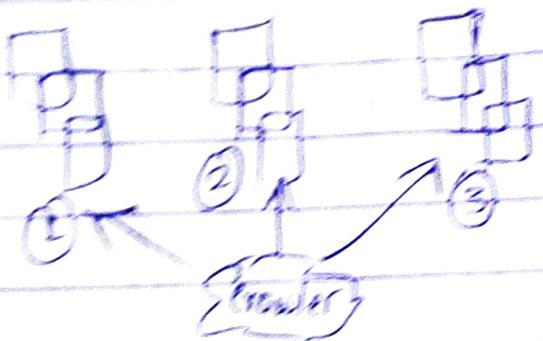


① Web Crawling:

It is a process by which we gather pages from the web, in order to index them and support a search engine.

Objective:

→ to quickly and efficiently gather as many useful web pages as possible, together with link structure that interconnects them



① Features:

① Robustness: crawlers must be designed to be resilient to traps

② Politeness: web servers

→ have both implicit and explicit policies regulating the rate at which a crawler can visit them

③ Distributed:

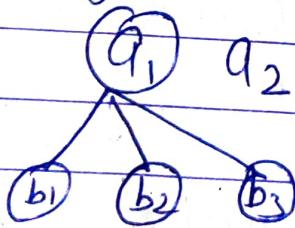
④ Extensible

⑤ Scalable:

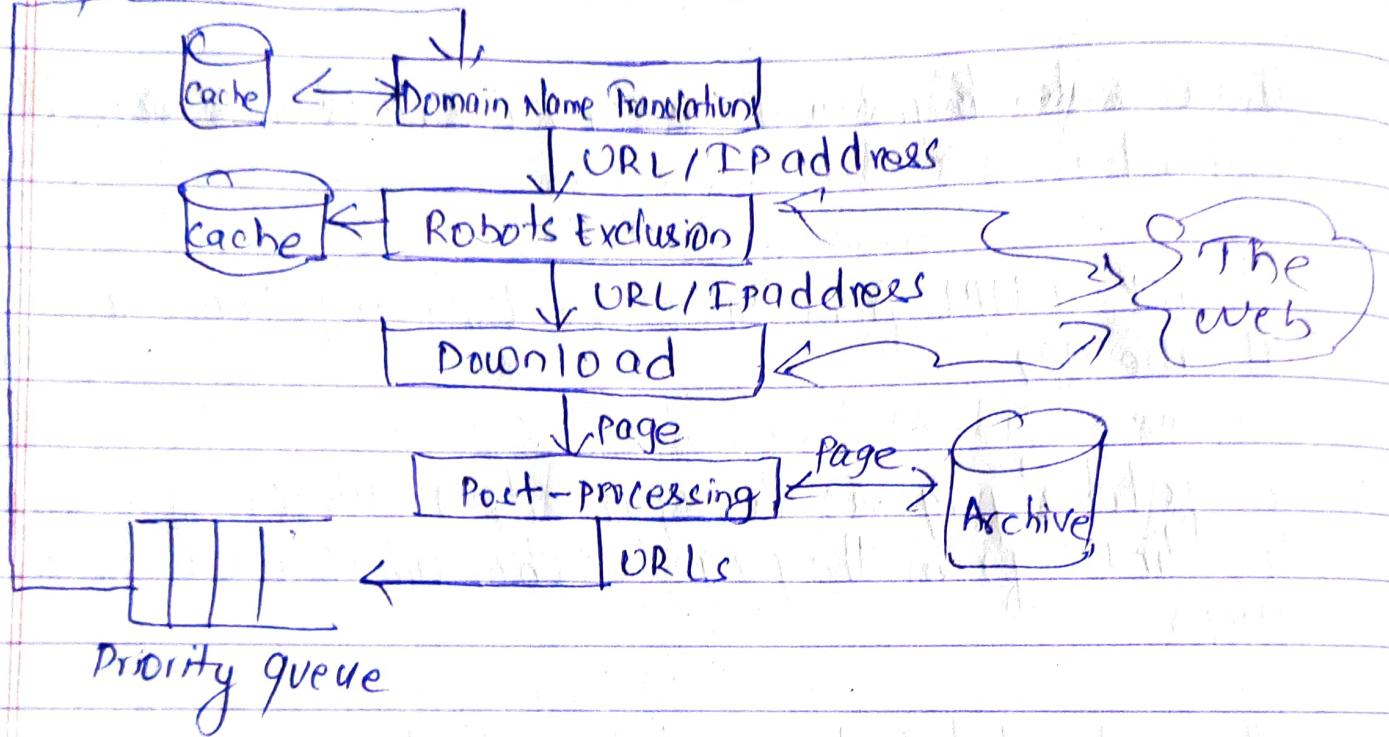
⑥ Performance and Efficiency: ⑦ Quality ⑧ Freshness

① Basic Operation:

- ① The crawler begins with one or more URLs that constitute a seed set
- ② It picks a URL from this seed set, then fetches the web page at the URL
- ③ The fetched page is then parsed, to extract both text and the links from the page (each of which points to another URL)
- ④ The extracted text is fed to a text indexer.
- ⑤ The extracted links (URLs) are then added to a URL frontier, which at all times consists of URLs whose corresponding pages have yet to be fetched by the crawler.



① Components:



② Domain Name Translation:

- Translating host name into 32-bit IP address
- for speed, maintains its own cache which requires thousands of translations per second.
- Cache maintains mapping between host names and IP addresses

③ Robots Exclusion:

- With IP address, available crawler must check that access to the page is permitted by the website.
- robots exclusion protocol
- access permissions are obtained by appending the path "robots.txt" to the host name

① Download :

- Accesses the website via the HTTP protocol and downloads the page
- Crawlers may be required to resolve directions during download
- Use of JS cause trouble.

② Post-processing :

- After download, the crawler stores the page in an archive for indexing by the search engine
- During indexing, post-processing, the crawler must respect conventions requesting that it not index a page or follow certain links:

```
<meta name="robots" content="noindex">  
<a rel="nofollow">
```

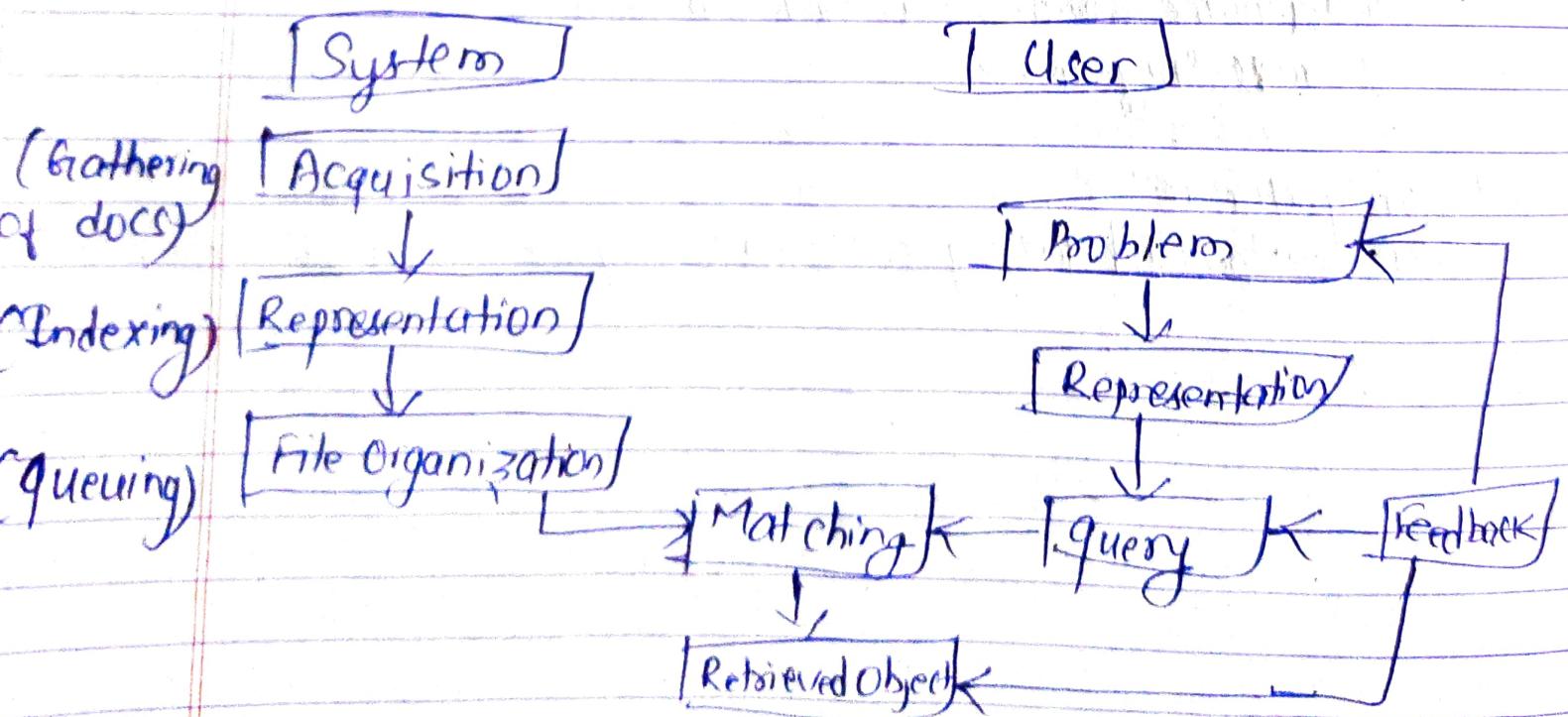
③ Priority queue:

- URLs extracted during post processing are inserted into pq

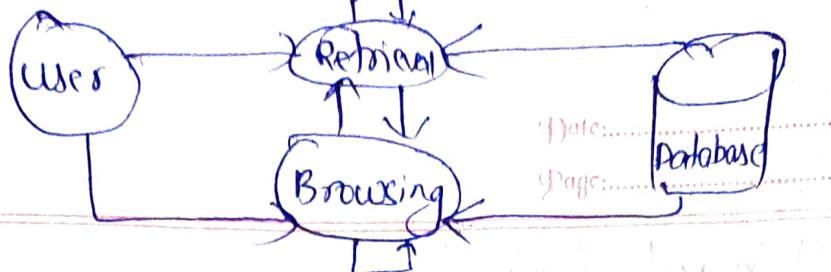
⑤ Web Information Retrieval:

- A software program that deals with the organization, storage, retrieval and evaluation of information from document repositories, particularly textual information.
- IR Mode:
 - Selects and ranks the document that is required by the user or the user has asked for in the form of query.

⑥ Components of IR model :



User Interaction:



① Difference :

IR retrieval

Data Retrieval

- ① from document repositories ① from a database management system

② retrieves information about a subject

② determines the keywords in user query and retrieves the data

③ small errors are likely to go unnoticed.

③ single error object means total failure

④ results are ordered by relevance

④ not ordered by relevance

⑤ probabilistic model

⑤ deterministic model

⑥ results are approximate matches

⑥ exact matches

⑦ does not provide solution to user of database system.

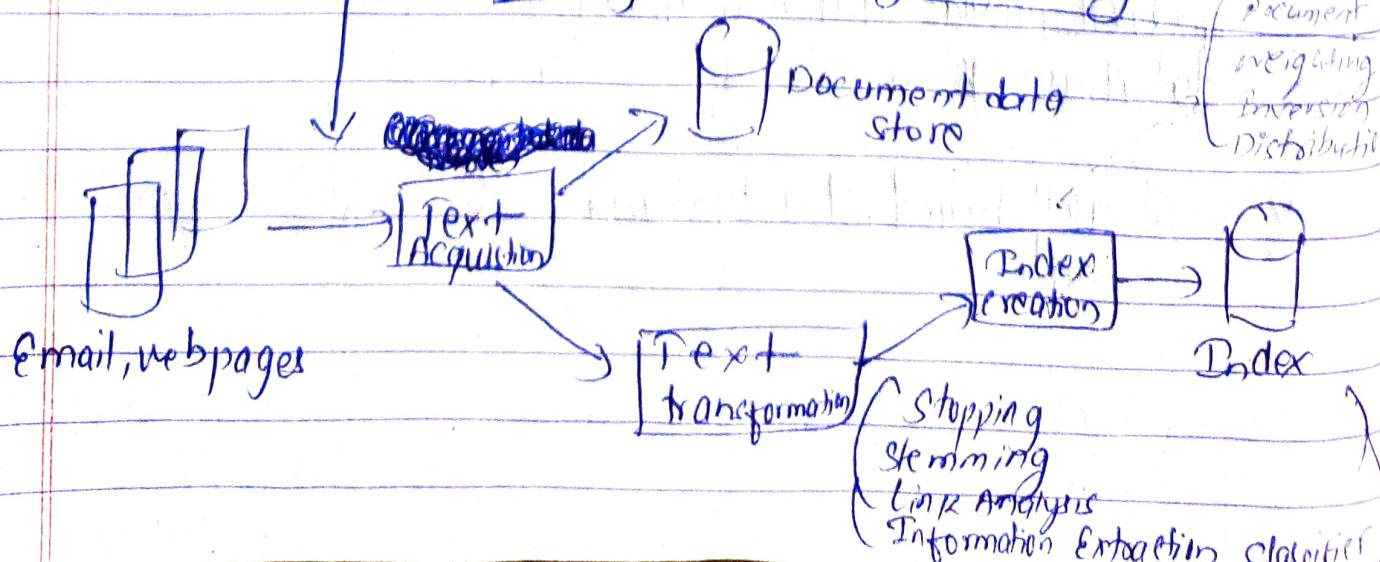
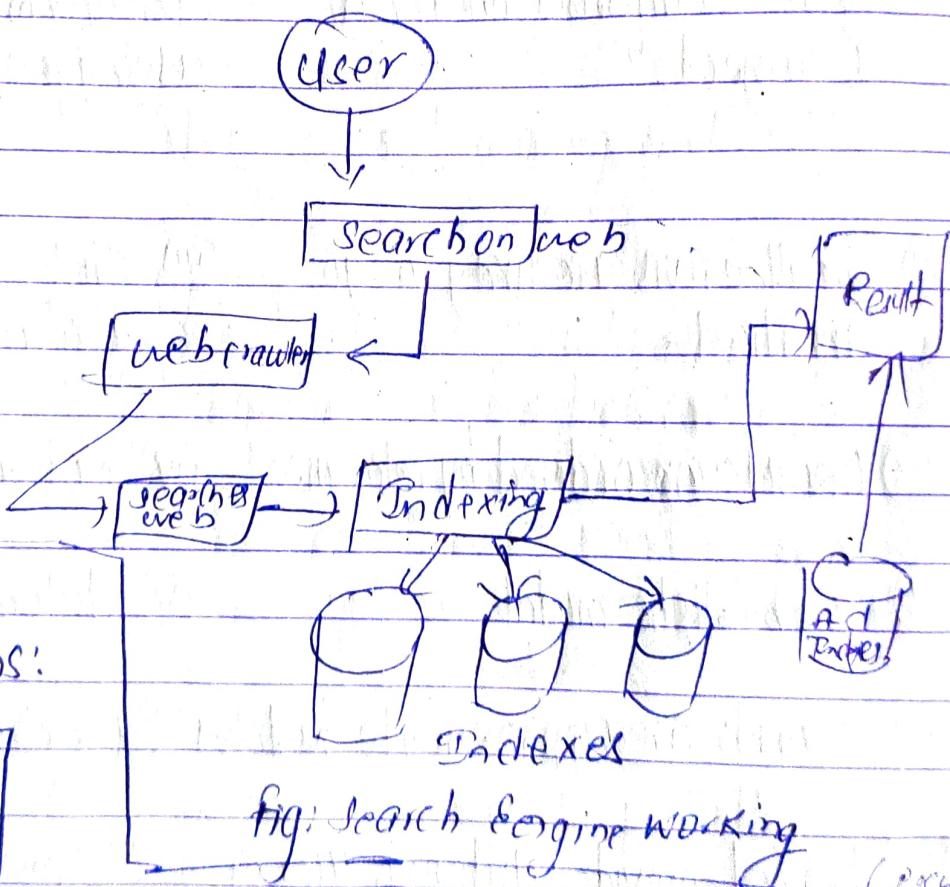
⑦ provides solution

⑧ not well defined structure and semantics

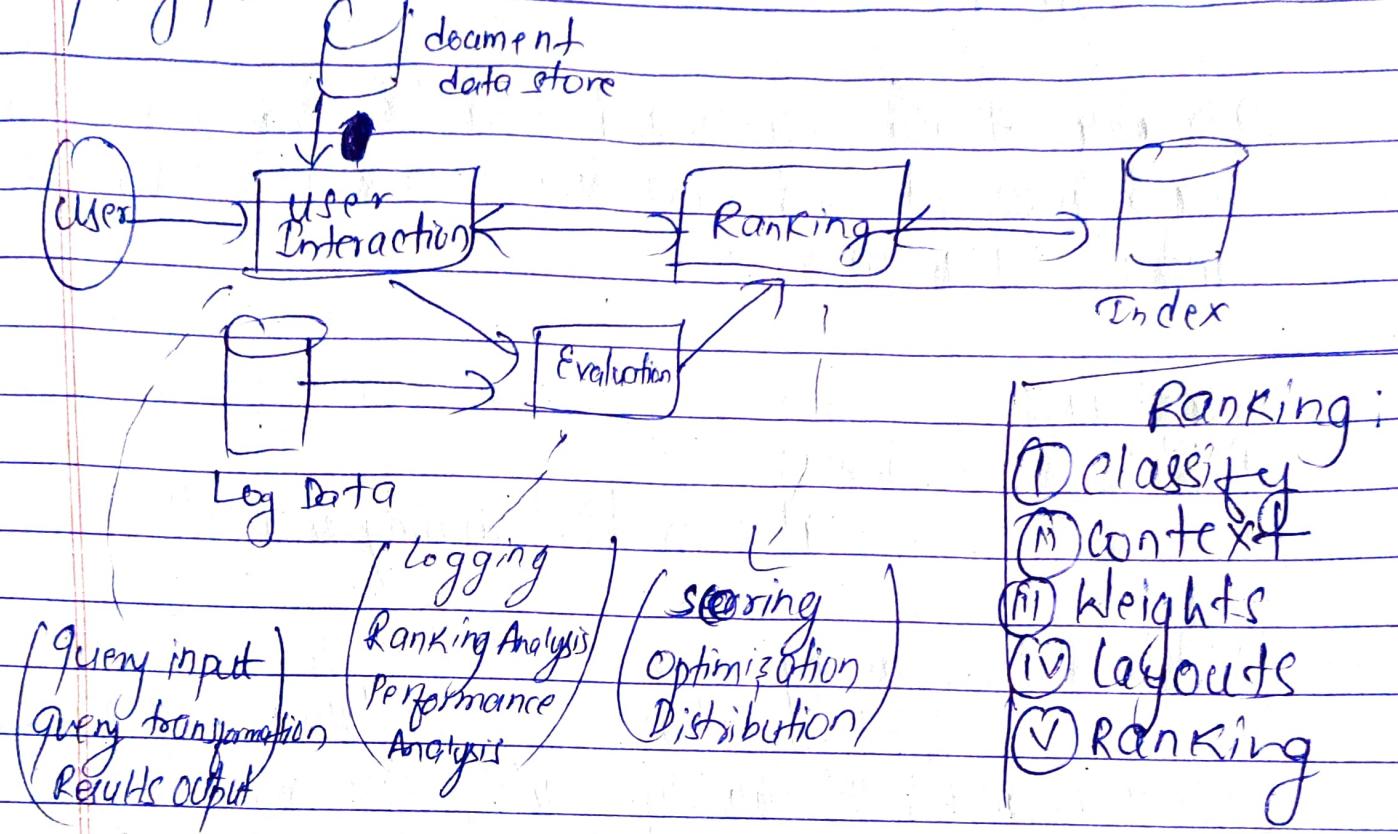
⑧ well defined

① Search Engine:

- the practical application of information retrieval techniques to large scale text collections
- Eg: Google, Yahoo



⑥ Query process:



⑦ Search Engines Features:

① Performance:

→ Efficient search and indexing

② Incorporate new data:

→ Coverage and freshness

③ Scalability:

→ Growing with data and users

④ Adaptability:

→ Tuning for applications

⑤ Specific Problems:

→ e.g. spam

① Web Mining:

→ To automatically discover and extract information from web documents and services i.e. to extract from www and its usage patterns

→ Techniques:

① Web Content Mining:

It relates to text mining and NLP.

② Web Structure mining:

Used to extract the structure of web pages e.g. pages as nodes and hyperlinks and edges connecting related pages, → describes the structure summary of a website.

→ to determine the connection between ^{to} commercial websites, web structure is useful.

③ Web Usage Mining:

→ Identifying interesting usage patterns from large datasets

→ It determines user behaviour

→ Users can collect data from logs, so it is called log mining

④ Applications:

① Help improve power of search engines by classifying web documents.

Date:

Page:

- ② Used for scheduling e.g. google, yahoo and ~~vertical~~
- ③ Predict User behaviours
- ④ Landing Page Optimization in particular website and e-service
- ⑤ Tools include Apache, Scrapy, logs.

⑥ Text Mining:

→ Essentially a subfield of data mining as it focuses on bringing structured to unstructured data and analyzing it to generate novel insights.

→ Types of data:
① Structured data → (SVM, Naive Bayes)
② Unstructured data
③ Semi-structured data

→ Steps involved:

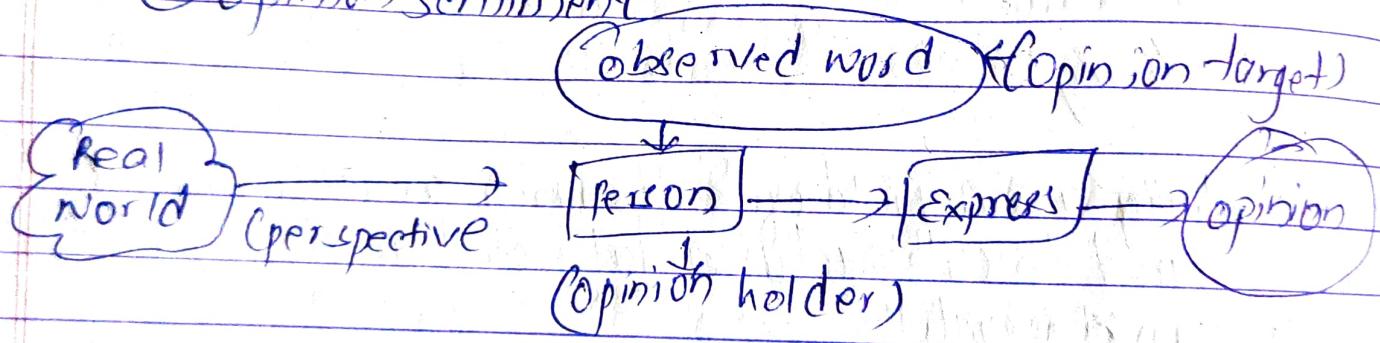
- ① Gathering Unstructured data (Info Extraction)
- ② Remove Anomalies (Information Retrieval)
- ③ Convert Unstructured to Structured data (↓)
- ④ Analyze the patterns within the data (Categorization and Clustering)
- ⑤ Store all the valuable information (Summarization) in a database

① Opinion Mining:

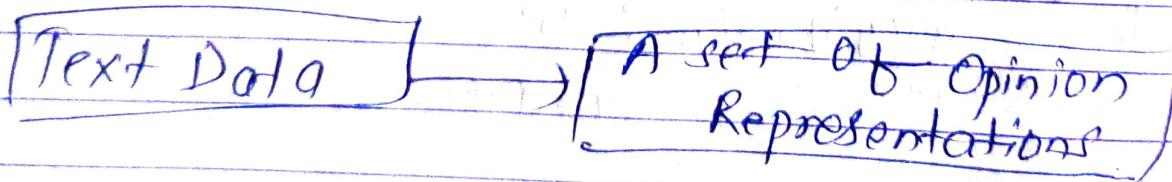
↳ a subjective statement describing what a person believes or thinks about something

Opinion representation

- ① Opinion holder
- ② Opinion target
- ③ Opinion content
- ④ Opinion context
- ⑤ Opinion sentiment



→ Task :



→ Why Opinion Mining:

① Decision Support

② Understand People

③ "Voluntary Survey"

① Social Web Mining:

→ the process of mining social data

→ involves collection, processing and analysis of raw data obtained from social media platforms such as fb, Instagram, Twitter etc. to uncover meaningful patterns and trends, draw conclusions and provide insightful and actionable information

② Data Collected:

- Age → clicks
- Gender → No. of followers
- location → likes

③ How it works:

→ involves:

- ① Combination of ① Statistical techniques
- ② ② Mathematics
- ③ ③ Machine learning

④ Used by ① Government Agencies

② Companies

③ Hotels

④ Retailers

⑤ Airlines

④ Used for:

① Trend Analysis

② event Detection

③ Social Spam detection

④ E-commerce

⑤ Brands

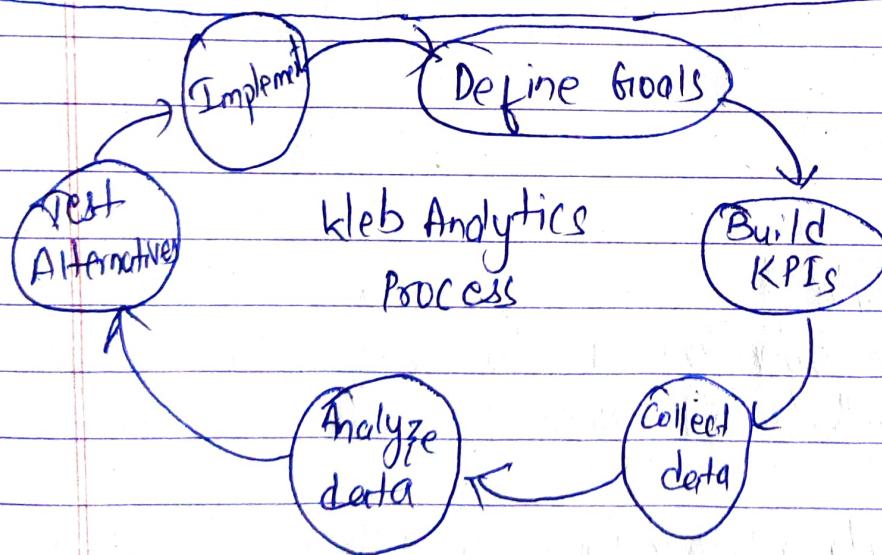
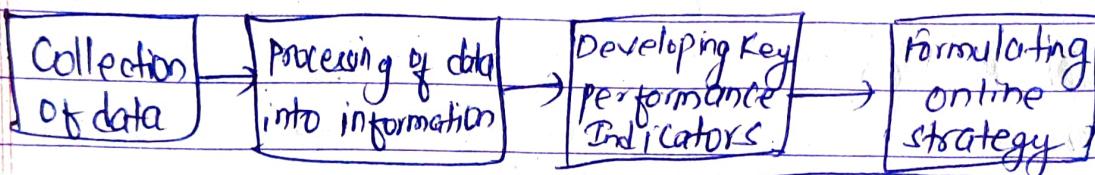
⑥ Research purposes

- Data Mining
- ⑥ Web Mining
 - ⑦ Extract info from web documents
 - ⑧ Extract info from large datasets
 - ⑨ Very useful for particular web site.
 - ⑩ Very useful for web page analysis
 - ⑪ access data publicly
 - ⑫ access data privately
 - ⑬ information from structured, unstructured and semi-structured web pages
 - ⑭ information from explicit structure
 - ⑮ Special tools like Scrapy, Page Rank
 - ⑯ ML algorithms
 - ⑰ Web Content Mining, Web Structure Mining
 - ⑱ Clustering, Regression etc

① Web Analytics:

- process of analysing the behaviour of visitors to a website
 - quantifiable measure
 - bounce rates,
- (i) Unique users,
 (ii) User sessions,
 (v) On-site search queries

② Process:



③ Terminology:

(i) Site References: Means from which website or source your readers come from

(ii) Keywords and Key Phrases:

(iii) Building Block Terms:

(Page, Page Views, Visits/Sessions, Return, Unique Visitors, New, Repeat)

① Visit Characterisation Terms:

- ⑩ Entry page
- ⑪ Landing Page
- ⑫ Exit Page
- ⑬ Visit duration
- ⑭ Referrer
- ⑮ Click through / click-through rate
- ⑯ Page View per visit

② Content Characterisation Terms:

- ⑰ Page Exit ratio
- ⑱ Single Page Visit
- ⑲ Bounce Rate

③ Conversion Metrics:

- ⑳ Events
- ㉑ Conversion

→ Types:

① Off site Web Analytics:

→ practice of monitoring visitor activity outside of an organization's website to measure potential audience

→ focuses on data collected from across the web, such as social media, search engines and forums.

→ provides industry wide analysis

③ Onsite Web Analytics:

→ Analytics to track the activity of visitors to a specific site to see how the site is performing

→ One's own site

④ Web Analytics Platforms:

- ① Google Analytics
- ② Adobe Analytics
- ③ Stat Counter
- ④ Mixpanel
- ⑤ Baidu Analytics
- ⑥ Crazy Egg
- ⑦ Hubspot

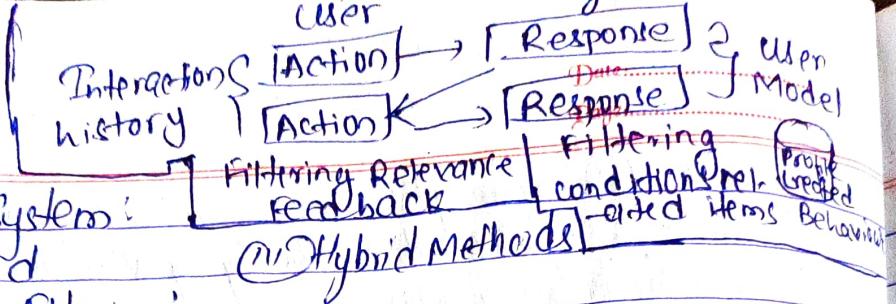
⑤ Need / Advantages:

- ① To understand the traffic
- ② To create better user experience
- ③ To understand customer behaviour
- ④ To track effectiveness of online marketing campaign
- ⑤ Help to understand ROI of social media campaign
- ⑥ Improve SEO
- ⑦ Identify pain points

⑥ Disadvantages:

- ① Time consuming
- ② Costly
- ③ Not 100% Accurate

→ User prof
interest



① Recommender Systems:

① Content Based

② Collaborative filtering:

- ↳ methods that are based solely on the past interactions recorded between users and items in order to produce new recommendations

↓

→ Two classes:

① Memory based → User-User
(Less Scalable) ↳ Item-Item

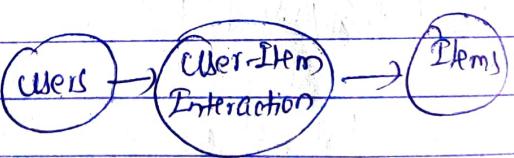
② Model based [Matrix factorisation]

→ Advantage:

① no pre information about user or item

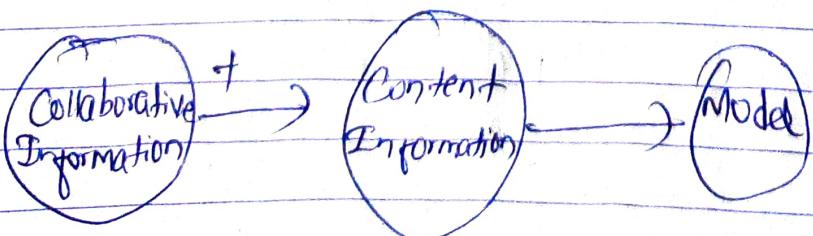
→ Disadvantage:

① "Cold Start Problem"



② Content Based Method:

→ build a model based on the available features, that explain the observed user-item interactions



→ Advantage

① Suffer less from cold start problem

→ Disadvantage:
→ totally new user can hard to adapt

① Java

② J

Driver

connecti

stati

reult

→ JDBC

-systems

③ Features

④ standa

⑤ datab

⑥ stor

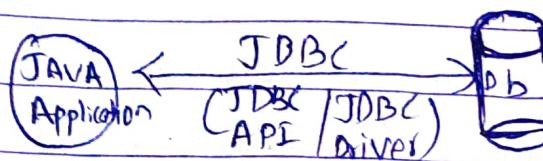
→ User profiling building refers to identifying user's interests and domains of relevance

Date.....

Page.....

① JDBC : (Java Database Connectivity)

→ A technology which is used to connect Java Application with database



→ part of JDK software

→ is an API which provides various classes and interfaces by which we can connect Java Application with DB.

→ JDBC API has two packages:

① `java.sql`

① ↓ ↓①
Driver Date
Connection Time
Statement Timestamp
ResultSet Rowset Listener
Types

② `javax.sql`

② ↓ ↓②
Datasource RowsetEvent
Rowset ConnectionEvent
Rowset Listener

⇒ JDBC is an abstraction which is provided by sun Microsystems and implemented by Database Vendors

① Features:

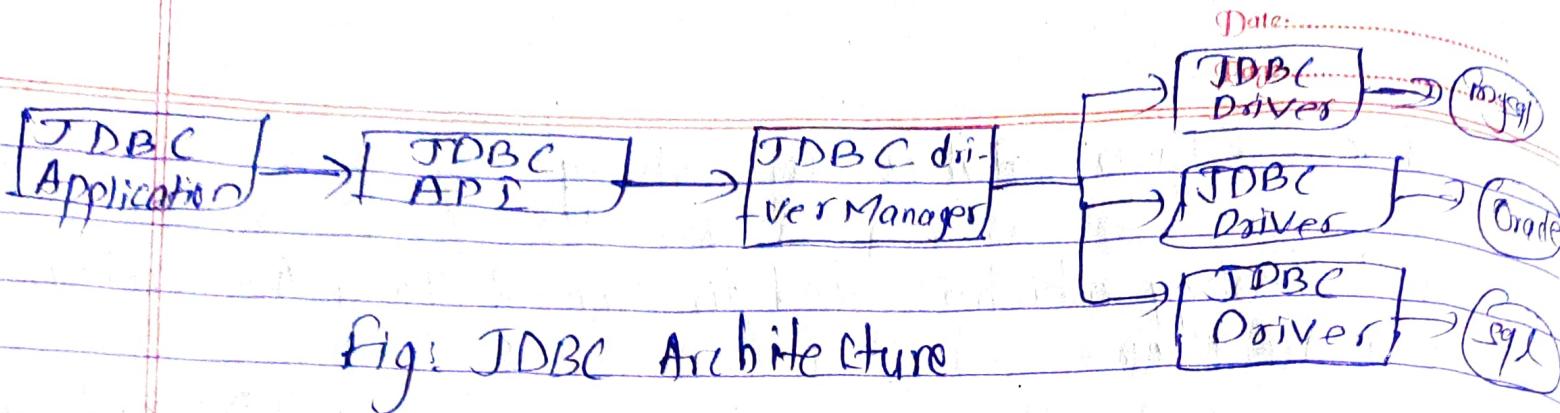
① Standard API

② platform independent technology

③ Database Independent API

④ CRUD or SCUD

⑤ stored procedure, triggers, views



① Why Use JDBC?

→ before functions were appended in front-end files provided by database vendors to connect to the database

② JDBC

① both platform and database independent

② Implemented in Java without pointer

③ Introduced by SUN Microsystems in 1993

④ performs comparatively better with Java applications

⑤ object oriented

⑥ stands for Java DB Connectivity

ODBC

① platform dependent and database independent

② Implemented in C with pointer

③ Introduced by MS in 1992

④ performs comparatively slower

⑤ procedural

⑥ stands for Open DB connectivity

- Q Steps to write JDBC program:
- ① Load the JDBC driver class or register the JDBC driver
 - ② Establish the connection
 - ③ Create a statement
 - ④ Execute the SQL commands on database and get the result
 - ⑤ Print the result
 - ⑥ Close the connection

⑥ ODBC:

- Open Database Connectivity
- standard database access method developed by MS Corporation
- makes it possible to access data from any application, regardless of which database management system is handling them
- based on CLI (Call-Level-Interface) specifications
- uses SQL as its database access language
- functionality of ODBC involves insertion of middle layer called a database driver, between application and DBMS
- purpose of database driver is to translate application's queries into commands that DBMS understands
- To access an ODBC db, you must have the appropriate driver for the db you wish to access
- To access a data source, connection to the data source must be first be established.

⑦ ODBC Architecture:

→ Four components:

① Application:

performs processing and calls ODBC functions to SP1 statements and retrieve results

(ii) Driver Manager: → loads and unloads driver on behalf of an application.
 → processes ODBC function calls or passes them to a driver

(iii) Driver: → processes ODBC function calls, submits SQL requests to a specific data source and returns results to the applications.

(iv) Data Source: → consists of the data user wants to access and its associated OS, DBMS, and network platform (if any) used to access the DBMS.

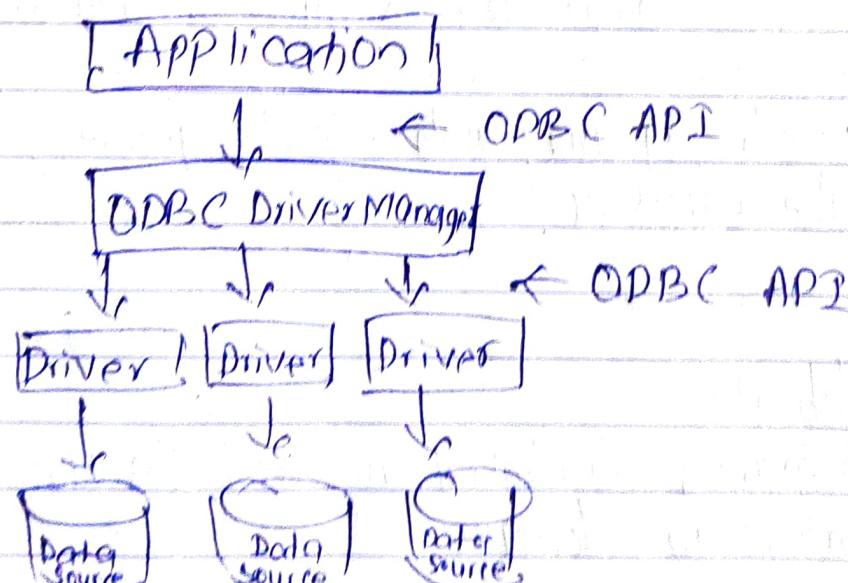


Fig: ODBC Architecture

⑤ Web Pages:

→ Single document written using HTML

(is a text, word, button, icon, which when clicked or hover on it, you will be directed to another web page)

→ To access a web page by entering its URL address using a web browser.

A typical URL looks like:

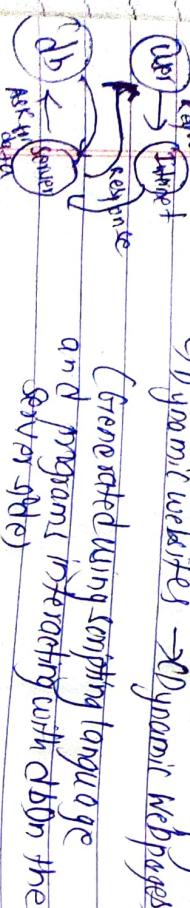
http://www.example.com/index.html

→ A web page may contain texts, graphics, hyperlinks and many other resources

⑥ Web sites:

→ A group of inter-linked and well-structured pages that exist on the same domain

→ Two types:
① Static websites → (Static web pages) the browser exactly as it is stored on server
② Dynamic websites → (Dynamic webpages) (not stored until changed manually)
(Generated using scripting language and programs interacting with db on the server)



⑥ Web Applications:

→ Software programs that exist on the server and runs using a web browser, through a web page

→ Created using a combination of web programming languages and web application frameworks

→ Uses RAM, allows user interactivity and is designed for many users

→ Google Mail, fb, YouTube

→ Web Applications

① based on user engagement
② Almost all content from user

③ Similar to desktop apps
④ more complicated and more skills - can websites do

Dynamic websites

① Allow the user to interact and content
② Create content is dominant
③ less complicated and skills needed are less

② Web Application development frameworks:

① Django:

(MVT)

→ high speed MVC framework for building web apps based on python

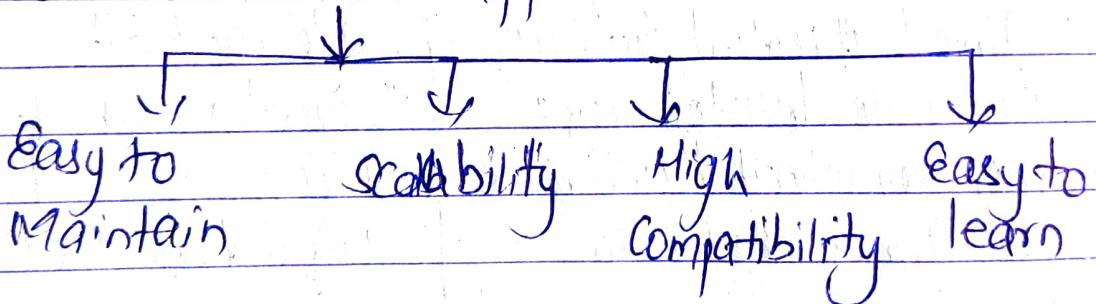
links the
M
odel architecture
(the database)

V
Presentation
(layer)
User interface of
the web application

F
Model
(layer)

→ excels at creating
web sites that can handle
high amount of traffic and
transactions

→ A web framework is a software framework that is created to support the development of dynamic sites, web services and web applications.



③ Features:

- ① Vault Community
- ② Open Source
- ③ Enhance Security
- ④ Rapid Development
- ⑤ Versatile

③ Companies:

- ① Instagram
- ② pininterest
- ③ Firefox
- ④ Netflix

⑥ Ruby on Rails:

→ Ruby:

- Popular programming language
- Acts as the fundamentals for working with framework like Ruby on Rails.
- Emphasizes on ensuring the language is expressive and easy to understand
- Free format language
- Case sensitive
- dynamic programming language

→ Ruby on Rails is an MVC-based framework used for web development and app programming at the server-side of the applications

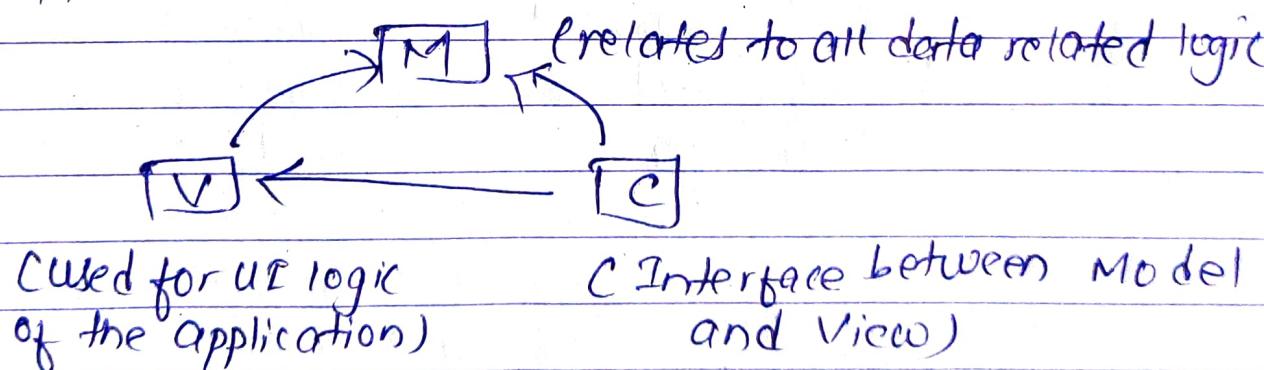


Fig: MVC architecture (Separates application layer into M, V, C)

→ Ruby on Rails developers prefer it for its time-saving qualities while writing code.

→ Offers ready-made code libraries for tasks such as building

forms, menu, and tables

⑥ Benefits:

① Cost effective

② Speed

③ Easy change management

④ Secure

⑤ Efficiency

⑥ Larger Developer Community

⑦ Responsibilities:

① Web App Development

② Data Integration

③ Server side functionality

④ Creation on Backend

⑤ App Connection with Services

⑥ Offering Solutions

→ Client side Technology:

Date:.....
Page:.....

- ① HTML (5) : (A standardized system for tagging text files to achieve form, colour, graphic, and hyperlink effects on World Wide Web pages)
- <html> text files to achieve form, colour, graphic, and hyperlink effects on World Wide Web pages
<head>
<title> <title>
<script>
~~~~~  
<script>  
</head>  
<body>      (because it is the set of markup symbols or codes inserted into a file intended for display on the internet)  
</body>  
</html>

## ② HTML Forms:

```
<form action="action.php">
<label for="fname">First Name! <label> <br>
<input type="text" id="fname" name="fname" value="John"> <br>
<input type="submit" value="submit">
</form>
```

First Name:

John

Submit

→ forms attribute:

① ~~method~~ action

② target = "-blank", "-Self", "parent", "top", "framename"

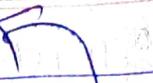
③ method = "get", "post", "put", "delete"

④ autocomplete = "On" ⑤ novalidate

① <select id="cars" name="cars">

<option value="Volvo"> Volvo </option> → (can use size

</select>



attribute to control no. of dropdown displayed)

② <textarea name="message" rows="10"

cols="30">

</textarea>

→ (multiple: for selecting multiple options)

\* ③ button:

<button type="button" onclick="alert('HelloWorld!')">  
Click Me! </button>

④ <input type="radio" id="html" name="fav\_language" value="HTML">

→ (checkbox)

HTML

⑤ <input type="date" id="birthday" name="birthday">

⑥ <input type="image" src="img\_submit.gif" alt="Submit" width="48" height="48">

⑦ <input type="file" id="myfile" name="myfile">

⑥  ⑦ Lists: 

- > 101
- > 102
- > 103
- > 104
- > 105

⑧ About links:

⑨ Absolute URLs: <a href="http://www.w3.org/"> W3C </a>

⑩ Relative URLs: <a href="/css/default.css"> CSS Tutorials </a>

⑪ Img: 

⑫ button: <button onclick="document.location='default.asp'"> HTML tutorial </button>

⑬

Tables:

<Table>

<tr>

<th> . . . </th>

<Table>

- |                                                                                                                                                                                                  |                                                                                                                                                                                                                                                       |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p><b>Advantages</b></p> <ol style="list-style-type: none"> <li>① Reduces file transfer size</li> <li>② Easy to customize online page</li> <li>③ Need few lines to improve site speed</li> </ol> | <p><b>Disadvantages</b> which creates confusion</p> <ul style="list-style-type: none"> <li>→ multiple levels which creates confusion</li> <li>→ cross-browser issues</li> <li>→ scarcity of security</li> <li>→ confusion for web browsers</li> </ul> |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

⑩ **HTML CSS:** → used to format the layout of a webpage friendly  
 → can be added three ways:

- ⑪ **less**
- ⑫ **Complex**
- ⑬ **Inline: style**
  - ⑭ **Internal: <style>**
  - ⑮ **External: using <link>**

## ⑩ JavaScript:

→ a light-weight, interpreted, or just-in-time compiled programming language with first-class functions. Both server-side and client-side language

→ for External JS:

```
<script src="myScript.js"></script>
```

Eg: <script>

```
function myFunction()
```

\$

```
document.getElementById("demo0").innerHTML = "Paragraph changed";
```

g

</script>

<body>

```
<p id="demo0">A Paragraph</p>
```

```
<button type="button" onClick="myFunction()>Try it</button>
```

</body>

⑩ **window.alert()**

⑪ **console.log()**

⑫ **(x == y)**

⑬ **const cors**

⑭ **const objc**

⑮ **let x = n**

function  
{  
 return  
};

⑯

⑰

⑱

⑲

⑳

㉑

→ let does not allow multiple initializations

Date.....

Page.....

① window.alert()

② document.write()

③ console.log()

④ window.print()

⑤ (x == y)

⑥ const cars = ["saab", "Volvo", "Bmw"]      ⑦ for( ; ; )

{ }

⑧ const object = { firstName: "John" }      do sth;

age: 50  
};  
};

⑨ let x = myfunction(4, 3);

⑩ object.name.firstName  
or

function myfunction(a, b)

{ return a \* b }

};

objectname["firstName"]

changed");

⑪ Events:

⑫ onChange

⑬ onClick

⑭ onmouseover

⑮ onmouseout

⑯ onKeyDown

⑰ onLoad

⑪ if...else

e.g.

if(hour < 18)

{

greeting = "Good Day"

};

② switch (value)

? case x: --- break;

? case y: --- break;

default:

## ① Class:

```
class Car
```

```
{
```

```
constructor(name, year)
```

```
{
```

```
this.name = name;
```

```
this.year = year;
```

```
{
```

```
}
```

```
let mycar = new Car("Ford", 2014);
```

## ② JSON:

→ JavaScript Object Notation

→ lightweight data interchange format

→ Inside curly braces,

## ③ XML schemas:

① describes the structure of XML document

② XML schema definition (XSD)

```
<?xml version="1.0"?>
```

```
<xsd:schema xmlns:xsd="http://www.w3.org/2001/XMLSchema">
```

```
<xsd:element name="note">
```

```
</xsd:element>
```

```
</xsd:schema>
```

Child  
elements

→ XML based and more powerful than DTD

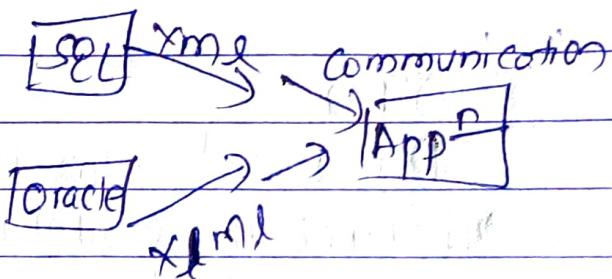
→ hundreds of XML formats are in daily use

## ① XML:

- Extensible Markup Language
- used to store and transport data
- self-descriptive
- used to carry data
- self-defined tags
- platform and language independent
- helps in easy communication b/w two platforms

### → Features:

- ① Separates data from HTML
- ② Simplifies data sharing
- ③ " " " " transport
- ④ Increases data availability
- ⑤ Simplifies platform change



### → XML example:

↳ Hierarchical structure

```
<?xml version="1.0" encoding="ISO-8859-1"?>
```

`<College>` ← Root element is must

`<Class1>`

Child  
element

`<Name>Amit</Name>`

`<Class1>`

`<College>`

→ Specify restrictions on data → XML Schemas → Keeps much better control over types of data → DTD → — → — → does not allow creating customized data

## HTML

- ① Display data
- ② Markup language

XML

- ① Transport and store data
- ② provide framework to define markup language

- ③ Not case sensitive
- ④ pre defined tags
- ⑤ static

- ⑥ Case sensitive
- ⑦ can create own tags
- ⑧ Dynamic

## ⑥ DTD:

⇒ Document Type Definition

⇒ defines the structure and legal elements and attributes of an XML document

⇒ With a DTD, independent groups of people can agree on a standard DTD for interchanging data

⇒ An application can use a DTD to verify that XML is valid.

⇒ Can be declared inside XML file or outside

```
<?xml version "1.0"?>
```

```
<!DOCTYPE note>
```

```
<!Element note (to,from,heading,body)>
```

```
>
```

```
<!Element to (CDATA)>
```

<note>

```
<to>Tove</to>
```

```
<note>
```

## ① TF-IDF :

- Stands for Term Frequency-Inverse Document frequency
- two matrices that are closely related and search and figure out the relevancy of a given word to a document given a larger body of document.
- All TF means is how often a given word occurs in a given document so within one web page one wikipedia article, how common is a given word within that document, what is the ratio of that word occurrence rate throughout all the words in that document that's it.
- DF means how often a word occurs in an entire set of documents i.e. all of wikipedia or every web page.

## ② IDF :

- Inverse Document frequency
- We use log of IDF since word frequencies are distributed exponentially.

$$\text{TF-IDF} = \text{TF} * \text{DF}$$

→ how often the word appears in a document over how often it just appears everywhere.

→ measure of how important and unique this word is for this document.

→  
Spe

سیده زین تبار

**B**uilt-in Functions in JS:  
⑥ `parseFloat` ⑦ `isNaN`

## ② PageRank Algorithm :

① `parent!` ② `Array` ③ `String(object)` ④ `!N`



$$PR_{t+1}(P_i) = \sum_j \frac{PR_t(P_j)}{CC_j}$$

$$PRCA = \frac{1/4}{1/2} = \frac{1}{2}$$

Go on similarly

$$P_{RC}(B) = \frac{1}{4} + \frac{1}{4}$$

- ⑤ HITS:
  - An algorithm used in link analysis.
  - It could discover and rank web pages relevant for a particular search.

**Hub**: A node is "Authority". A node is high-quality if it has many links to many high-quality nodes.

~~New hub is the sum of authority of its children~~

Normalize new authority and

Authority of every node in the graph  
⑦ New authority is the sum of hub of its parent.