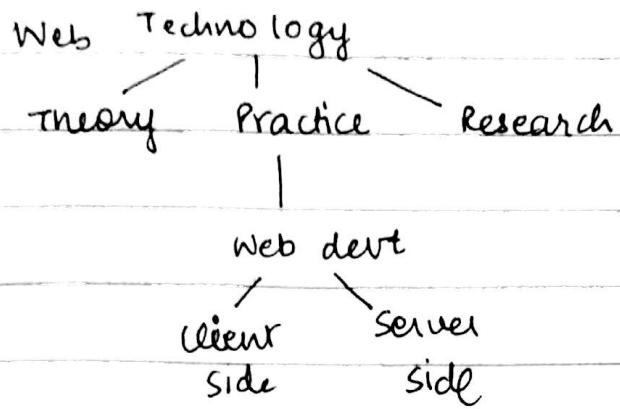


WEB TECHNOLOGY



frameworks → Django, Ruby on rails, JS Node.

<u>Data</u>	<u>Information</u>	<u>Knowledge</u>
(chunk, raw unstructured)	(specific, structured)	(what we extract)
- details of every student	- DOB of every student	- How many are adults vs minors.

INTER NETWORKING :

Acronym

- Need of networking :
- (1) communication
 - (2) sharing resources
 - (3) sharing information.

Metadata : sth which tells us about data. (data about data)

Meta: replicate of sth - eg. Data about data (meta-data)
Network of networks. (meta-n/w)

protocols

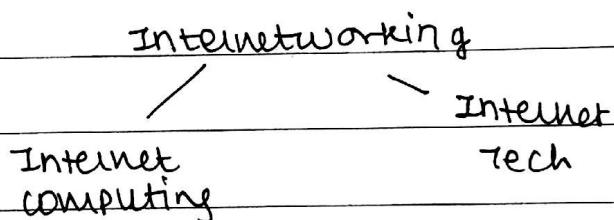
HTTP
SMTP
FTP
IMAP
VOIP

Internet:

- (1) Network of networks (aka metanetwork)
- (2) Collection of small LANs to form one big WAN.

components of Internet which help internet to work:

- (1) Client - server (Architecture)
- (2) Protocol
- (3) IP address
- (4) DNS
- (5) URL



Every service of In is either a client or server.

(who seeks info)

(provides service or responds to request of client).

90% of services provided by internet have client - server interaction

e.g. we access a website

↓
client

↓
server

Skype chat (client uses server to get address of other person)

Protocols

HTTP

Web

SMTP

Email

FTP

File transfer

IM

Instant messaging

VOIP

Internet Technology



SERVICES PROVIDED
BY INTERNET

Protocols: ^{Set of} Rules / regulations that help client & server interaction.

http:// www.

IP Address: Every machine on the internet has unique address. (address where sth is stored)

Microsoft: 207.46.92.254

We can't remember IP addresses, so we remember URLs.

~~Map~~

Mapping b/w IP address & URL is done by DNS.

DNS: Domain Name Server.

ANATOMY OF URL:

protocol:// host / path / filename.

URL: Uniform Resource Locator

Every host machine has unique address → IP address

Path of every file on that host → URL.

~~http://^{www.}dtu.ac.in / academics/syllabus~~

Protocol host ↗ path filename
(www.)

Hypertext Transfer Protocol

More than text

(videos, audios,

advertisements,

copyrights are
present APART

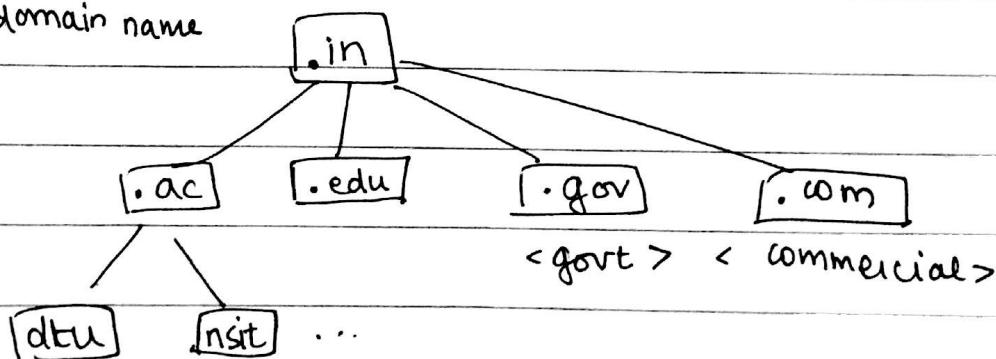
from text.

Multimedia + text)

^{2nd LEVEL} | ^{1st LEVEL} | TOP LEVEL DOMAIN NAME .in
www. dtu. ac. in .jp
Academic Country India .uk
Institution :

Hierarchical tree:

- Top level domain name
- 1st level
- 2nd level



http://www.dtu.ac.in/academics/syllabus
Protocol Host ~~Path~~ DIRECTORY Path Filename

Bookmark: URL stored by web browser.

Browser: Used to access URLs.

services provided by internet can be used only because of an ISP.

e.g. MSN, AOL

Intranet → Within an organisation. → Helps for training, leave access, tutorials, etc.

Extranet → Intranet + Access given to ppl outside organisation (extra clients < collaborating ppl >)

→ for client-side projects, ie, outsourcing.

e.g. we access Bank of America's website's server.

MAJOR ISSUE: higher security required.

Website - Collection of web pages.

e.g. Facebook (website), & our profiles are webpages.

IETF → Manages net

No one owns the internet → intangible entity (can't touch or feel)

Layers of Network



OSI Reference Model

Application: HTTP, FTP, VoIP are application layer protocol. Client & server ke beech mein protocol nota hai.

TCP / IP → connection protocol

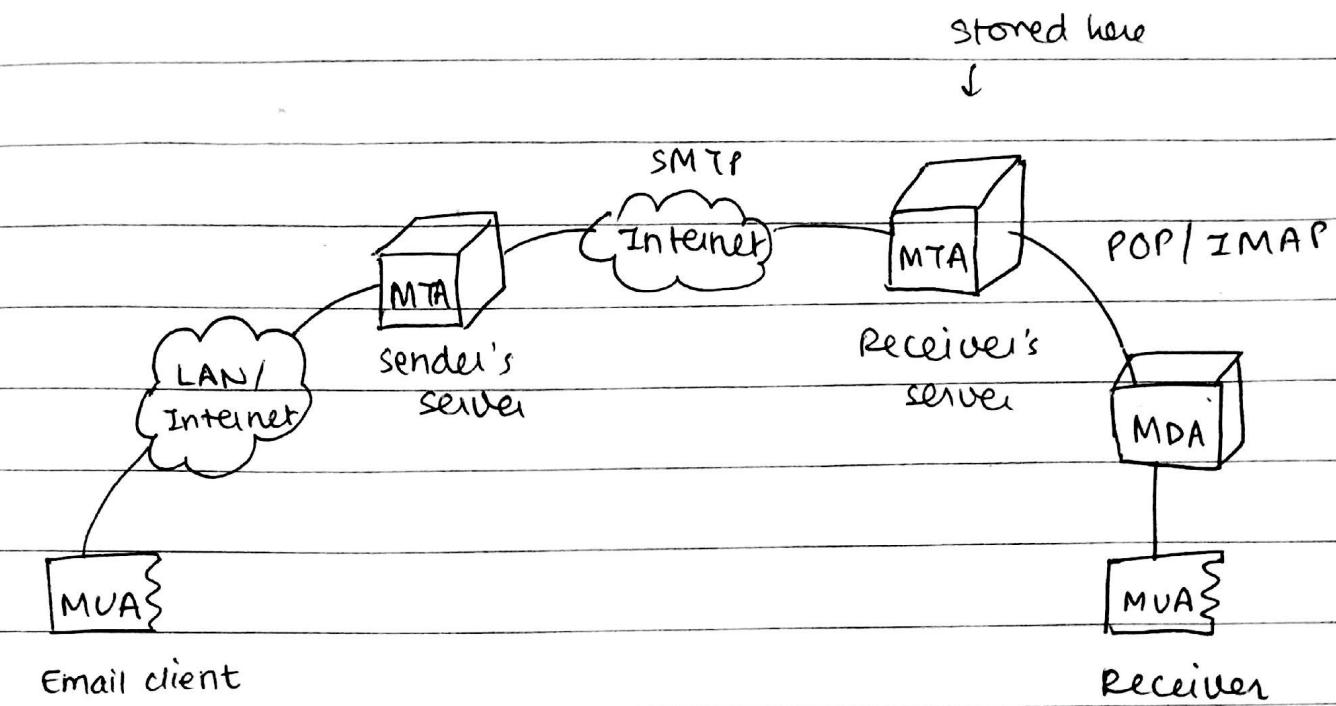
Transport
Router
Connection: TCP / IP

Network layer → which way / path to choose
physical → which cables / fibres to use.

Data link → packet.

Services① EMAIL :

Follows store-and-forward model



Mail User Agent w/ help of LAN / Internet sends to Mail Transfer Agent (sender's server)

SMTP used & POP / IMAP used → 2 protocols.

2nd MTA → Receiver's server

Mail Delivery Agent: jo deliver karega doosre ko.

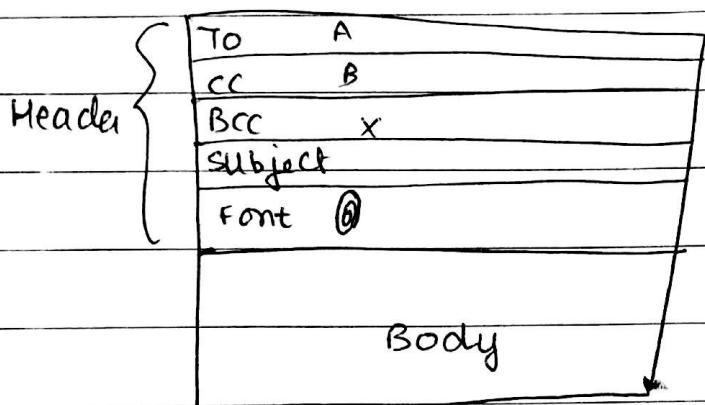
POP : Post Office Protocol

IMAP : Instant Message Access Protocol → maintains a copy of mail on server using IMAP

POP : Every msg when retrieved may not be stored.

- (1) Email Client
- (2) Web mail.

- (1) → On Android/ iOS phones → the app
- (2) → on a browser → using web address of email provider.



abc.def@gmail.com
 user @ domain

X can see A & B have been sent

A & B can't see X.

② FILE TRANSFER: Platform independent

Architecture → client - server, Protocol → FTP

Telnet / Usenet

Telnet: Permission to use someone else's identity or computer to access a resource. → providing remote login instead of user logging in. <Protocol hai>
eg. We access library's catalogues.

Usenet: News groups & discussion forum. Global bulletin board.

Torrent → Peer to peer.

VOIP → Hybrid (client-server AND peer to peer)
client tells server I want to video chat. Server
agrees & connects 2 peers (2 client machines).

for a common man, internet would be useless w/o web.
we use web for doing everything like web mail, etc.
∴ It is most important service provided by internet.

(3) HTTP (HyperText Transfer Protocol)

HTTP's characteristics to define it :

- (a) Application layer protocol
- (b) request - response protocol
- (c) stateless

Client - server architecture present.

NON-PERSISTENT HTML 1.0 : \Rightarrow ^{2 response times + 1} transmission time:

STEP 1 (a) Client initiates a connection to the server

(b) Server accepts the connection.

STEP 2 Client initiates a request msg.

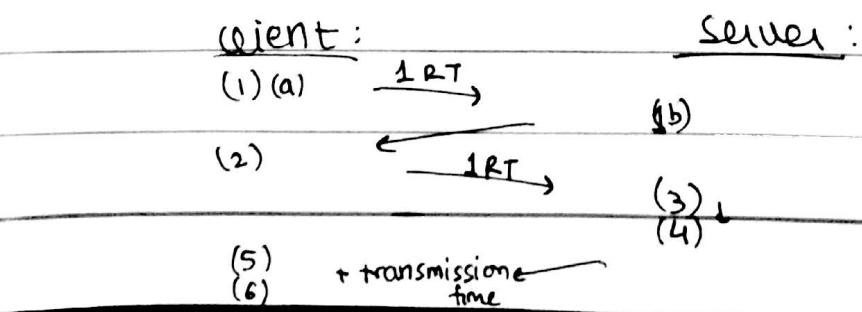
STEP 3 Server receives request. Forms a response msg.

STEP 4 ~~Server~~ closes the connection.

STEP 5 Receives the response (html file)*

STEP 6 Steps 1-5 repeated as many times as no. of referenced objects.

* if html file has 10 imgs \Rightarrow 10 objects. \therefore step 1-5 repeated 10x.

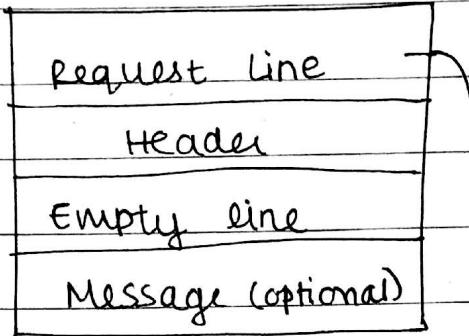


HTTP
 Non-persistent → HTTP 1.0 → $2rt + 1tt$
 Persistent without pipeline ① → $1rt + 1tt$
 with pipeline. → HTTP 1.1 ②

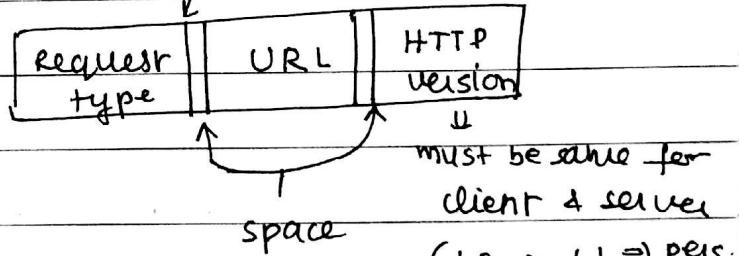
HTTP
 / \
 Request Messages response

①: sends request msgs continuously w/o checking if response is being sent ⇒ MIXING of msgs occurs. (Issue).

Request



②: send next request only one prs response is received.



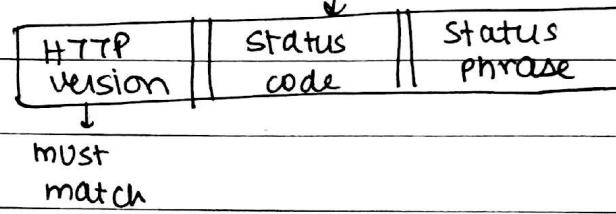
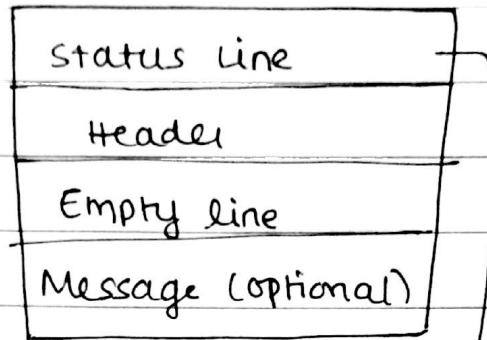
Methods / Commands:

GET, POST, DELETE, PUT.

Example:
 request type URL HTTP version
 GET / somedir / page. htm HTTP / 1.1

Header {
 Host : www. some school. edu
 User-agent : Mozilla / 4.0
 Connection : close
}

Response:



Status code: 1XX (One hundred series)

 2XX (2 hundred series)

3XX

4XX → Bad request

5XX

1XX → Informational series (provisional responses)

100 → If half → page loaded, tells it to CONTINUE

100 : CONTINUE

101 : SWITCHING

102 : PROCESSING

2XX → Success series

200 : OKAY → we don't see this ∵ it means

201 : CREATED website is displayed.

202 : ACCEPTED

Response msg } 404 → Page not found
from server } 301 →

Burk Pagers

3XX → Redirection series.

301 → Moved permanently

4XX → Client Error

400 → Bad request

408 → Request timeout

401 → Unauthorised

403 → forbidden

404 → Not found

5XX → Server error

500 → Internal server error

502 → Bad gateway

504 → Gateway ~~not found~~. timeout.

505 → HTTP version not supported

1 marker → status code & its use.

Example: HTTP / 1.1 200 OK

Header { connection: close

Date:

Server: Apache

Last modified:

Content-length:

Content-type:

Then msg.

PERSOCIAL NETWORK

Stateless: No state info of client stored in se
(of any previous request → not called by server)

Cookies: Not softwares, cannot be programmed.
or add
not malware & can't contain virus → not corrupting system
can be used as spyware to keep track of browsing activity/history. ∴ safer to remove them.

Cookies → to save info in a stateless system.

User used net for these purposes:

- (1) Informational
- (2) Navigational
- (3) Transactional

TYPES OF COOKIES:

(Based on above 3 purposes)

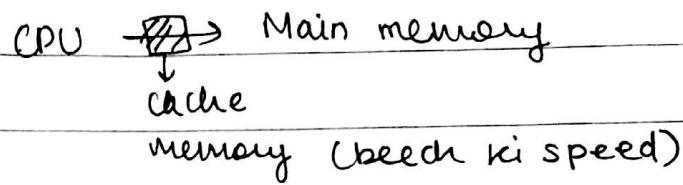
- (1) Session cookies
- (2) Persistent cookies → still exist after session.
- (3) Secure cookies → banking → ^{when} redirected to other pg.
- (4) Zombie cookies → recreated after user is gone

server overloading → ^{sewer} Bottleneck.

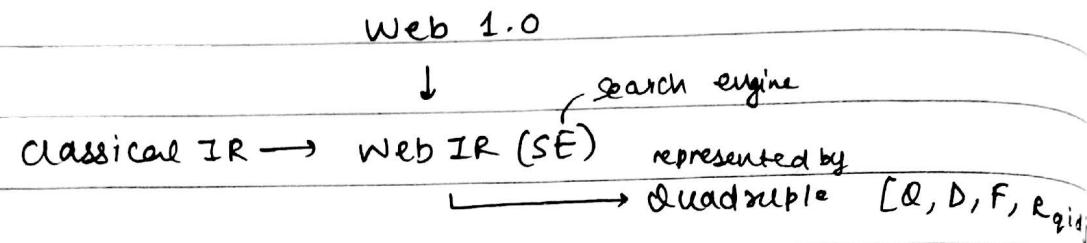
Proxy server → web cache.

Web caching.

Cache Memory → High speed memory which was employed b/w main memory & CPU ~~for~~ to compensate for speed difference b/w these 2. CPU → fast MM → slow.

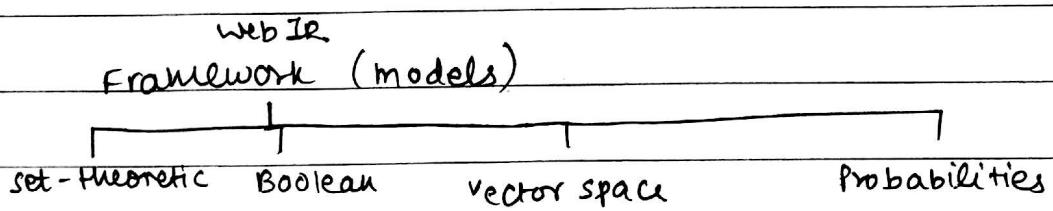


Web caching → Small servers used to save some user activity. Proxy servers → Name of these servers.
 ↓
 False server. tells you falsely that it is a server.



[query, documents, framework, ranking mechanism]
 which will be retrieved to sort materials

Database records SQL queries	Documents (web pages, things)



Search Engine ≡ Web IR

SEARCH ENGINE ARCHITECTURE (SEA) : 2 modules -

- (1) Indexing : → Text acquisition] OFFLINE
 → Text transformation] PHASE
 → Index creation

- (2) Querying:
- User Interaction
 - Ranking
 - Evaluation

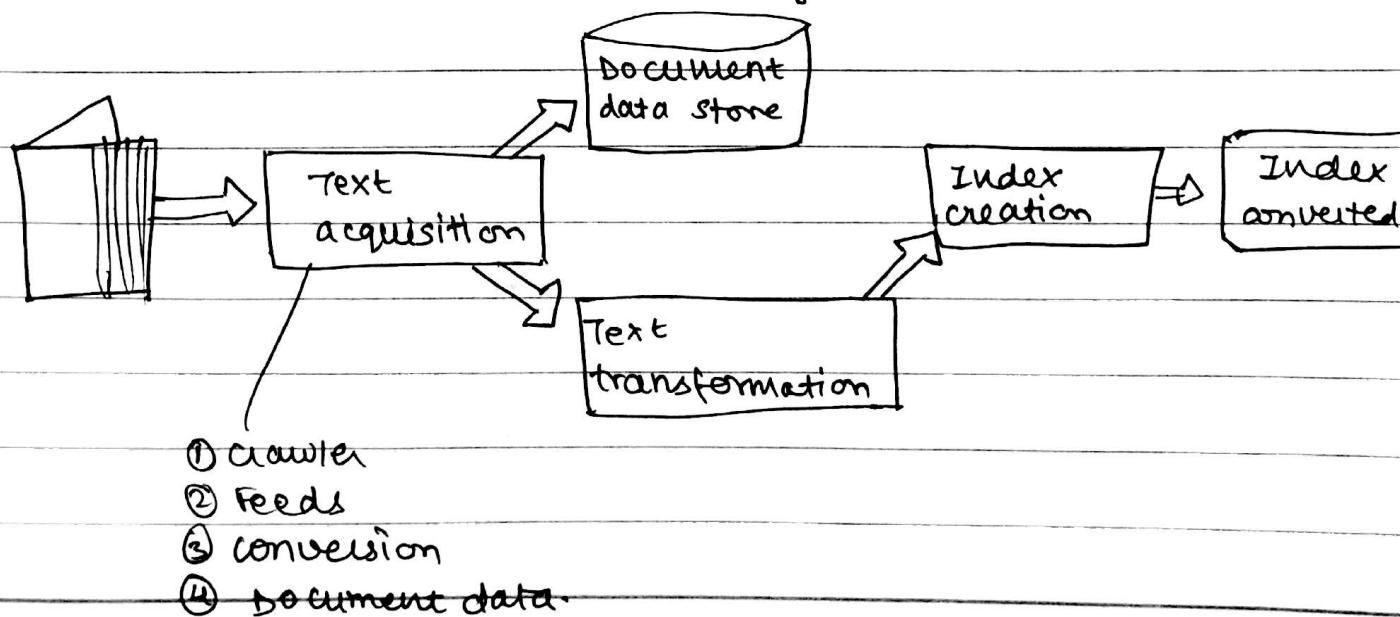
]} ONLINE
PHASE

Term	Present in - Doc ID
abc	11, 3, 4, 41, 200, 5001
def	3, 2, 17, 400, ...
xyz	...

Each Doc has its own table

DOC No.	URL	Title	Snippet
			which is seen.

→ Unstructured text (compressed)
→ metatags



Crawler - program which helps to crawl the web graphs to acquire the info.

vertical crawler: topic related / wordical searches.

RSS feeds: rich site summaries → realtime streaming of documents (like news / man-

∴ info retrieved through either crawlers / feeds.

conversion - all formats have to be converted to consistent formats.

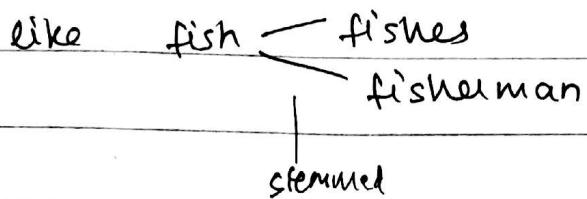
Text transformation: identifying imp. terms & features in query to help w/ indexing

e.g. how to get admission in dtu = dtu admission

5 substeps in t.t → ~~①~~ • parsing

- ~~②~~ • stemming
- ~~③~~ • stoppping
- ~~④~~ • link analysis
- ~~⑤~~ • info extraction
- ~~⑥~~ • classification, classifier

Stemming: Natural language processing ~~part-of-speech~~ procedure (transforming text)



Stopping: Removal of common words like punctuation or how, to, get, in, etc.

Info Extraction:

Name Entity Recogniser → identify nouns related to place, name, company, etc. To identify imp/indexed terms in query documents.

Classifier: Pre-defined set of classes, to categorise the data.

• Maintaining statistical data facilitate smooth ranking mechanism (online part)

Index creation: → Doc statistics → term freq (TF) calculated from keywords collected before.
→ Weighting → to store ranking using maintained statistics
→ Inversion
→ Index distribution.

Inverse document frequency (IDF)

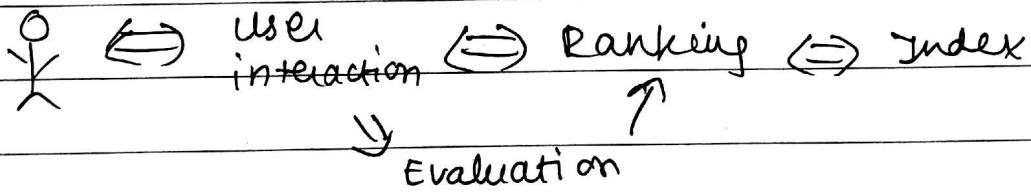
within the set of doc. what is freq of ~~term~~^{term.}

Google's pagerank works by counting how many links point to a particular webpage \Rightarrow rank it higher & show it on top.

TF: Term kitni baar aayi doc. mein.

IDF: Agar term zyada baar na aaye but page is overall more imp. so find an avg value.

Querying Doc. data store



input query } : 2-way
output result

- (1) input query
 - (2) query transformation
 - (3) output result
- suggestions \rightarrow DTU \Rightarrow (1) admin's
(2) placements etc
- spell-check
- expands acronyms (DTU)

results ↗ paid results (ads on top)
 ↗ organic results

Recall }
 Precision } evaluation's effectiveness measured in 2 ways

Recall → how many relevant docs retrieved from doc. set closely

Precision → how many are actually precisely relevant from total relevant retrieved docs from doc. set.

Efficiency measured in response time (google shows it too).

PageRank → PageRank™

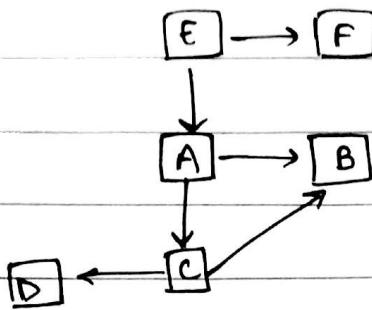
Made by Larry Page & Sergey Brin

Google sees web as a web graph & uses it as a RANDOM SURFER MODEL.

PageRank → Based on probability distribution \Rightarrow probability of user clicking on a particular link.

* PageRank ranks webpages & not websites.
 " works on in-links & out-links

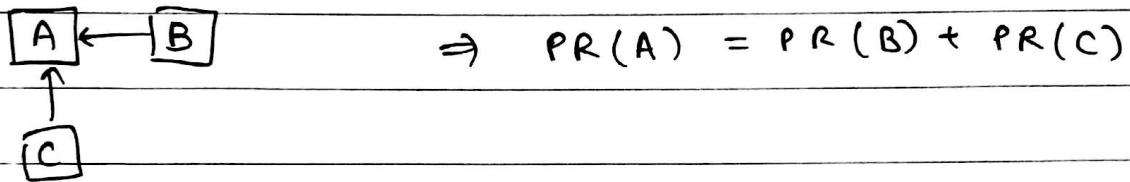
Each link that points to a page acts as a vote. More the pointing link, higher rank.



PageRank formula is :

$$PR(A) = (1-d) + d \left(PR\left(\frac{T_1}{C_1}\right) + \dots + PR\left(\frac{T_n}{C_n}\right) \right)$$

d: damping factor , $d = 0.85$

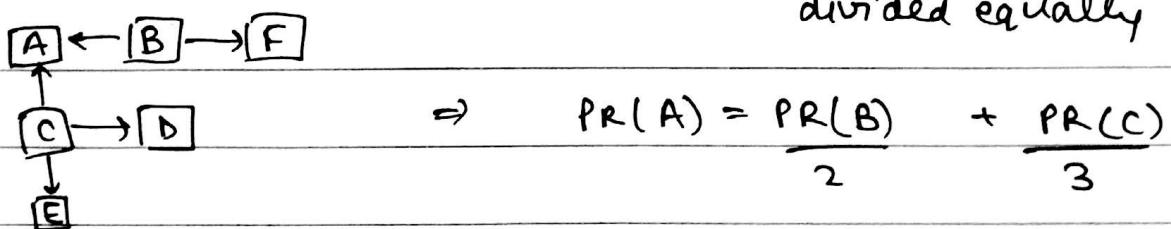


$$\Rightarrow PR(A) = PR(B) + PR(C)$$

$$\Rightarrow PR(A) = PR(B) + \frac{PR(C)}{2}$$

$\therefore C$ points to 2 others \Rightarrow vote share

divided equally



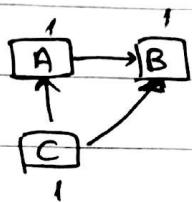
$$\Rightarrow PR(A) = \frac{PR(B)}{2} + \frac{PR(C)}{3}$$

damping factor was introduced : if no link points to a page, it should still be displayed

$$PR = (1-0.85) + 0.85(0)$$

$$= 0.15 \quad \checkmark \Rightarrow \text{seen.}$$

Initial probability is taken as 1 for all.
∴ it keeps changing ∵ all are co-dependent.



$$PR(A) = 1 - 0.85 + 0.85 \times \frac{PR(C)}{2}$$

$$PR(B) = 0.15 + 0.85 \left(PR(A) + \frac{PR(C)}{2} \right)$$

$$PR(C) = 0.15$$

Google says it only uses link analysis & not content analysis, but it does actually happen ∵ keywords are shown ⇒ May be it occurs in backend.

Q1 Explain Web & Internet.

Q2 Explain store & forward model of email.

Q3 Briefly describe search engine architecture.

Q4 Differentiate b/w Explain need of web caching.

Q5 Explain the following:

(a) URL (b) HTTP (c) HTTP 1.0