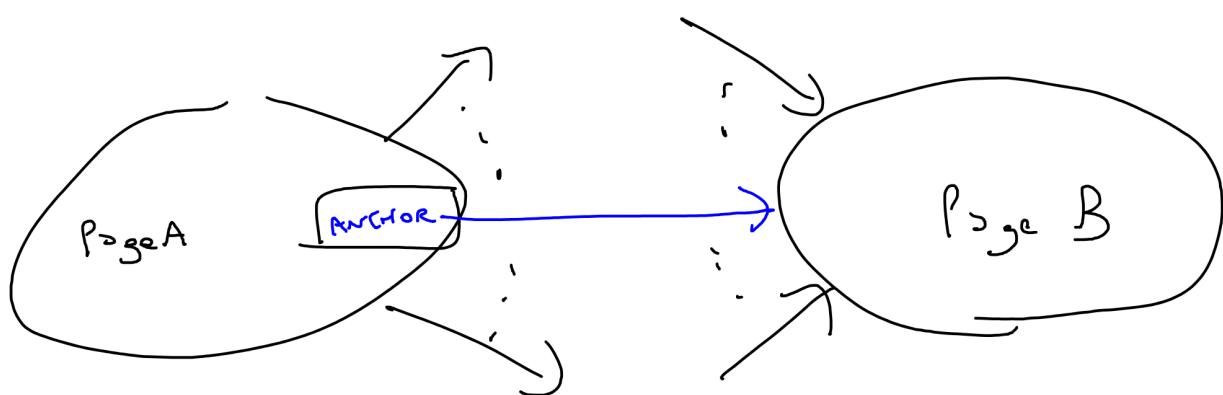


Link-based ranking

Second generation, used by modern web search engines, which ranks a document (usually, since we are talking about the web, we talk about web nodes) based on the hyperlinks that represent the webgraph.

The assumptions are that on hyperlink between two pages (web nodes) denotes a perceived relevance by the author (**quality signal**) and that the text in the anchor of the hyperlink describes the target page (**textual context**)



USS

When indexing a webpage, you should also include the anchor texts pointing to it, because sometimes the publishers don't provide an accurate description of themselves (marketing strategy?).

EXAMPLE

The IBM homepage doesn't contain the term "computer" anywhere, but the pages citing it may contain more useful informations.

OSS of USS

You should also add a weight representing the trust factor of the anchor page's website, because it might use wrong descriptions to associate the cited page with bad terms.

Under the assumptions shown above, a random walker that traverses the web graph would step over some nodes more often than others, based on the number of links pointing to them. The idea of web ranking functions such as PageRank is that the pages visited more often in this walk are more important.

How do we model this random walk on a graph?

MARKOV CHAINS

Given a graph

- Adjacency matrix A = represent the graph in memory
 - Transition probability matrix P = represent the probabilities of moving from a state (node) to another
- $$P_{ij} = P_r(j|i) = \text{probability of moving from } i \text{ to } j \in [0,1]$$

Note = the total transition probability from a node i to all the other nodes must be 1

$$\forall i \quad \sum_j P_{ij} = 1$$

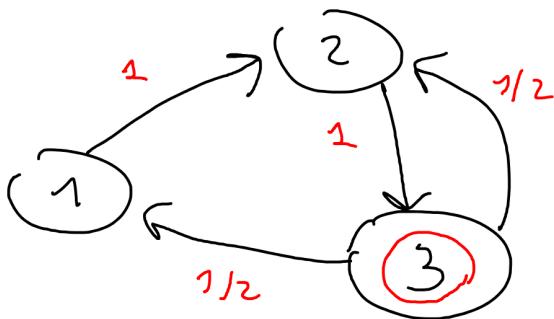
P is STOCHASTIC MATRIX

The probability distribution of the next state (node)
depends only on the current state (node) and not
on how the Markov chain (graph) arrived
at that current state (node)

- $X_t(i)$ = probability that the random surfer is at node i at time t

$$\begin{aligned} \cdot \underline{X_{t+1}(i)} &= \sum_j X_t(j) \cdot P_{ji} = x_t \cdot P = \\ &= x_{t-1} \cdot P \cdot P = \\ &= \dots = \\ &= x_0 \cdot P^{t+1} \end{aligned}$$

Example



I'm at node 3 at time t (x_t known),
compute X_{t+1}

$$X_{t+1} = X_t \cdot P = [0 \ 0 \ 1] \cdot \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1/2 & 1/2 & 0 \end{bmatrix} = \begin{bmatrix} 1/2 \\ 1/2 \\ 0 \end{bmatrix}$$

And I would be able to compute it even if
I was given only the starting position (x_0)

We talk of stationary probability distribution
when the process (the random walk) goes on for
a long time and the distribution does not change
any more

$$X_{t+1} = X_t \Rightarrow X_t^T P = 1 \cdot X_t$$

left eigenvector
of eigenvalue 1

Therefore, it exists a stationary state probability vector π , such that

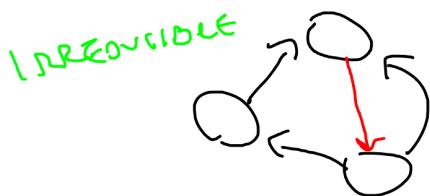
$$\lim_{t \rightarrow +\infty} x_0 \cdot P^t = \pi$$

Does a stationary distribution always exist?
is it unique?

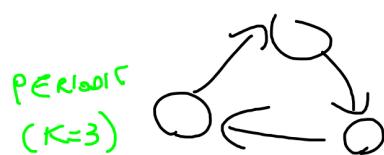
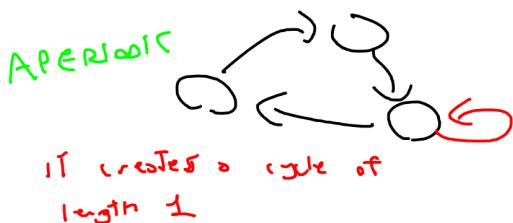
Yes, but only if the graph is "well-behaved",
which means that the Markov chain is
irreducible and aperiodic

- irreducible = A graph is irreducible if there is a path from every node to every other node

EXAMPLES = strongly connected directed graph / connected undirected graph



- aperiodic = A graph where the length of every cycle is not divisible by $K > 1$ (The Greatest Common Divisor is 1)



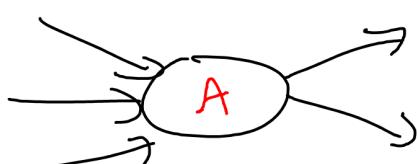
We can finally talk about how we can score web pages using links

1^o generation

Simply count the links to measure the popularity of a webpage. In particular there are two types of popularity:

- undirected popularity, where each page gets a score given by the number of in-links plus the number of out-links

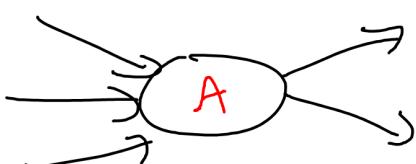
Example



$$r(A) = 5$$

- directed popularity, where the score is given only on the number of in-links

Example



$$r(A) = 3$$

IT'S VERY EASY TO TRICK THIS RANKING METHOD:
JUST CREATES PAGES FULL OF LINKS TO MAKE IT RANK HIGHER.

ALSO, IT DOESN'T EXPLOIT THE GRAPH STRUCTURE EXPLAINED ABOVE

2^o Generation: PageRank

To overcome the problem, every link is assigned with an importance value, which depends on how much is "valuable" the pointing webpage, which is determined by how many nodes link to it, and the pagerank will be computed using them.

As you can see, we can apply the random walk idea and model the web graph as a Markov chain

HOWEVER WE ARE INTERESTED IN THE STATIONARY PROBABILITY DISTRIBUTION, BECAUSE WE DON'T WANT OUR RANKING TO CHANGE

AT EVERY TIME STEP

But does the web guarantees it?

Actually NO, because the web has GUT nodes with no hyperlinks in it,

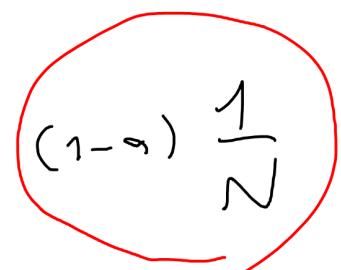
therefore it is not irreducible

BUT GOOGLE SOLVED THE PROBLEM INTRODUCING THE RANDOM JUMPS

This way we transformed the web graph in
→ strongly connected graph: From anywhere you
might go everywhere.

OSS: If you think about it, if you are "stuck"
in a page with no hyperlinks, you can just
type a new URL in the URL box and
get away with it

$$\pi(i) = \alpha \cdot \sum_{j \in B(i)} \frac{\pi(j)}{\#out(j)} + (1-\alpha) \frac{1}{N}$$



R RANDOM
JUMP

(OR TELEPORTATION)
STEP

- $\pi(i)$ = pagerank of i
- N = total number of pages in the webgraph
- $B(i)$ = set of pages linking to i
- $\#out(j)$ = number of out-links from j
- α = damping factor

NOTE IF WE IMAGINE THAT EACH NODE CONTAINS A TEXT CONTENT,
WE MIGHT WANT TO USE THE SIMILARITY $sim(T_i, T_j)$ INSTEAD OF $\#out(j)$
TO REPRESENT THE RELEVANCE OF node i

NOTE WE ASSUME RANDOM JUMP WITH UNIFORM
RANDOM PROBABILITY $p(\text{jump to node } i) = \frac{1}{\#nodes} = \frac{1}{N}$

principal eigenvector $\pi = [\alpha \pi^T + (1-\alpha) e e^T] \cdot \pi$
where $\pi_{i,j} = \begin{cases} 1/\#out(j) & \text{if } i \rightarrow j \\ 0 & \text{else} \end{cases}$

MANCHE ROBA PER CAPIRE E INTANZIARE CHE SUCCEDÈ

Page Rank in web searches

- PRE-PROCESSING PHASE = Given the Web graph, build the transition prob matrix P and then compute its principal eigenvector r :
 $r(:)$ is the pagerank of page i .
Since computing r is extremely costly, it is usually approximated as $\boxed{r = e^* P^t}$, where $t = 0, 1, \dots$ is the number of steps and $e = [\frac{1}{N} \ \frac{1}{N} \ \dots]$ is the starting prob.

The personalized pagerank let you decide the preferred jump pages by substituting e with a preference vector: this way you bias the jumps towards certain pages.

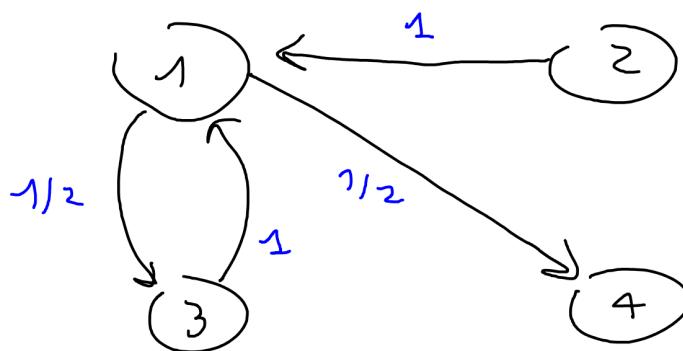
- QUERY PHASE = Retrieve the pages containing the query terms and rank them by their pagerank (the ranking is completely independent from the query%)

Exercise 1

Given a directed graph G with edges $\{(1,3), (3,1), (1,4), (2,1)\}$, simulate the execution of one step of PR algorithm, starting with $PR = 1$ (unnormalized)

every node and assuming a teleportation step that jumps only to node 3.

normalized



$$n(1) = \frac{1}{4}$$

$$n(2) = \frac{1}{4}$$

$$n(3) = \frac{1}{4}$$

$$n(4) = \frac{1}{4}$$

$$\varrho = \text{teleportation vector} = (0 \ 0 \ 1 \ 0)$$

$$\alpha = 1/2 \quad (\text{since the text doesn't say anything})$$

$$n(1) = \frac{1}{2} \left(\frac{1}{\frac{1}{2}} + \frac{1}{\frac{1}{2}} \right) + \frac{1}{2} \cdot 0 = \frac{1}{4}$$

$$n(2) = \frac{1}{2} \cdot 0 + \frac{1}{2} = 0$$

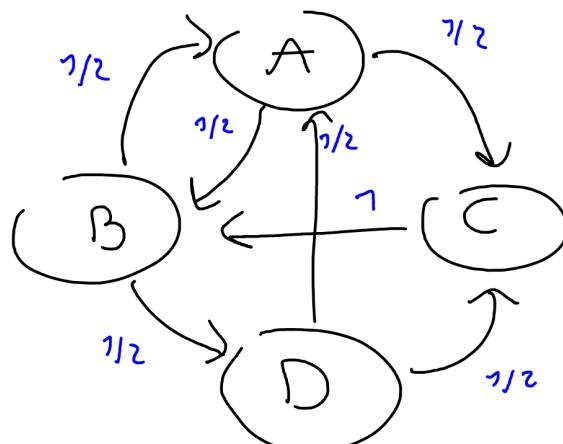
$$n(3) = \frac{1}{2} \left(\frac{1}{\frac{1}{2}} \right) + \frac{1}{2} \cdot 1 = \frac{1}{16} + \frac{1}{2} = \frac{9}{16}$$

$$n(4) = \frac{1}{2} \left(\frac{1}{\frac{1}{2}} \right) + \frac{1}{2} \cdot 0 = \frac{1}{16}$$

Exercise 2

Given the directed graph G of edges $\{(B, A), (A, C), (D, C), (A, B), (C, B), (D, A), (B, D)\}$; compute two steps of PageRank assuming uniform starting probabilities and $\alpha = 3/4$ (teleportation step with $1/4$); compute one step of personalized PR for the node A assuming starting probability $[1, 0, 0, 0]$ and $\alpha = 1/2$.

(1)



$$\pi(A) = 1/4$$

$$\pi(B) = 1/4$$

$$\pi(C) = 1/4$$

$$\pi(D) = 1/4$$

Step 1

$$\pi(A) = \frac{3}{4} \left(\frac{\frac{1}{4}}{2} + \frac{\frac{1}{4}}{2} \right) + \frac{1}{4} \cdot \frac{1}{4} = \frac{3}{4} \cdot \frac{1}{4} + \frac{1}{16} = \frac{1}{9}$$

$$\pi(B) = \frac{3}{4} \left(\frac{\frac{1}{4}}{2} + \frac{\frac{1}{4}}{1} \right) + \frac{1}{4} \cdot \frac{1}{4} = \frac{3}{4} \cdot \frac{3}{8} + \frac{1}{16} = \frac{11}{32}$$

$$\pi(C) = \frac{3}{4} \left(\frac{\frac{1}{4}}{2} + \frac{\frac{1}{4}}{2} \right) + \frac{1}{4} \cdot \frac{1}{4} = \frac{1}{4}$$

$$\pi(D) = \frac{3}{4} \left(\frac{\frac{1}{4}}{2} \right) + \frac{1}{4} \cdot \frac{1}{4} = \frac{3}{4} \cdot \frac{1}{8} + \frac{1}{16} = \frac{5}{32}$$

Step 2

$$r(A) = \frac{3}{4} \left(\frac{\frac{11}{32}}{2} + \frac{\frac{5}{32}}{2} \right) + \frac{1}{16} = \frac{3}{4} \cdot \frac{1}{4} + \frac{1}{16} = \frac{1}{4}$$

$$r(B) = \frac{3}{4} \left(\frac{\frac{1}{4}}{2} + \frac{\frac{1}{4}}{1} \right) + \frac{1}{16} = \frac{11}{32}$$

$$r(C) = \frac{3}{4} \left(\frac{\frac{1}{4}}{2} + \frac{\frac{5}{32}}{2} \right) + \frac{1}{16} = \frac{3}{4} \cdot \frac{13}{64} + \frac{1}{16} = \frac{55}{256}$$

$$r(D) = \frac{3}{4} \left(\frac{\frac{11}{32}}{2} \right) + \frac{1}{16} = \frac{3}{4} \cdot \frac{11}{64} + \frac{1}{16} = \frac{95}{256}$$

(2)

SAME GRAPH, BUT

$$r(A) = 1$$

$$r(B) = r(C) = r(D) = 0$$

Therefore the teleportation step can only
reach node A (at the start)

Step 1

$$r(A) = \frac{1}{2} \left(\frac{0}{2} + \frac{0}{2} \right) + \frac{1}{2} \cdot 1 = \frac{1}{2}$$

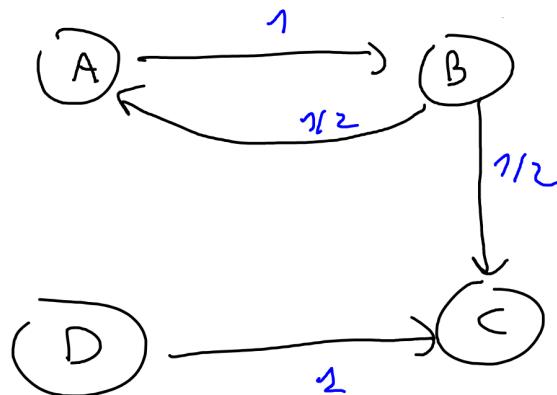
$$r(B) = \frac{1}{2} \left(\frac{1}{2} + \frac{0}{2} \right) + \frac{1}{2} \cdot 0 = \frac{1}{4}$$

$$r(C) = \frac{1}{2} \left(\frac{1}{2} + \frac{0}{2} \right) + \frac{1}{2} \cdot 0 = \frac{1}{4}$$

$$r(D) = \frac{1}{2} \left(\frac{0}{2} \right) + \frac{1}{2} \cdot 0 = 0$$

Exercise 3

Given the following graph



How can you estimate the similarity of B with all the other nodes, using

Personalized PageRank? Apply your logic.

For one step with $\alpha = \frac{1}{2}$.

Since we are not given a starting distribution, put it equal to the jump vect. $e = (0 \ 1 \ 0 \ 0)$

$$r(A) = \frac{1}{2} \left(\frac{1}{2} \right) + \frac{1}{2} \cdot 0 = \frac{1}{4}$$

$$r(B) = \frac{1}{2} \left(\frac{0}{1} \right) + \frac{1}{2} \cdot 1 = 1/2$$

$$r(C) = \frac{1}{2} \left(\frac{0}{1} + \frac{1}{2} \right) + \frac{1}{2} \cdot 0 = 1/4$$

$$r(D) = \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 0 = 0$$

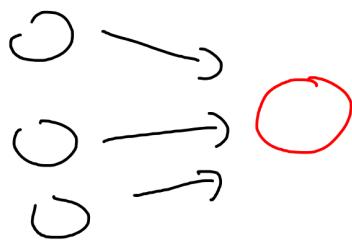
The most similar nodes are A and C
(To find out which one, at another step)

HITS = Hypertext Induced Topic Search

A different approach of scoring, where the score depends on the query.

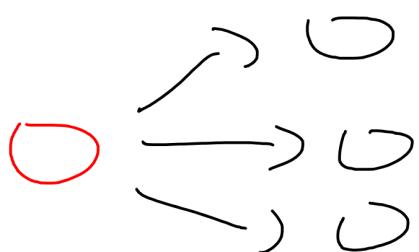
Actually the HITS score is made of two values:

- Authority score $a(i)$ = this value is high if page i has a good authority for a certain topic, which means that many good hubs for that topic point to him



$$a(v) = \sum_{v \rightarrow v} h(v)$$

- Hub score $h(i)$ = this value is high if page i is a good hub for a topic, which means that it points to many authoritative pages for that topic



$$h(v) = \sum_{v \rightarrow v} a(v)$$

Formally

$$\left. \begin{array}{l} \alpha = A^T h \\ h = A \alpha \end{array} \right\} \Rightarrow \begin{array}{l} \alpha = A^T A \alpha \\ h = A A^T h \end{array}$$

where A is the adjacency matrix

$\rightarrow h$ is an eigenvector of $A A^T$

$\rightarrow \alpha$ is an eigenvector of $A^T A$

Note

We could add some weights to the scores

to evaluate more some links

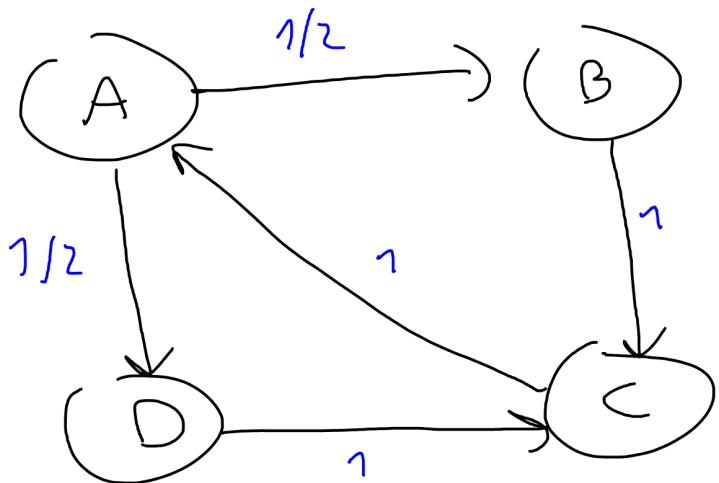
$$h(u) = \sum_{v \rightarrow u} w(u, v) \alpha(v)$$

$$\alpha(u) = \sum_{v \rightarrow u} w(u, v) h(v)$$

This ranking method is too costly at query time

Exercise

Given the graph



And the following initial scores (h, r)

$$A \left(?, ? \right), B \left(\frac{1}{4}, \frac{1}{2} \right), C \left(\frac{1}{4}, 0 \right), D \left(\frac{1}{9}, 0 \right)$$

Do the computation for node A

$$r(B) \quad r(D)$$
$$h(A) = \frac{1}{2} + 0 = \frac{1}{2}$$

$$r(A) = \frac{1}{4} \quad h(C)$$

