

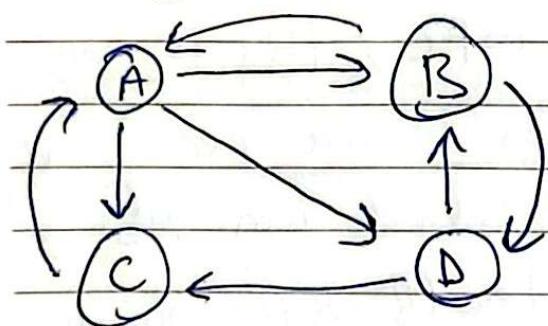


LINK ANALYSIS

I do have a relation which is a subset of 'Cartesian Product' btw 2 different sets.

* PAGERANK algorithms finding out the importance of webpage
one of the 2 founders of 'Google' Larry Page!

↳ They proposed specific algorithm that tried to find out among the huge number of pages which are typically returned when you query something on the web... We are not going to 2nd google page anymore
↳ Showing the first most reliable results! Google



→ This graph as description of a portion of the web (very small one)

→ each node of the graph refers to a webpage

→ Edges are Directed & refer hyperlinks between pages

→ Directed Edges because links in webpages are not bi-directional

→ The number of webpages estimated as $A \approx 10^9$

→ There's no universal definition of 'Big Data'. We deal with Big Data when we have an amount of data to be processed which is radically different from what we called small data approach. So, Big Data for Medicine field is radically different from Big Data in fields of Humanities or Computer Science etc. Especially, for what concerns Computer Science, we speak about Big Data, when we cannot reasonably use a single computer! (process or even store the data)

↳ So, if I have n total webpages & I'll deal with a matrix

$$M = [m_{ij}]_{n \times n}$$

Square matrix
'n' times 'n'

→ The generic entry ' M_{ij} ' will tell me something about the relation that exist btw 2 of these webpages ('i' and 'j')

If 'n' is the order of ' 10^9 ', M will have 10^{18} elements (Since it's a Square Matrix). Let's do the Math, and check how much space I need in order to fully store this matrix.

→ Before 'Google', on the search engines it was easy to cheat for showing their webpages at top since those relied on content. After 'Google', they changed 'order of importance' from content to 'structure'. Relies on the structure of the web, hyperlinks from one page to another one. Webpage is important if it is linked by important webpages! OFC, this is not acceptable definition! It leads to an infinite recursion!

→ webpage 'A' (A) nodesun gennie: bennet ian A'ya gelen koller baku (B) ve (C)'yi baku. Ama (B) ve (C)'nin important olup olmadiira bennet ian oyni process sonuz kez dekrh (infinite recursion olur!)

→ Basic Idea, 'Random Surfing' is just figured out that I have a population of huge number of individuals and I placed uniformly individuals in the various webpages and each individual will uniformly select at random one of the outgoing links that stand from the page where he/she sits! So, we have sort of 'Discrete-time' process, at each time, all individuals select uniformly at random one of the outgoing links and follow it. And we will iterate this process. (Stochastic Process). Also, we want to figure out if there is some 'convergence' in this process. Ideally, if I have a convergence of the probability distribution of individuals to webpages

that prob. dist will identify the importance I'm seeking!



UNIVERSITÀ DEGLI STUDI DI MILANO

Mathematically:

CONNECTION MATRIX

$$M = [m_{ij}]_{n \times n} \Rightarrow$$

A B C D Source & row-wise

'STOCHASTIC BY COLUMN'

(Ex)
 $\text{col}(A) = 0 + \frac{1}{3} + \frac{1}{3} + \frac{1}{3} = 1$
column 'A'

$$M = \begin{pmatrix} 0 & 1/2 & 1 & 0 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 1/2 & 0 & 0 \end{pmatrix}$$

A $m_{AA} = 0$ 0
B self-loop 1/2
C $m_{BA} = 1/2$ 0
D 0 0 0 0
B has out-degrees 2

Destination & column-wise

$$m_{ij} = P(j \rightarrow i)$$

$$v_j(t) = P(\text{sitting in } j \text{ at time } t)$$

Geliyoruz

$$v_j(t+1) = \sum_i P(i \rightarrow j \mid \text{sitting in } i \text{ at time } t).$$

P(sitting in i at time t)

which makes probability independent on time

$$= v_i(t)$$

→ We are considering the structure of the web as 'Immutable' that there will be no 'deletion' or 'addition' of new links so it becomes:

Now it becomes Matrix vector Product

$$v_j(t+1) = \sum_i P(i \rightarrow j) \cdot v_i(t) \Rightarrow v_j(t+1) = \sum_i m_{ji} \cdot v_i(t)$$

we can swap!

$$\Rightarrow \underline{v}(t+1) = M \underline{v}(t)$$

vector of probabilities

connection matrix 'M'

same vector but at previous time

POWER METHOD is Straightforward tool that allows to find out one eigenvector of a square matrix.

$$A(\underline{x}) = (\lambda)\underline{x} \rightarrow \text{Eigenvalue } (\lambda) \\ \text{Eigenvector } (\underline{x})$$

$$A = [a_{ij}] \xrightarrow{\substack{\text{Square matrix} \\ n \times n}} \\ \text{General element } (a_{ij}) \text{ of } A \text{ is } i^{\text{th}} \text{ row } j^{\text{th}} \text{ column}$$

Eigen Pair (λ, \underline{x})
(43. datatada)

$$(\lambda_1, \underline{x}_1), \dots, (\lambda_n, \underline{x}_n)$$

biggest eigenvalue is $\lambda_1 \geq \lambda_i \forall i$

$\underline{x}_1, \dots, \underline{x}_n$ is a basis

$\underline{v}_0 = \alpha_1 \underline{x}_1 + \alpha_2 \underline{x}_2 + \dots + \alpha_n \underline{x}_n$ (v_0 is weighted sum of the vector in the basis)

$$\underline{v}_1 = A \underline{v}_0 = A(\alpha_1 \underline{x}_1 + \alpha_2 \underline{x}_2 + \dots + \alpha_n \underline{x}_n)$$

$$= \alpha_1 A \underline{x}_1 + \alpha_2 A \underline{x}_2 + \dots + \alpha_n A \underline{x}_n$$

Can be written as

$$= \underbrace{\alpha_1 \lambda_1 \underline{x}_1 + \alpha_2 \lambda_2 \underline{x}_2 + \dots + \alpha_n \lambda_n \underline{x}_n}$$

$$\underline{v}_2 = A \underline{v}_1 = A(\alpha_1 \lambda_1 \underline{x}_1 + \dots + \alpha_n \lambda_n \underline{x}_n)$$

$$= \alpha_1 \lambda_1^2 \underline{x}_1 + \alpha_2 \lambda_2^2 \underline{x}_2 + \dots + \alpha_n \lambda_n^2 \underline{x}_n$$

$$\underline{v}_k = \alpha_1 \lambda_1^k \underline{x}_1 + \alpha_2 \lambda_2^k \underline{x}_2 + \dots + \alpha_n \lambda_n^k \underline{x}_n$$

$$\text{factor-out } \lambda_1^k \\ \underline{v}_k = \lambda_1^k \left(\alpha_1 \underline{x}_1 + \alpha_2 \left(\frac{\lambda_2}{\lambda_1} \right)^k + \dots + \alpha_n \left(\frac{\lambda_n}{\lambda_1} \right)^k \underline{x}_n \right)$$

$$\lim_{k \rightarrow \infty} \underline{v}_k = \lim_{k \rightarrow \infty} \lambda_1^k \alpha_1 \underline{x}_1$$

Result is that as k gets bigger as \underline{v}_k tends to have same direction of \underline{x}_1



UNIVERSITÀ DEGLI STUDI DI MILANO

which happens to be eigenvector associated to the first eigenvalue. (first eigenvector)

Thus, this shows 'Power Method' works! If I want to find out which is the direction associated to the first eigenvector of a matrix I can simply take the matrix, multiply it by generic vector (\underline{v}_0 in our case). Take the result and multiply it by 'A' and repeat it until some form of convergence takes places. The result is a vector direction identified first eigenvector!

$$\det(A) = \sum_i a_{ij} \underset{\text{MINOR}}{\overbrace{c_{ij}}} = \sum_j a_{ij} c_{ij}$$

minor: submatrix I obtained after I eliminated i^{th} row and j^{th} column

$$\text{Ex: } M = \begin{pmatrix} A & B & C & D \\ 0 & -1/2 & 1 & 0 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 1/2 & 0 & 0 \end{pmatrix} A \quad \text{let's say } C_{12} = \begin{pmatrix} 1/3 & 0 & 1/2 \\ 1/3 & 0 & 1/2 \\ 1/3 & 0 & 0 \end{pmatrix}$$

$$\det(A^T) = \sum_i a_{ij}^T c_{ij}^T = \sum_i a_{ji} c_{ji} \quad \xrightarrow{\text{Transpose makes them}} \quad \Downarrow \text{If we swap } i \leftrightarrow j$$

$$\det(A^T) = \sum_j a_{ij} c_{ij} = \det(A)$$

Facts : ① $\det(A^T) = \det(A)$

② A and A^T have same λ_s (eigenvalues)

$$\det(A - \lambda I) = 0 \xrightarrow{\text{should}} \det(A - \lambda I)^T = 0$$

↓
Unknown value of the equation
Identity Matrix (has same dimensions of A)

$$\downarrow \quad \det(A^T - \lambda I) = 0$$

∴ characteristic eqn of A and A^T are equal.

→ If A is row-stochastic matrix →

$$\sum_j a_{ij} = 1$$

(j) 'i' is fixed and 'j' is running

→ A row-wise stochastic is 'a' element

i^{th} row fixed kahr sadece o row'a koziki

gelen column'larda toplam ('j' is running) ve toplam 1'e esitdir.

$$A \cdot \underline{1} = \left[\sum_j a_{ij} (\underline{1})_j \right] \rightarrow j^{th} \text{ element of vector of } \underline{1}'s$$

$$\left(\begin{array}{l} \\ \end{array} \right) = \left[\sum_j a_{ij} \right]_n \quad \leftarrow \text{multiplying with } \underline{1} \text{ changes nothing, it is neutral!}$$

$$\left(\begin{array}{l} \\ \end{array} \right) A \cdot \underline{1} = [\underline{1}]_n = \underline{1} \rightarrow \text{vector of } \underline{1}'s$$

which means
 $\underline{1}$ is eigenvector
of 'A'

$$\left(\begin{array}{l} \\ \end{array} \right) A \cdot \underline{1} = \underbrace{1}_{\text{Scalar}} \cdot \underbrace{\underline{1}}_{\text{vector of } \underline{1}'s} \Rightarrow (1, \underline{1}) \text{ is eigen pair of 'A'}$$

We proved the 3rd fact

③ $(1, \underline{1})$ is eigenpair of row-wise stochastic

④ Transpose of row-wise stochastic matrix necessarily a column-wise stochastic matrix. So, 1 is eigenvalue of column-wise stochastic matrices.

→ One of the possible eigenvalues of the connection matrix (A) will be 1.

If A is row-wise stochastic.

by induction

base: $k=1$ obvious

Step: A^k is row-stochastic

then $\rightarrow A^{k+1}$ is row-stochastic matrix

$$A^{k+1} = A^k \cdot A$$

$$\left\{ \begin{array}{l} a_{ij}^{k+1} \\ \text{if } i \text{ row and } j \text{ th column of } A^{k+1} \\ \text{not } A \\ \text{generic element of } A^{k+1} \text{ not } A \end{array} \right\} a_{ij}^{k+1} = \sum_s a_{is}^k \cdot a_{sj}$$

$$\sum_j a_{ij}^{k+1} = \sum_j \sum_s a_{is}^k \cdot a_{sj}$$

$$= \sum_s a_{is}^k \sum_j a_{sj} = \sum_s a_{is}^k = 1$$

(5) If A is row-wise stochastic matrix also A^k (for any k)
is a row-wise stochastic matrix $\forall k$

\rightarrow The first eigenvalue of A column-wise stochastic is 1

by absurdity initial hypothesis

$\exists \lambda > 1$ is eigenvalue of A

$\rightarrow \lambda$ is eigenvalue of A^T

$$A^T v = \lambda v \rightarrow (A^T)^k v = \lambda^k v$$

$$\sum_j (A^T)^k_{ij} v_j = \lambda^k v_i$$

$$v_{\max} = \max_i v_i \rightarrow \sum_j (A^T)^k_{ij} v_{\max} > 6$$

set '6' as HIGH

$$\underbrace{\sum_j (A^T)^k_{ij}}_{\text{This sum}} > \frac{6}{v_{\max}}$$

This sum
is 1

$$1 > \frac{6}{v_{\max}}$$

Contradiction

Absurdity $\forall k$

Then hypothesis become $\lambda \leq 1$ (Method works) $\lambda > 1$ (works absurd)

the first eigenvalue is 1

$$\rightarrow M \underline{v}(0) = \underline{v}(1)$$

high enough for 'Power Method' works

column-wise
stoch.
matrix

$\vdots \quad \} \text{ after } t \text{ iteration}$

$\vdots \quad \underline{v}(t) = M \underline{v}(t-1)$

$\vdots \quad \vdots \quad \vdots$

$\underbrace{\underline{v}}_{\text{converges}} = \underbrace{M \underline{v}}_{\underline{v}}$

After a suitable number
of iterations
we won't see any remarkable
differences.

2nd part → The matrix will be always the same matrix whereas the vector changes each time. Thus we need an initial vector \underline{v}_0 .

1) $\underline{v}_0 = [1/n]_n$ where n : number of webpages, portion of web we are considering

2) $t=0$ { we set variable counts current time so, \underline{v}_0 is describing a particular kind of prob. dist overall the webpages we are considering. (Uniform dist)

3) $\underline{v}_{t+1} = M \underline{v}_t$

4) $t+1$ { we increased time index

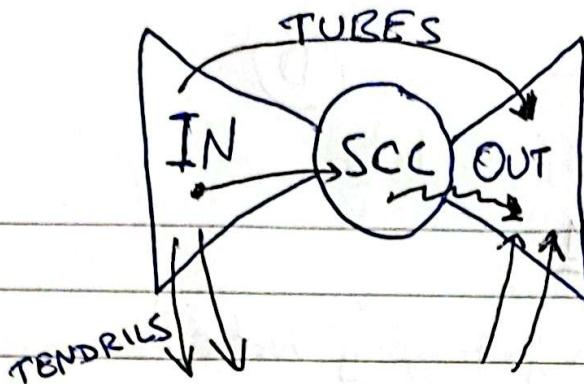
5) if $\|\underline{v}_{t+1} - \underline{v}_t\| > \epsilon$ go to Step 3, otherwise:

6) return \underline{v}_{t+1}

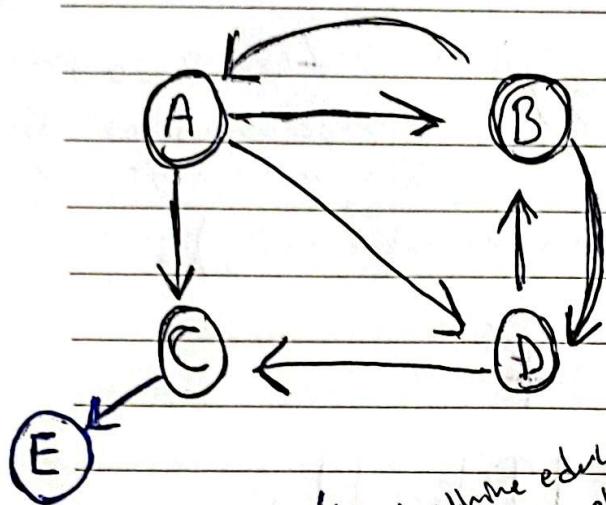


BO W-TIE representation of web

NCC



- Deadends E on the directed graph
- Spider traps

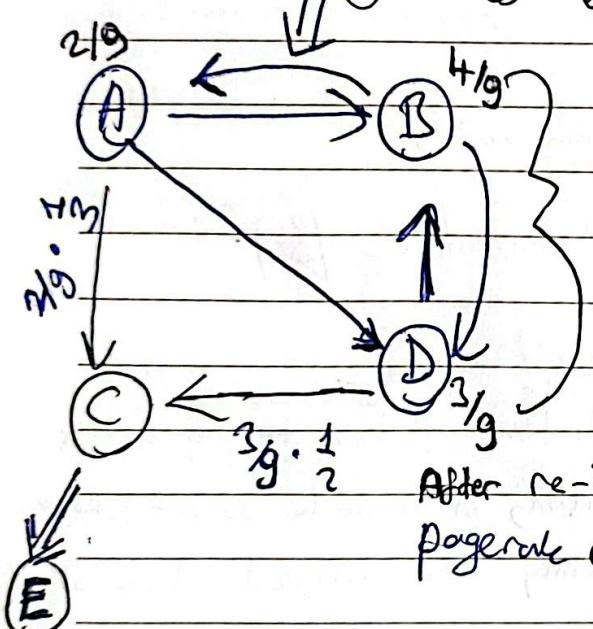


Deadends E in A ne begletsiget
B ve D'ye de gitmeyen

A	B	C	D	E	Source row
0	1/2	0	0	0	A
1/2	0	0	1/2	0	B
1/2	0	0	1/2	0	C
1/3	1/2	0	0	0	D
0	0	1	0	0	E

$= M$

E outdegree is zero (0) y contradiction to be a column-wise stack



Deadends'lerin ortaklari holding
ve ($\text{return } \frac{1}{f+1}$) algoritma mantikli bir
sayisal mekanizma ortak
verdigii page ranklerini özerince yazdirme

After re-introducing C and E to the graph:

$$\text{page rank of } \text{C} = \frac{2}{9} \cdot \frac{1}{3} + \frac{3}{9} \cdot \frac{1}{2} = \frac{13}{54}$$

$$\text{page rank of } \text{E} = \frac{13}{54} \text{ since only contributor is } \text{C} \text{ to } \text{E}$$

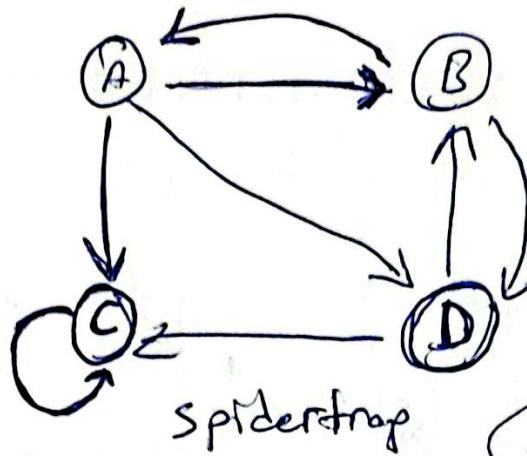
E

Problems PageRank'linin döplerini golden 1'i getiri C ve E'yi
of deadends eklemeler once bille 1 idig

Spider traps

Spider TRAP

TELEPORT / TAXATION



Spidertrap

$$v_{t+1} = \beta M v_t + (1-\beta) \frac{1}{n} \cdot \underbrace{\frac{1}{n}}_{v_0: \text{triggering Vector}}$$

$$\beta \in (0, 1)$$

GOOGLE USES
0.85

$\beta = P(\text{keeping on with random surfing})$

$$\beta \approx 0.8$$

efter ola sot yahı bir deger olursa 6 egeneli algoritma (bir önceki sayfa)
return $v_{t+1} \Rightarrow v_0$ gibiler
seye output verir. 3) ve
5) arasi loop girmez!

$M v_t \leftarrow P(\text{sitting in } i \text{ at time } t+1 \mid \text{keeping on with random surfing})$.

$\beta \leftarrow P(\text{keeping on with random surfing}) +$

$\frac{1}{n} \leftarrow P(\text{sitting in } i \text{ at time } t+1 \mid \text{Teleport})$.

$\frac{1}{n}$ \leftarrow $P(\text{Teleport})$
 ↓ not keeping random surf

$(1-\beta)$ (complimentary event of β)

Why do we have also Taxation name here?

⇒ Another Possible interpretation is that each time I have to pay a tax that is some portion of the individuals who are sitting in a node do not follow the random surfing process but they are automatically re-assigned democratically in a ideal democratic society, taxes are uniformly distributed over all the population?

Why this Teleport/Taxation algorithm avoid the SpiderTraps?

→ It will not totally/completely eliminate the problem but it will somehow relieve the problem!

How Do I derive the PageRank creation process



UNIVERSITÀ DEGLI STUDI DI MILANO

so that the result is biased towards

one of those topics?

→ Let's say: I have fixed one topic and also I know whether topic S or not this page says something about that topic.

I may build a vector let's call it \underline{S} that is boolean

$\underline{S} \rightarrow$ (all components will be 0 or 1 accordingly it has that topic on the webpage \underline{S}_1 otherwise \underline{S}_0)

$$S_i = \begin{cases} 1 & \text{if Page is about the topic fixed} \\ 0 & \text{otherwise} \end{cases}$$

general component vector \underline{S}

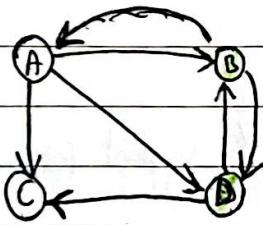
$$\underline{V}_{t+1} = \beta M \underline{V}_t + (1-\beta) \frac{1}{|S|} \cdot \underline{S}$$

dividing $|S|$ cardinality of S
by instead total number of webps
'n'

$|S| = \sum S_i$

The idea is that when Teleport occurs the individuals will not be uniformly assigned at random to all the available pages, rather they'll be uniformly assigned to one of the pages that speak of a particular subject. → So, If I know that utmost interest of a user is about cars, whenever that specific user gets teleported somewhere it'll be teleport to webpages speaks/related about cars!

Ex:



$$S = \{B, D\}$$

B and D speaks

about a specific topic
(topic in which my user interested in)

Topic-sensitive
PageRank

$$TSPR \quad .26 \quad .28 \quad .18 \quad .28$$

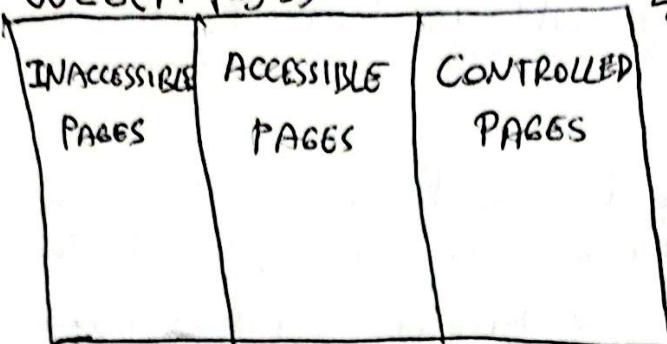
$$PR \quad .33 \quad .22 \quad .22 \quad .22$$

A user's
PR: depends (TSPR)

That's how Google does it

Web-Spammers & Artificially Increases the Page Rank of their webpages?

n : Total webpages Web(n pages)



We divided Venn-diagram in 3 subgroups of World Wide Web

$$PR(t) = y$$

$$PR(sp) = \beta \frac{y}{m} + (1-\beta) \frac{1}{n}$$

PR of t
evenly distributed
among m webpages (outgoing links of t)

β denotes Target Page

m support pages (SP) : denotes any of the support webpages

Once I deducted part pertaining to TAXATION so to obtain the part that is not deducted

will be evenly distributed to the all 'm' support pages. Again, because of Taxation each of 'm' support pages will receive some pageRank of the individuals that are randomly teleported to it. $(1-\beta)$ portion of the overall population is evenly distributed over all the 'n' webpages.

Let's focus again t : $PR(t) =$

3 components

$$\begin{aligned} & \text{from accessible pages: } \left(\beta \frac{y}{m} + (1-\beta) \frac{1}{n} \right) \cdot \beta \\ & \text{support page (for each sp)} \\ & \text{NEGLIGIBLE: } (1-\beta) \frac{1}{n} \end{aligned}$$

So:

$$y = x + m \cdot \beta \left(\beta \frac{y}{m} + (1-\beta) \frac{1}{n} \right)$$

$$y = x + \beta^2 y + \left(\beta (1-\beta) \right) \cdot m$$

$$y(1-\beta^2) = x + \beta (1-\beta) \cdot \frac{m}{n}$$

Suppose $\beta = 0.8$
 $PR(t) = (3.6)x + (0.18) \frac{m}{n}$

SPAM MADE ≈ 4 times better!