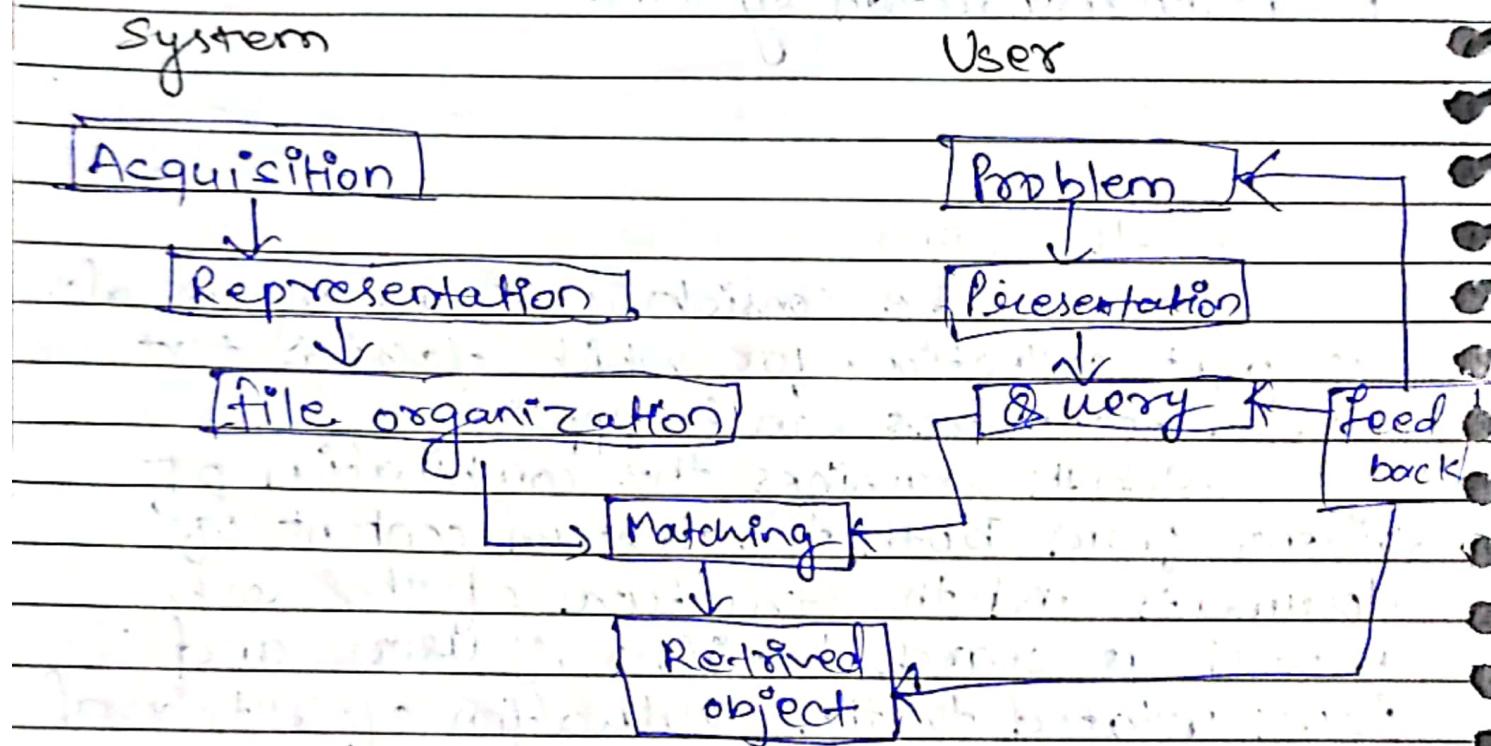
Software programs that deals with the organization, storage, retrieval and evaluation of information from document repositories, particularly textual information.

### IR Model

Selects & ranks the document that is required by the user or users asked for in the form of query.

### Retrieval status value (RSV)

Components of IR model



system

User

Acquisition

Representation

File organization

problem k

presentain

feedback

Query k

Matching k

Retrieved  
information

- Acquisition

Section of document or other object takes place.

- Representation:

contains indexing that contains free-text terms, controlled vocabulary, manual & automatic techniques as well.

- File Organization:

two types of file organization

- Sequential

- Inverted.

- Query :-

An I/O process starts when a user enters a query into the system

Good Write

## Data Retrieval

Information Retrieval  
Software that deals with storage organisation, retrieval and evaluation of information from document repositories particularly textual information.

Retain info about a subject

Small errors are likely to go.

Does not provide a solution to the user of the database system.

- Results obtained are approximate matches
- Result are ordered by relevance.

- probabilistic model

Good Write

## Data retrieval

Deals with obtaining data from DBMS for example, DBMS based on the query provided by user or application.

Determines the key words in the user query and retrieve the data.

A single error object means total failure.

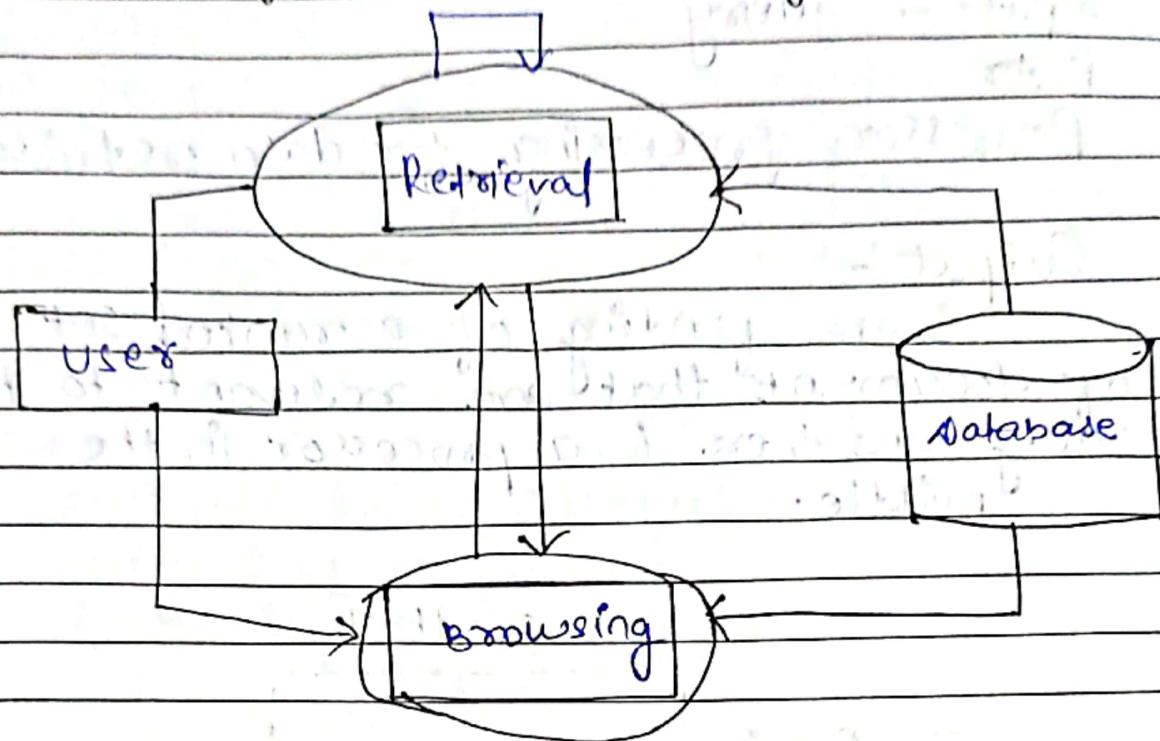
Provide

The results obtained are exact matches

Results are unordered by relevance.

A deterministic model.

## User Interaction with Information Retrieval System.

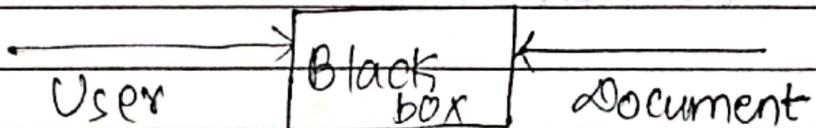


Past present & future of Information Retrieval.

- Early Development
- Information Retrieval in libraries.
- The web & digital libraries.

The logical view of the document is provided by representative keywords or index terms which are frequently used historically to represent documents in collection.

IR as a bridge.



## IR System Architecture.

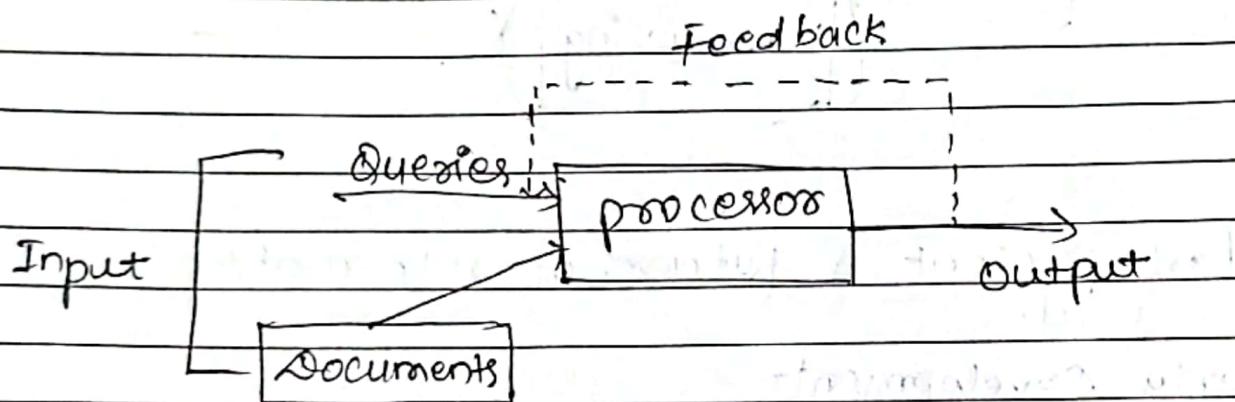
Input: query.

Reader

Processor processing for data retrieval.

Output:-

is the finding of a ranked set of documents that are relevant to the query and there is a processor in the middle.



## Search Engine:-

Tool designed to search for information, when this occurs on the world wide web it is then called a web search engine.

Search Results are usually presented in a list and are commonly called hits.

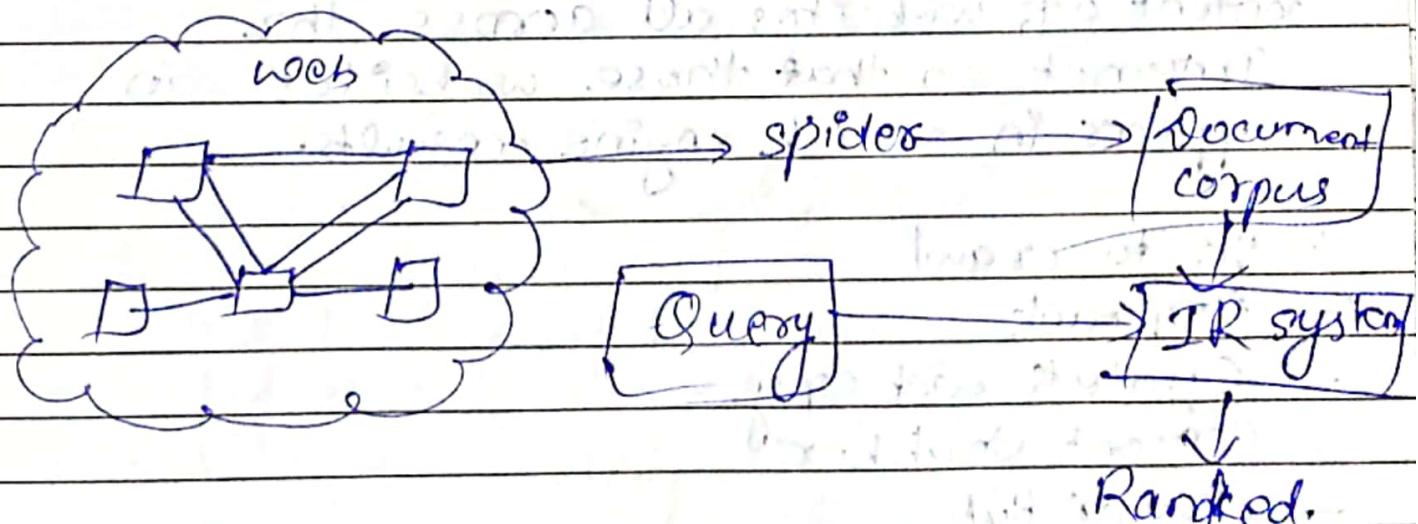
Common search engines include google, Yahoo, AOL Search, Ask.com, Bing & LookSmart.

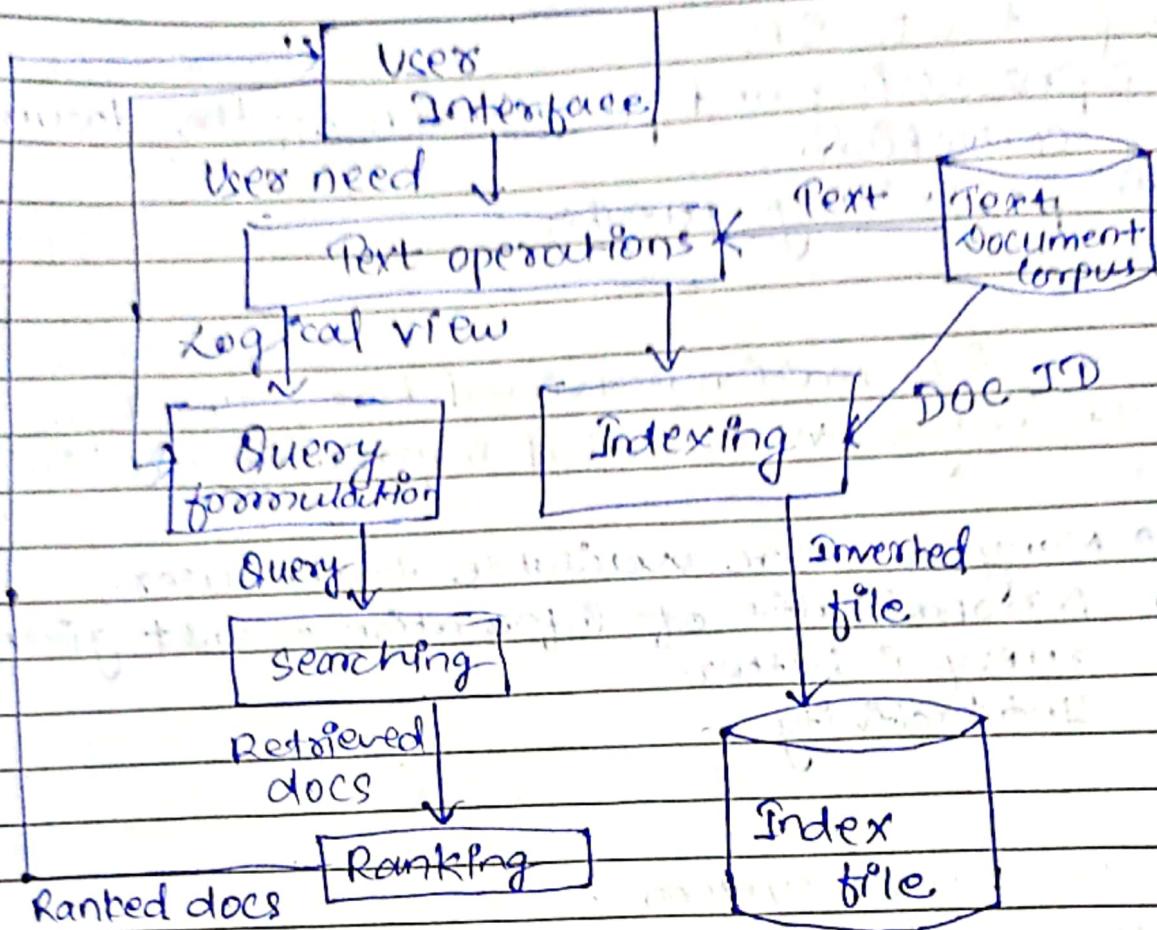
## Aspect of IR

- i) processing and presentation of the document collection.
- ii) processing queries.

Web information material has following advantages over classical information material.

- Many tools are available to the user.
  - Personalization of information result given a query is better.
  - Interactivity.
- 
- Collection or system
    - Hyper links are available to link one document to the other.
    - Statistic is easy to gather even in large sample sizes.
    - Interactivity.





### Web crawling :-

- A web crawler or spider is a type of bot that is typically operated by search engine like google and bing.
- Their purpose is to index the content of websites all across the Internet so that those websites can appear in search engine results.

### Ways to crawl

- HTTrack
- Cytetek web copy
- Content Grabber
- Parser Hub
- Outlink Hub

## Search engine optimization:-

→ process of improving the ranking or visibility of website in search engines.

→ higher the appearance of website in search engine list, the more visitors it is expected to receive.

→ considers how search engine works, what people search for, and which search term are typed.

Promoting a site to increase the number of links, is another SEO tactic.

## Importance of SEO :-

- To help gain more visitors
- Important for social promotion of website.
- It play important role in improving the business of a commercial site.
- Improving users experience.

## Basic principles in the working of a SE

- Crawling :- Search web pages linked to the website
- Indexing
- Processing
- Calculating Relevancy
- Retrieving Result.

## Customer Analytics

act of studying customer behavior or interests and interpreting them to take important business decision.

### Web analytics

→ defined as 'the measurement, collection as well as analysis and reporting of the data for optimization of the web page usage'.

→ provides understanding of user behaviors over the web page or website.

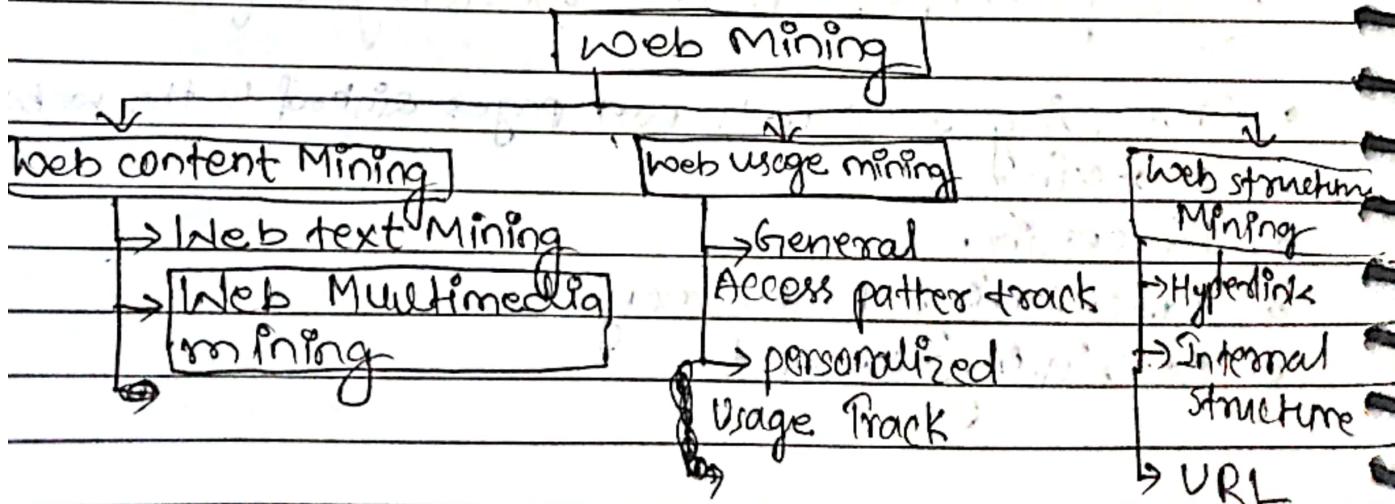
→ tool for marketing research as well.

→ keeps track of:-

- number of users visiting the website.
- duration of their stay.
- the pages they visit.

Hence used to improve effectiveness and performance of the website.

### \* Web mining taxonomy:-



Web mining is the process of data mining techniques to automatically discover and extract information from web documents & services.

### Application of Web Mining

- o Used for searching
- o predict user behavior.
- o Improve power of web search engine by classifying the web documents and identifying the web pages with additional info.
- o Improve other performance metrics.
- o useful of a particular website and e-service eg.: landing page optimization.

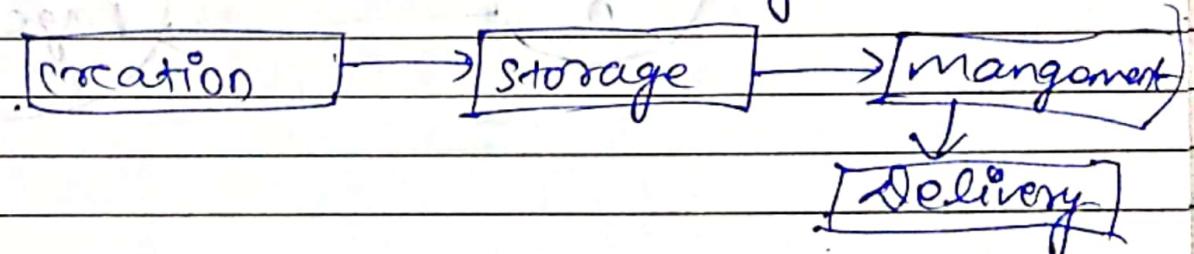
### Web Content Mining:

→ application of extracting useful info from the content of the web documents.

→ contains several types of data - text, image, audio, video etc.

→ provides interesting patterns about user's needs

→ Also known as text mining



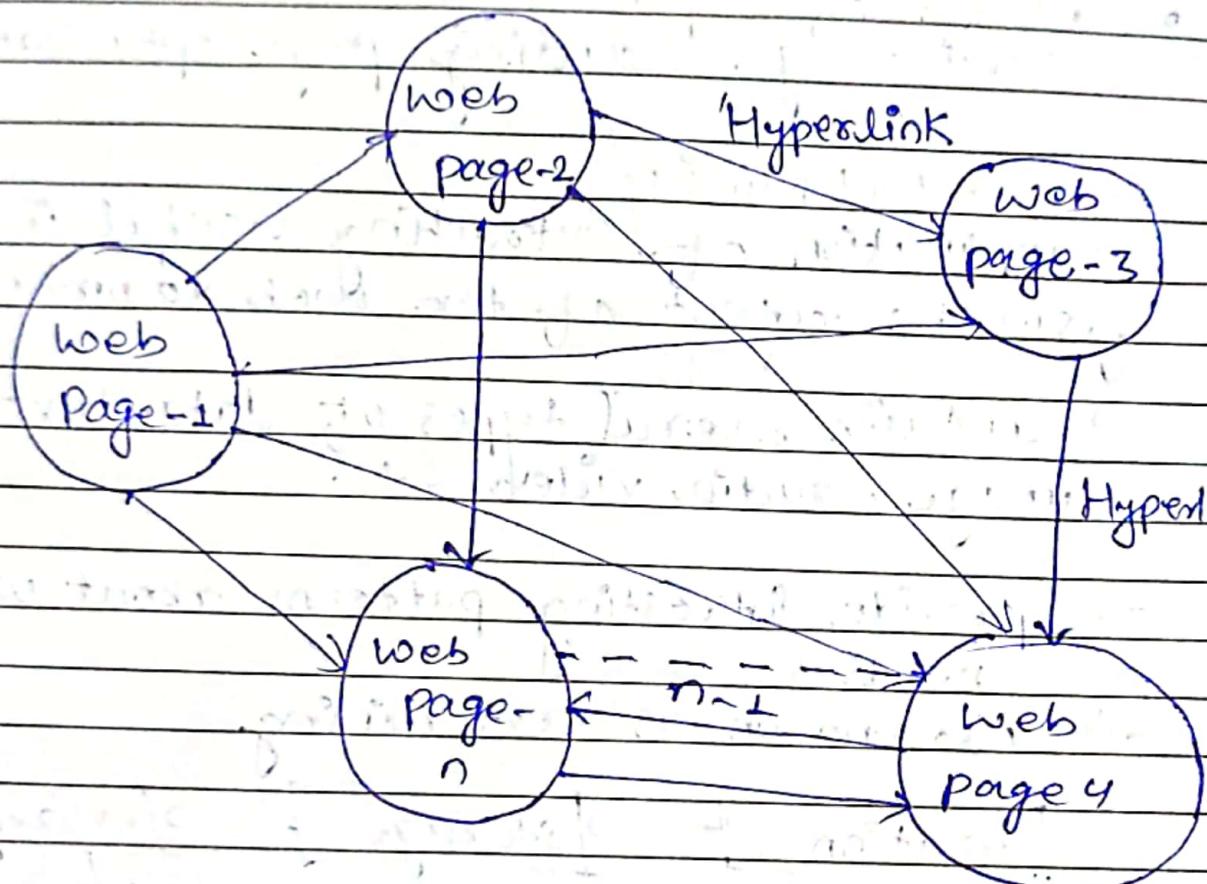
o Web Structure Mining:

→ Application of discovering structure information from the web.

→ consists of web pages as nodes, and hyperlinks as edges connecting related pages.

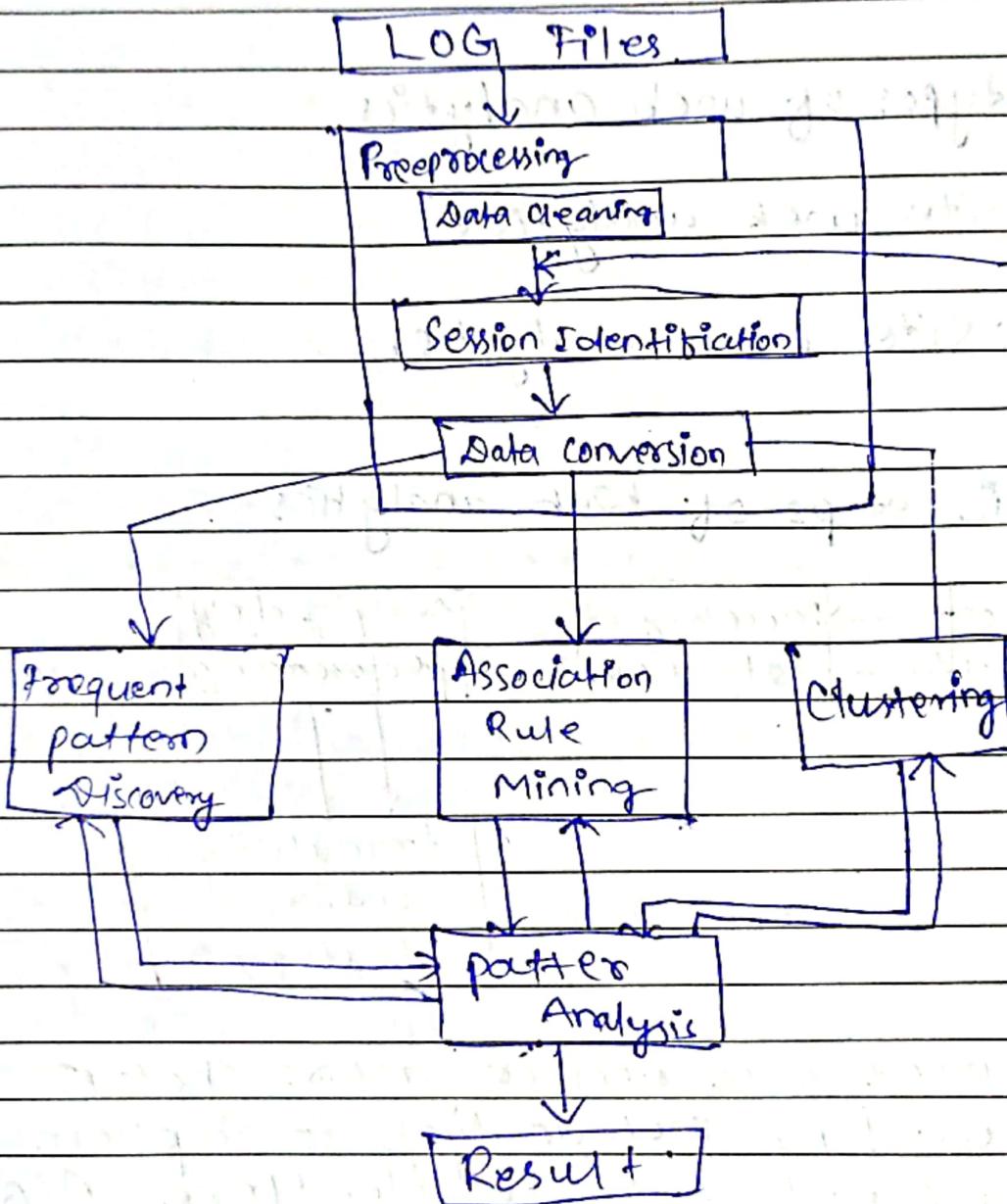
→ Shows the structured summary of a particular website.

→ To determine the connection between two commercial websites.



## Web usage mining :-

Application of identifying or discovering interesting patterns from large data sets



→ process of gathering, processing and evaluating data from website.

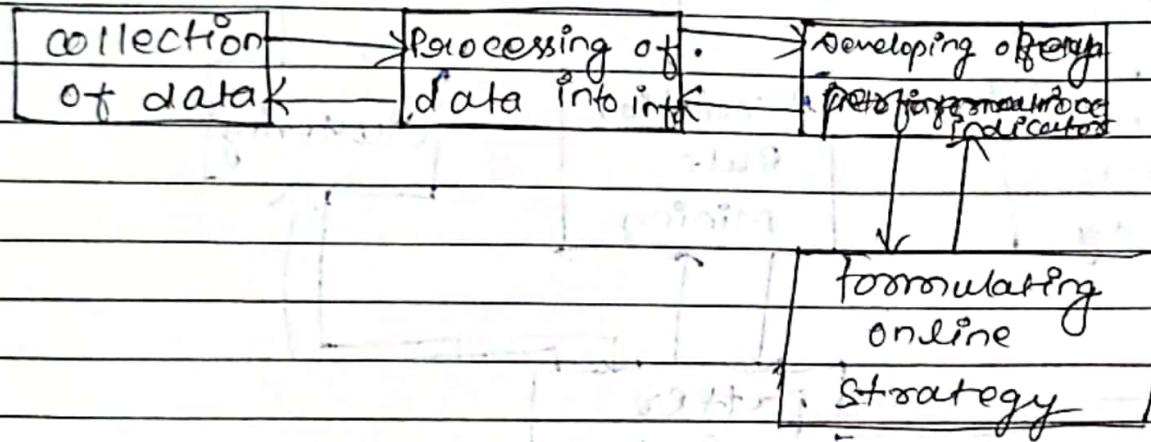
Web analytics has many applications like it shows web navigation, time-specific website traffic, search engine data, visitor data.

Two types of web analytics.

On-site web analytics.

Off-site web analytics.

Basic steps of Web analytics



Experiments are another factor that ~~need~~ is introduced by web analysts to improve the performance of website. Here A/B testing is done.

## \* Tools used in web analytics

- Google analytics
- Yahoo web analytics
- Google website optimizer
- facebook insights
- 4Q by iPerceptions
- Click table
- Optimizely.

## Advantage of web analytics:-

- Measure online traffic
- Tracking bounce rate
- Optimizing and tracking of Marketing campaign.
- finding right target audience and its capitalization.
- Improves and optimize websites and web services.
- Conversion rate optimization
- Tracking business goal online.

Good Write

## Sentiment analysis

→ an approach do NLP that identifies the emotional tone behind a body of text.

→ It involves the use of data mining, machine learning, and artificial intelligence(AI) to mine text for sentiments and subjective information.

i) fine grained :- further categorization to provide a more precise level of polarity.

ii) Emotion detection:- specifies emotion rather than positivity.

iii) Intent based action behind a text in addition to opinion.

iv) Aspect based analysis:- specific component being positively or negatively mentioned rather than whole product.

## Challenges of sentiment analysis

- i) Inaccuracies in training model.
- ii) sentiments written in neutral manner.
- iii) when system cannot understand the tone or context.

- Use of emoji's and irrelevant information.
- contradictory sentiments.

Recommendation system:-

Mainly deals with the likes and dislikes of the users, by learning from the large set of repositories generated day by day.

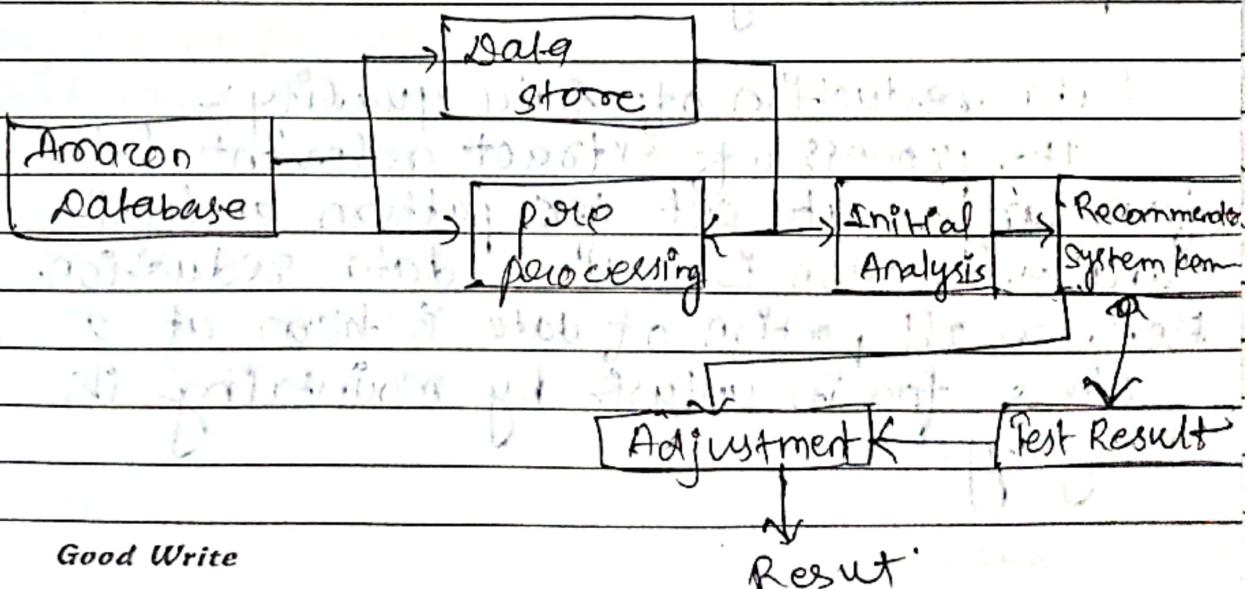
Types:-

User based recommendation:-

Here, we calculate person's similarity measure, which is needed to determine the closely related users, i.e. whose likes and dislikes follows the same pattern.

Item based Recommendation:-

Initial aim is to obtain the mean adjusted matrix. Here we use cosine similarity matrix.



## Working steps of Data mining:-

1. Data cleaning (removing noisy data & filling gaps)
2. Data Integration (combining different data sets)
3. Data reduction for data quality (extract relevant info)
4. Data transformation
5. Data mining
6. Pattern Evaluation
7. Representing knowledge in Data mining

### Data cleaning :-

After collecting data from different sources, it need to be cleaned using methods like binning. This method cleans up the noisy data and make it eligible for data integration. Dirty data leads to poor insights and increase the cost & time.

### Data Integration :-

Combining different data sets for proper data analysis is called data integration. Some inconsistency of the data also cleaned up in this stage.

### Data reduction of Data quality :-

The process of extract relevant data from large data set for pattern analysis and evaluation is called data reduction. Here small portion of data is taken at a time, for its analysis by maintaining its integrity.

## 2. Data Mining

### Data Transformation

In this standard process, engineers transform data into an acceptable form to align with mining goals.

Techniques to do so are smoothing or eliminating noise from data, aggregation, normalization or discretization.

### Data mining

To extract useful trends and optimize knowledge discovery to generate business intelligence. Use methods like clustering, classification or other modelling techniques to ensure accuracy.

### Pattern evaluation

From these data specialist can generate patterns of business knowledge. Visualization data mining is use for better understanding.

### Representing knowledge in datamining:-

Use combination of data visualization, reports and others mining tools to share information with others. These insights presented in reports help the higher authority to take decisions, generate new business, eliminate waste, and create better advertising campaigns.

## Difference between Data Mining & Web Mining..

### Data mining

- Statistical technique of processing raw data in a 'structured' form.
- Pre-processing database and spreadsheets were used to gather information.
- Processing of data is done directly.
- Statistical techniques are used to evaluate data.
- Data stored in structured format.
- Data is homogenous and easy to retrieve.
- Mining of mixed data.
- AI, ML and statistics.
- Marketing, medicine, health care.

Good Write

### Web mining

part of data mining which involves processing of text from documents

The text is used to gather high quality information.

Processing of data is done linguistically.

Computational linguistic principles are used to evaluate text.

Data is stored in unstructured format.

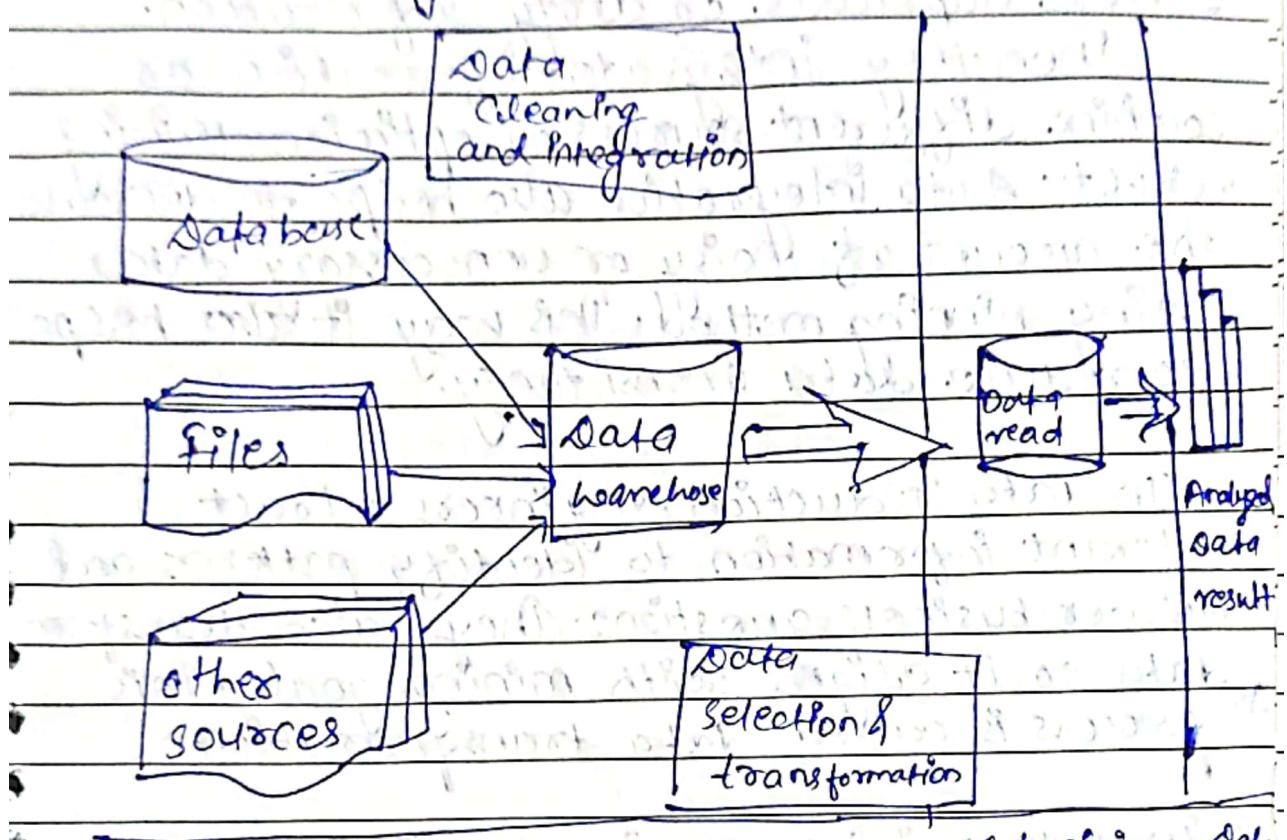
Data is heterogeneous and is not easy to retrieve.

Mining of text data is only done.

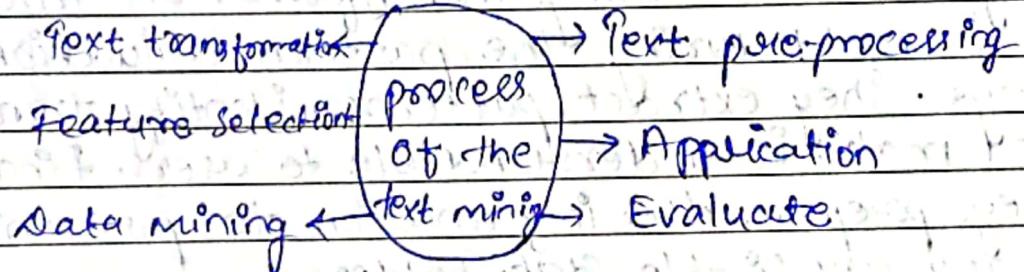
Pattern recognition & NLP

Business customer profile analysis

## Data mining



### Data pre-processing



## Data mining process

First, specialist need to cleanse data to remove duplicate or dirty information. Then they integrate information or combine different sources to optimize mining result. Data integration also helps to decrease the amount of noisy or unnecessary data, using pruning method. This way it also helps to reduce data inconsistency.

In data reduction, engineers extract relevant information to identify patterns and answer business questions. They also transform data so it aligns with mining goals. This process is called data transformation.

In data mining, engineers assign relevant patterns to each data set before they extract it. They then generate models ~~with~~ with clustering or classification techniques.

Engineers then bring the information into the real world during the pattern evaluation stage. They extract patterns, identify trends and make it understandable to users. Finally they prepare the information to present to any applicable stakeholders.

Business owner use data mining insights to optimize decision making, increase sales and learn more about customers.

## Unit 4

web database

Database connectivity

JDBC

ODBC

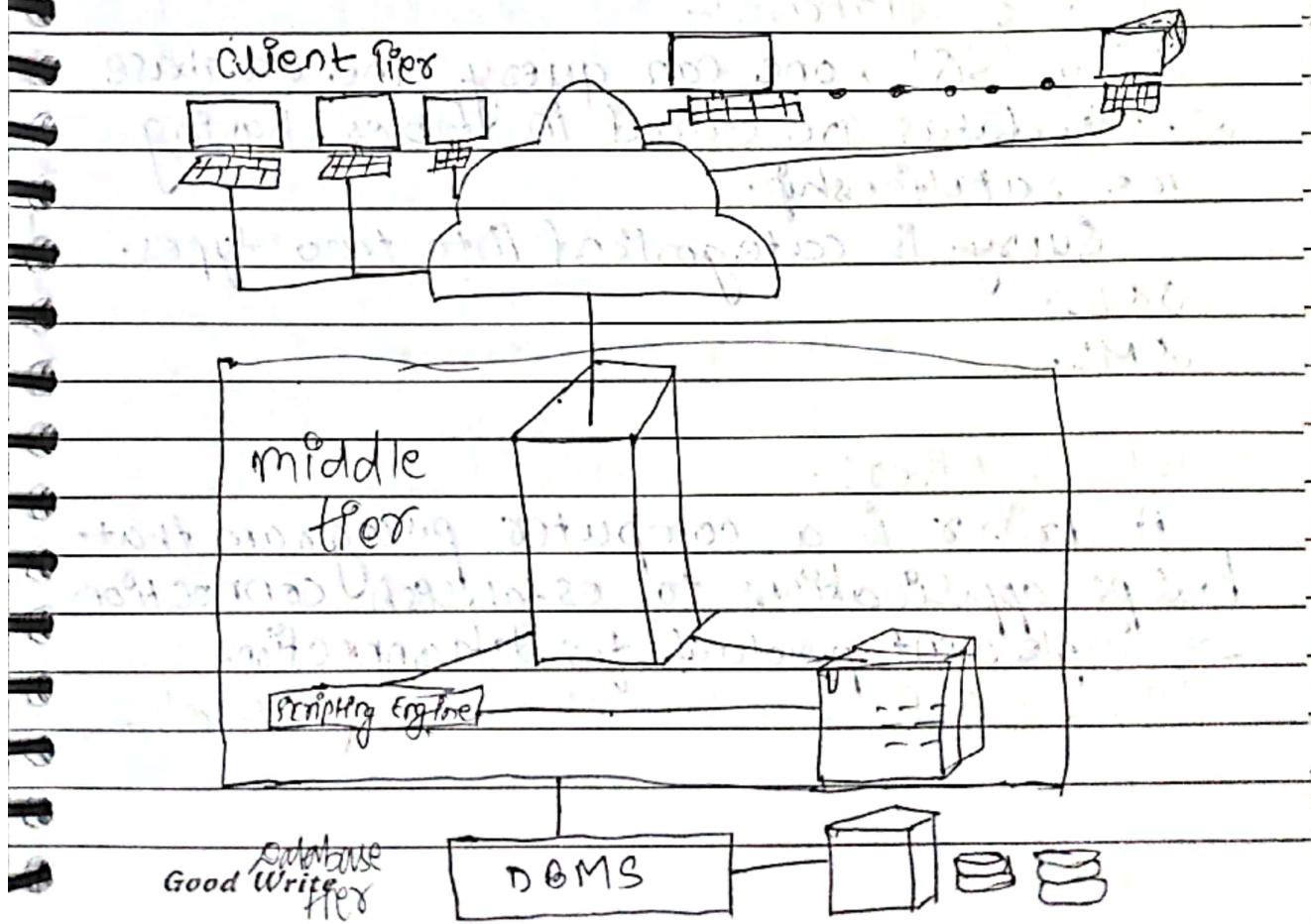
Database to web connectivity.

Web database connectivity

A web database is a database application designed to be managed and accessed through the internet.

Website operators can manage this collection of data and present analytical result based on the database application.

Three tier web database architecture.



## Database connectivity

Database connectivity allows the client software to communicate with the database server software.

Elements of frontend application  
Websites like buttons, fonts or menus need to be connected to the database to deliver relevant information to the end-user.

Database connectivity can be done using different programming languages, Java, C++ and HTML. By using these language, one can connect a program to a particular database and can query data to access & manipulate them.

### Querying database:-

Using SQL, one can query the database where data are stored in tables having some relationship.

Query is categorized into two types.

DQL.

DML.

### Database drivers:-

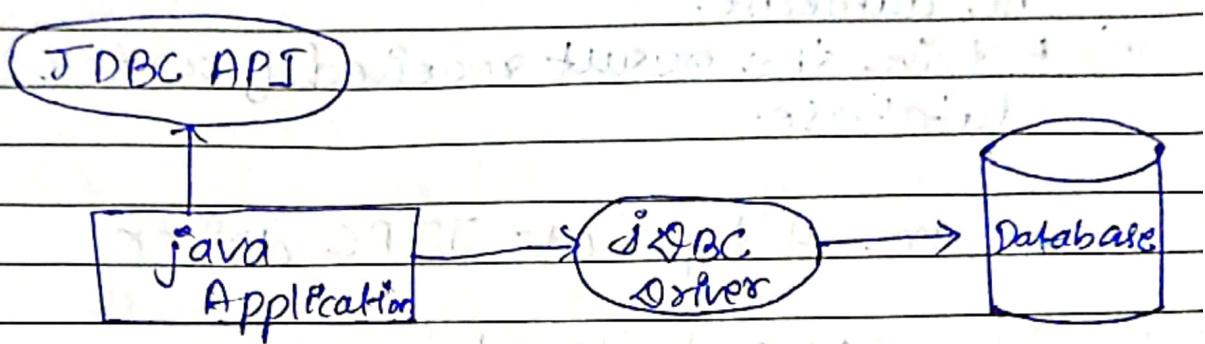
A driver is a computer program that helps applications to establish connection. It implements protocol for data connection.

JDBC & ODBC are two standard protocols.

The driver works like an adaptor which connects a generic interface to a specific database vendor implementation.

## Java Database Connectivity (JDBC)

It's an API for java programming language to connect and execute query with the database.



It access data from the relational database table. JDBC use drivers to connect with database. Some of the popular JDBC API interfaces include ResultSet interface, CallableStatement interface, Driver Interface, Connection Interface, Statement Interface and Rowset Interface.

## popular classes of JDBC API

Driver Manager API

Blob class

Clob class

Types class

• ODBC API uses ODBC driver which is written in C language. i.e. platform dependent and unsecured.

JDBC API can do:-

- i) Connect to the database.
- ii) Execute queries and update statements to the database.
- iii) Retrive the result received from the database.

There are 4 types of JDBC drivers

i) JDBC-ODBC bridge driver.

• Uses ODBC drivers to connect with database.  
• The JDBC-ODBC bridge converts the JDBC method calls into the ODBC function calls.

JDBC API

Data Application

JDBC-ODBC  
bridge driver

ODBC  
driver

vendor  
database  
library



## Advantage:-

- easy to use

- can be easily connected to any database.

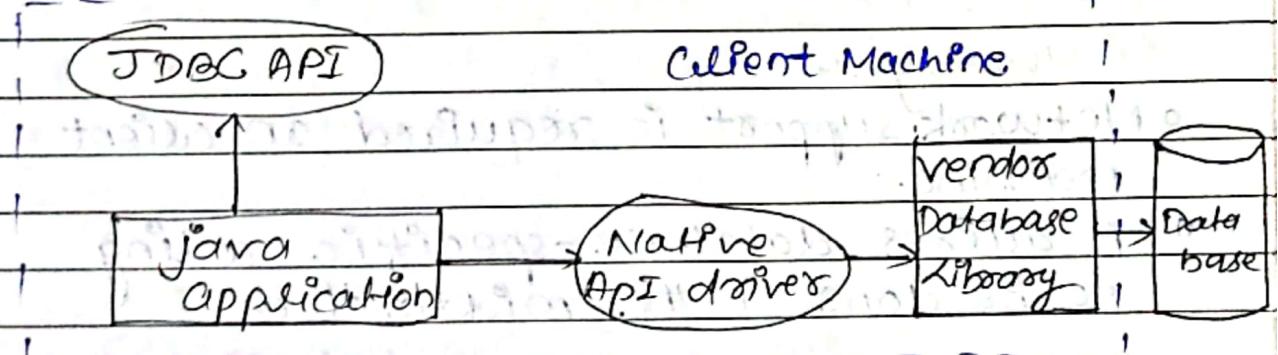
## Disadvantage:-

- Performance degraded because JDBC method calls is converted to ODBC function calls.

- Since its platform based, so it should be installed on the client machine.

## 2. Native API driver:-

The native API driver uses the client-side libraries of the database. It's not entirely written in Java, converter JDBC method calls into native calls of the database API.



## Advantage:-

- performance upgraded than JDBC-ODBC bridge driver.

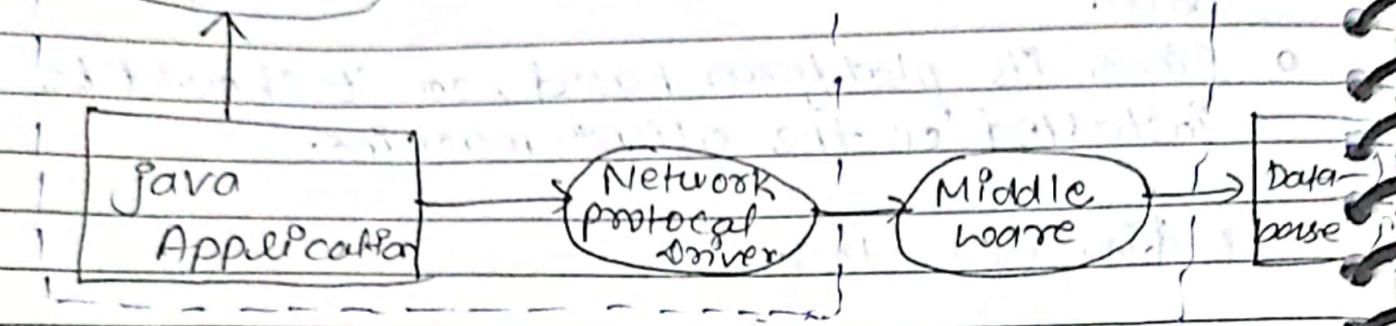
## Disadvantage:-

Native driver and vendor client library needs to be installed on client machine.

### 8) Network Protocol driver:

The Network protocol driver uses middle ware that converts JDBC calls directly or indirectly into the vendor-specific database protocol, fully written in Java.

JDBC API



Advantage:-

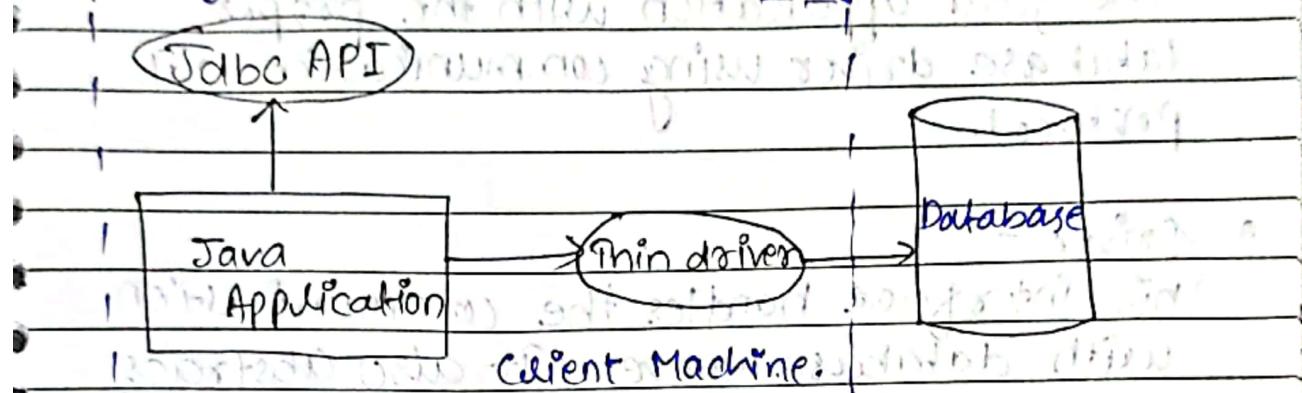
- No client library is required because of application server.

Disadvantage:-

- Network support is required on client machine.
- Requires database-specific coding to be done in the middle tier.
- Cost is high because it requires database specific coding to be done in the middle tier.

## 4) Thin driver

The thin driver converts JDBC calls directly into the vendor-specific database protocol. That's why it's known as thin drivers. It is fully written in Java language.



### Advantages

- Better performance than all drivers.
- No software is required at client side.
- or server side.

### Disadvantage

- Drivers depend on the Database.

### JDBC connectivity with 5 steps:

- Register the Driver class (Driver Manager)
- Create connection (Connection interface)
- Create statement (Connection interface)
- Execute queries (Statement interface)
- Close connections (Connection interface)

## JDBC components.

- **DriverManager:**

This class manages a list of database drivers. Matches connection request from the Java application with the proper database driver using communication sub protocol.

- **Driver:**

This interface handles the communication with database server. It also abstracts the details associated with working with Driver objects.

- **Connection:**

This interface with all methods for contacting a database. All the communication with database is through connection obj. only.

- **Statement:**

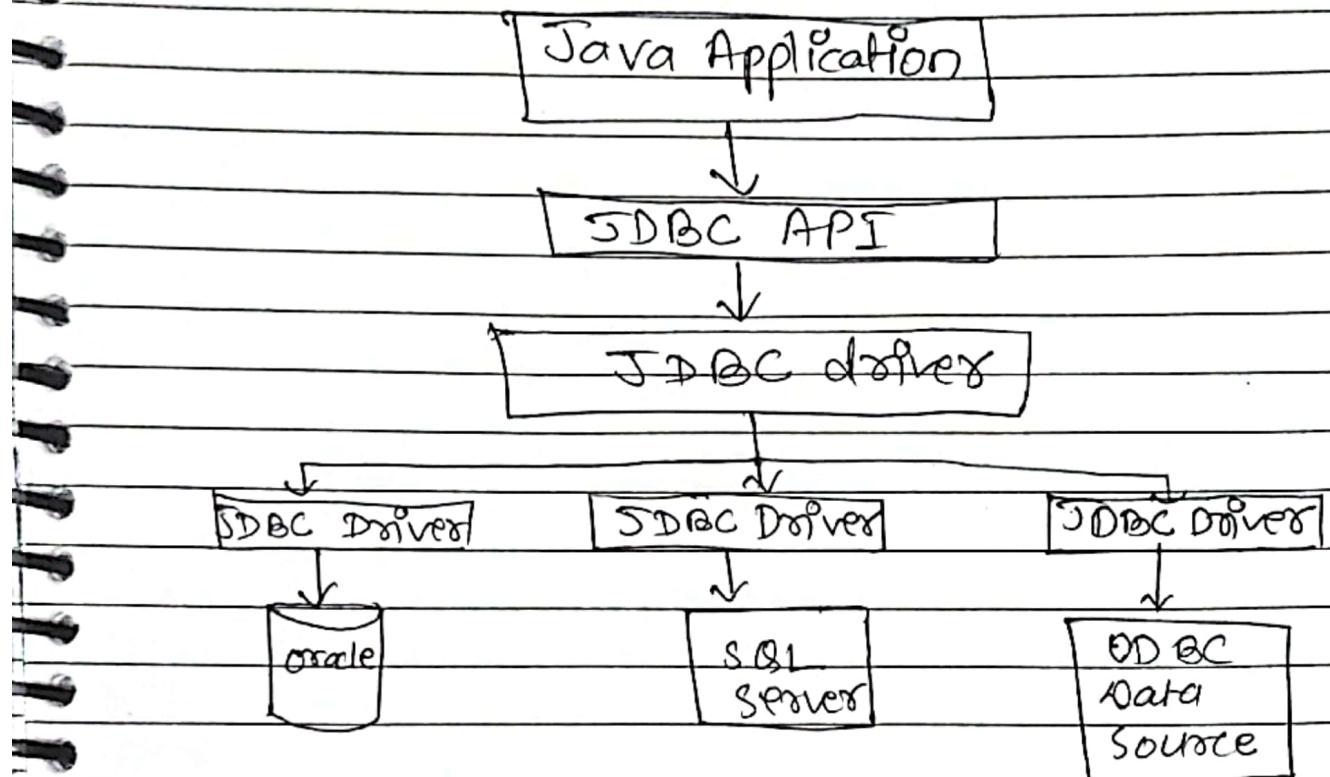
Here we use objects created from this interface to submit the SQL statements to the database.

- **ResultSet:**

These objects hold data retrieved from a database after we execute an SQL query using Statement object.

- SQL exception.

This class handles errors that occurs in a database application.



### Database

### ODBC

Open Database Connectivity (ODBC) is a standard (API) for accessing DBMS. The designer of ODBC make it