

DMMML, 28 Mar 2019

Information retrieval on the Internet

Internet documents are not "refereed"

Commercial value attached to search

Ranking

HTML makes it easy to add
invisible misleading content to
misdirect search

Self published documents do not contain

"obvious" terms

IBM webpage did not mention "computer"

Additional structure available - hyperlinks

Links within documents refer to other docs

` Some text `

↑
location

↑
anchor text

↖ describes contents

Can trust anchor text more than document itself

- Use anchor text to index target doc

Different approach to indexing

World Wide Web as a gigantic graph

- Reason about the graph as a whole

Social Network Analysis

e.g. Film industry

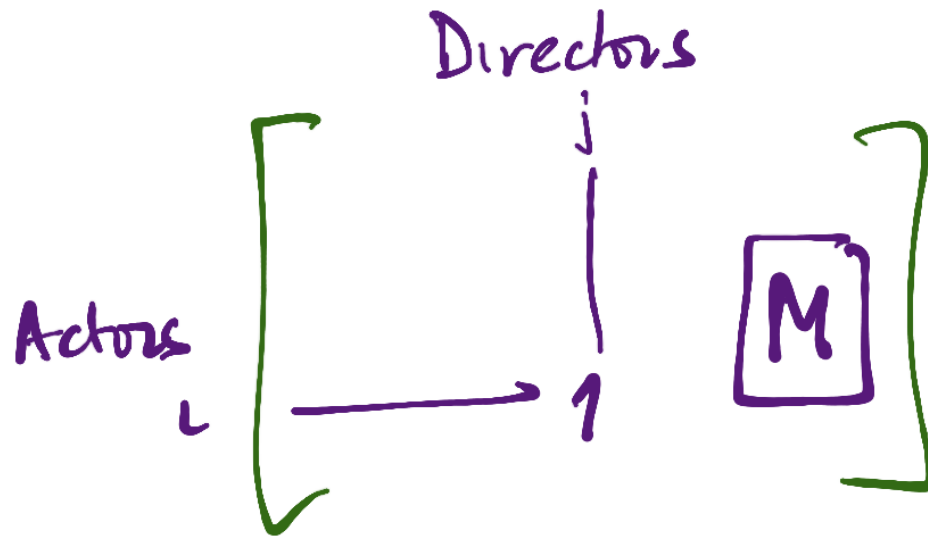
└ Actors — Star?

└ Directors — Famous?

Good actors want to work with famous
directors

Famous directors get stars to work in
their films

Graph.



$M[i, j] = 1$ if actor i has worked in
a movie of director j

Each actor has a star rating $S[i]$

Each director has a "famous" rating $F[j]$

Stars derive their rating from the directors they work with:

$$S[i] = \sum_j M[i,j] \cdot \underline{F[j]}$$

Symmetrically

$$\underline{F[j]} = \sum_i M[i,j] \cdot S[i]$$

$$S[i] = \sum_j M[i,j] \cdot \sum_i M[i,j] S[i]$$

$$S \approx (M \cdot M^T) \cdot S$$

$$F \approx (M^T \cdot M) \cdot F$$

Solve for $S, F \Rightarrow$ computes ratings

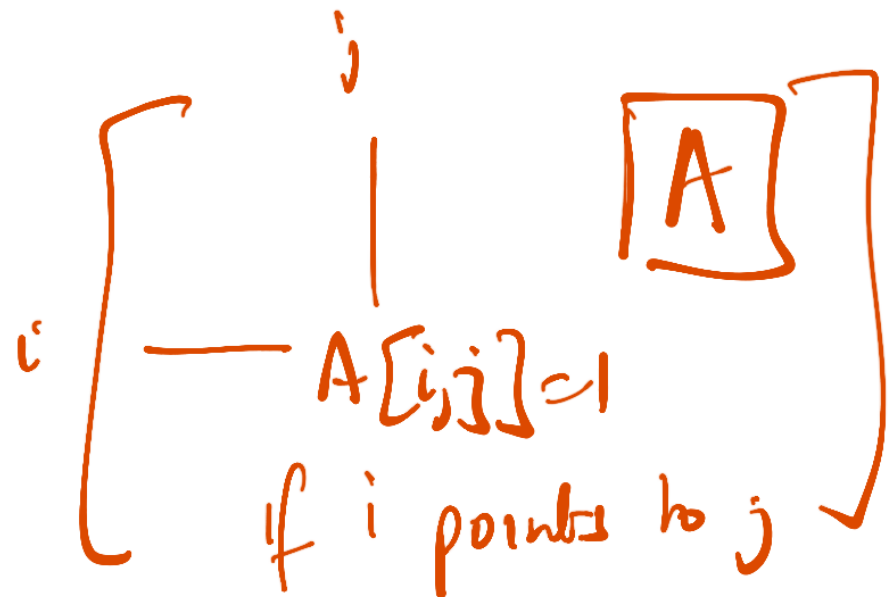
Similarly, use graph structure to derive
some conclusions about document ranks

Every document has some "prestige" : $P[i]$

$P[i]$ is shared (equally) across all outgoing
links

Document i derives prestige from incoming
links

Adjacency
matrix of
Internet



If document i has n outgoing links (i.e. n 1's in row i of A) - each gets $\frac{1}{n}$ of $P[i]$

3 documents.

$$\begin{array}{c} P \\ [1 \ 1 \ 1] \end{array} \begin{array}{c} A^* \\ \left[\begin{array}{ccc} 0 & 1/2 & 1/2 \\ 1 & 0 & 0 \\ 1/2 & 1/2 & 0 \end{array} \right] \end{array} = [1.5 \ 1 \ 0.5]$$

Stable solution. \swarrow row entries $\frac{1}{n}$, not 1

$$P^T \cdot A^* = P^T$$

Page Rank — Larry Page

A^* (henceforth A) is a stochastic matrix

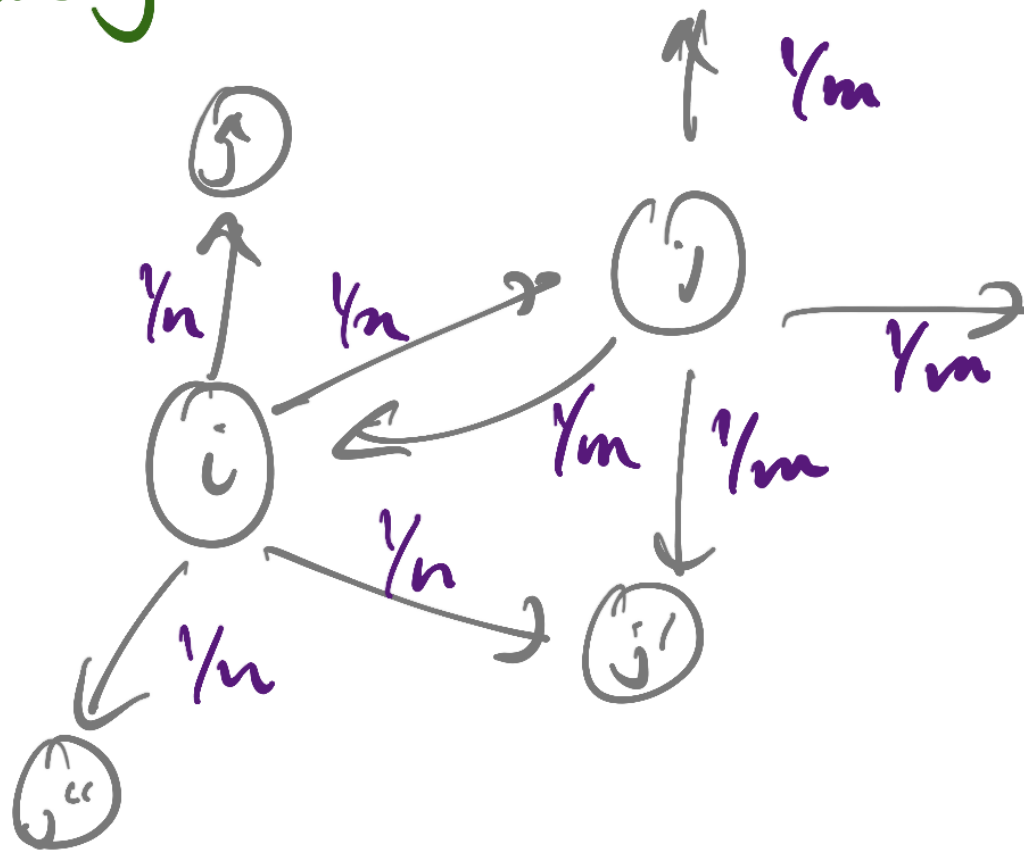
Every row adds up to 1

$$\forall i : \sum_j A[i,j] = 1$$

Interpret the entries as probabilities

"Random surfer" model

Probability $A[i,j]$, move from doc i
to doc j



Markov Chain

Finite collection of "states"

Transition probabilities between states

Start in document 1

$$[1 \ 0 \ 0] \overset{A}{\begin{bmatrix} 0 & 1/2 & 1/2 \\ 1 & 0 & 0 \\ 1/2 & 1/2 & 0 \end{bmatrix}} = [0 \ 1/2 \ 1/2]$$

$$\begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} = \begin{bmatrix} \frac{3}{4} & \frac{1}{4} & 0 \end{bmatrix}$$

\downarrow
 P

After k iterations,

$P[j]$ is probability of being in state j after k steps

In our example

$$[1 \ 0 \ 0] \rightarrow [0 \ \frac{1}{2} \ \frac{1}{2}] \rightarrow [\frac{3}{4} \ \frac{1}{4} \ 0]$$

↓

$$\underline{\underline{[\frac{7}{16} \ \frac{1}{8} \ \frac{1}{8}]}} \leftarrow [\frac{1}{4} \ \frac{3}{8} \ \frac{3}{8}]$$

are all these
numbers non-zero
after some point?

Markov chain is ergodic if there
is some t_0 s.t

$$(P A^t)[i] > 0$$

for all i , for all $t > t_0$, for all P

Not ergodic?

- ① Go from i to j_1, j_2, \dots, j_k , no path back
- ② Cycle $i \rightarrow j \rightarrow k \rightarrow i \rightarrow j \rightarrow k$

In an ergodic Markov chain

- there is a stationary distribution π s.t.

$$\pi A = \pi$$

- for any starting P

$$\lim_{t \rightarrow \infty} P \cdot A^t = \pi$$

Guarantee ergodicity

1. Irreducibility \Rightarrow strongly connected graph

Any i, j have paths in both directions

2. Aperiodic \Rightarrow "no cycles"

Any i, j , lengths of all paths $i \rightsquigarrow j$

gcd should be 1

Web graph need not satisfy 1, 2

Also dead ends - docs with no outgoing links

Fix this - "teleportation"

Allow a random jump anywhere!

$$T: A[i,j] = \frac{1}{N} \text{ everywhere,}$$

$N = \# \text{ of docs}$

Transition matrix M :

└ probability of teleportation

$$M = \alpha T + (1 - \alpha) A$$

Check that M is stochastic

By construction

M is strongly connected

M is aperiodic - paths of any length

No dead ends

Can solve

$P = MP \rightarrow P$ is Page Rank