

Pagerank

... picking up from where the slides ran out...

What is pagerank?

→ It's just a measure of "importance" of a page, based purely on graph structure.
(It doesn't depend at all on the query, and so can be precomputed.)

Q: What possible measures of importance can you think of?

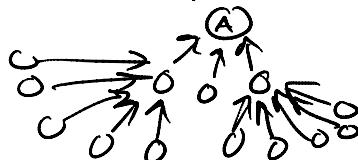
Strawman Answer: Degree

Q: Is the in-degree of a webpage a good measure of importance?

A: No!

→ Not all links are equal!

e.g. A is likely the most important in this graph



→ It's very easy to manipulate.

Q: What can we do to improve on degree?

Idea: Weight the edges by the importance

Details: Let $r_i(t)$ be the rank of page i at iteration t .

Set $r_i(0) = \frac{1}{n} \leftarrow \# \text{ of nodes}$

$$\therefore r_i(t+1) = \sum_{j \in N(i)} \frac{r_j(t)}{d_i^{\text{out}}} \quad \begin{array}{l} \text{out degree} \\ \text{of } i \\ \uparrow \text{pages that} \\ \text{link to } j \end{array}$$

In words, each node splits its rank evenly among its outgoing links, & each node's rank is the sum of what it gets from its neighbors.

We can write this nicely in matrix form:

$$r(t+1) = r(t)P = r(t) \begin{bmatrix} p_{ij} \end{bmatrix}$$

$$\text{where } p_{ij} = \begin{cases} \frac{1}{d_i} & \text{if } i \rightarrow j \\ 0 & \text{else} \end{cases}$$

(this is just a variation of
the adjacency matrix of the graph)

Note : $\sum_j r_j(t) = 1 \quad \forall t.$

Pagerank Algorithm (take 1):

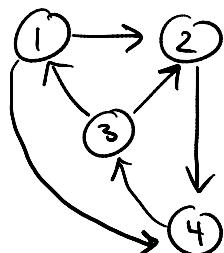
- 1) Assign each node an initial pagerank $\frac{1}{n}$.
- 2) Repeat $r(t+1) = r(t)P$ until it converges

* Think of a "random web surfer" starting at each node with probability $\frac{1}{n}$ and then following each link out of each node with equal probability. Then $r_i(t)$ is the

probability you're at node j at time t .

example:

Consider the graph



$$P = \begin{pmatrix} 0 & k_2 & 0 & k_2 \\ 0 & 0 & 0 & 1 \\ k_2 & k_2 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

$$r(0) = \left(\frac{1}{4}, \frac{k_2}{4}, \frac{k_2}{4}, \frac{1}{4}\right)$$

$$r(1) = r(0)P = \left(\frac{1}{8}, \frac{1}{4}, \frac{1}{4}, \frac{3}{8}\right)$$

$$r(2) = r(1)P = \left(\frac{1}{8}, \frac{3}{16}, \frac{3}{8}, \frac{5}{16}\right)$$

$$r(3) = r(2)P = \left(\frac{3}{16}, \frac{1}{4}, \frac{5}{16}, \frac{1}{4}\right)$$

$$r(4) = r(3)P = \left(\frac{5}{32}, \frac{1}{4}, \frac{1}{4}, \frac{11}{32}\right)$$

⋮

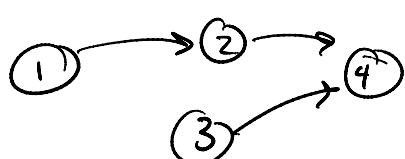
Two questions we want to answer about this:

1) Does this always converge?

2) Does it converge to what we want?

To get a feel, let's look at some more examples:

ex:



$$P = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad r(0) = \left(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}\right)$$

$$r(1) = r(0)P$$

$$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad r(0) = r(0) P$$

$$r(1) = (0, \frac{1}{4}, 0, \frac{1}{2})$$

$$r(2) = (0, 0, 0, 1)$$

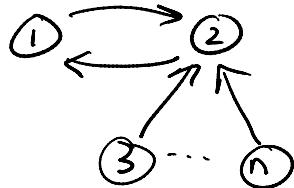
$$r(3) = (0, 0, 0, 1)$$

So, we converge...

But, is this what we want?

No!
we would answer
node 4 for every query!

ex:



$$P = \begin{pmatrix} 0 & 1 & 0 & \cdots \\ 1 & 0 & 0 & \cdots \\ 0 & 1 & 0 & \cdots \\ 0 & 0 & 1 & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

$$r(0) = (\frac{1}{n}, \dots, \frac{1}{n})$$

$$r(1) = r(0)P = (\frac{1}{n}, \frac{n-1}{n}, 0, \dots)$$

$$r(2) = r(1)P = (\frac{n-1}{n}, \frac{1}{n}, 0, \dots)$$

$$r(3) = (\frac{1}{n}, \frac{n-1}{n}, 0, \dots)$$

$$r(4) = (\frac{n-1}{n}, \frac{1}{n}, 0, \dots)$$

So it oscillates and never converges
and the difference between ranks goes
from ≈ 0 to ≈ 1 between iterations!

\Rightarrow So there are certainly some problems
with our initial version of pagerank!

1) periodicity in the links caused
trouble in the second example.

2) Having weakly connected nodes
caused problems in both examples.

(if we had been disconnected of
Google.com and Amazon.com have

(course more problems would have shown up!)

Thm. If the graph is finite, strongly connected and aperiodic, then $\lim_{t \rightarrow \infty} r(t) = \pi$ & π is

the unique solution to $\pi = \pi P$

"Stationary" rank vector

def G is aperiodic iff \forall vertices v the GCD of the path lengths from v to itself is 1

Q: Why is π called the stationary vector?

A: Because if we were to start with $r(0) = \pi$ then the ranks would never change.

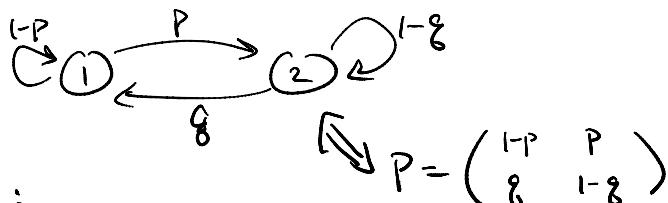
Q: Does anyone recognize what π is?

A1: left eigenvector of P w/ eigenvalue 1.

A2: Stationary distribution of the Markov chain defined by P . (Take CS47 or ACM216 to learn more about these)

Aside: 1 thing to notice is that this theorem gives us a second way to calculate pagerank: using $\pi = \pi P$ & $\sum_i \pi_i = 1$.

Ex: Consider



So $\pi = \pi P$:

$$\pi_1 = \pi_1(1-p) + \pi_2 g \Rightarrow p\pi_1 = g\pi_2$$

$$\pi_2 = \pi_1 p + \pi_2(1-g) \Rightarrow p\pi_1 = g\pi_2$$

$$\pi_1 + \pi_2 = 1$$

$$\pi_1 + \frac{P}{Q} \pi_1 = 1$$

$$\pi_1 = \frac{1}{1 + \frac{P}{Q}} = \frac{Q}{P+Q}$$

$$\& \pi_2 = \frac{P}{P+Q}$$

....back to pagerank...

★ Big Question Is the web graph finite, aperiodic, & strongly connected?

A: Of course not

So how do we adjust P so that it is aperiodic & strongly connected?



Add an edge from each node to every other node.

... make this edge have small weight compared to the real edges.

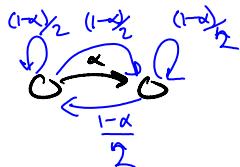
Specifically:

instead of P , use

$$G = \alpha P + (1-\alpha) \begin{bmatrix} \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

So $\textcirclearrowleft \rightarrow \textcirclearrowright$

Becomes



Q: what is an interpretation for this

Change!

A: Random web surfer follows links
for a while, but w/prob $(1-\alpha)$ starts
fresh...

so this can be viewed as
more "realistic" too!

Great! We now have an algorithm that we
know will provide convergence to something
reasonable.

‡ we have two ways to calculate
the ranks

(1) Iterate $r(t) = r(t-1) G$

(2) Solve $\pi = \pi G$
‡ $\sum \pi_i = 1$

(set of n equations with n unknowns)

Q: Which does google use?

A: (1) ...

Because of network structure

(small world & heavy tailed degree)

Convergence of the iterations happens
quite quickly. ≈ 100 iterations

whereas (2) requires inverting a very
large matrix...

PageRank Algorithm (take 2):

- 1) Assign each node an initial page rank $\frac{1}{n}$.
- 2) Repeat $r(t) = r(t-1) G$ for ≈ 100 steps

Some other issues:

- 1) Q: What should α be?

A: $\alpha \approx .85$ turns out to be good for both computation time & results quality.

Turns out to get w/in T of π
you need $\frac{\log T}{\log \alpha}$ iterations.

$\alpha \rightarrow 1$ lots of iterations
 $\alpha \rightarrow 0$ few iterations, but ignore web structure.

- 2) we made the matrix very dense by adding lots of links
→ This seems to make computation much harder!

* This can often be avoided, i.e. we can still compute with sparse matrices.

(Possible HW question!)

Now that we understand page rank ...

What are its problems?

- ① Takes lots of computation!

Right now it takes days to calculate
it needs to be updated weekly or faster!.

So → it cannot be personalized much to

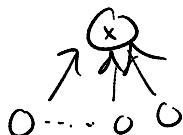
- type of user, or
- topic

... which both would lead to
big improvements in quality.

② Google Bombs

Page rank can be gamed by exploiting link
structure!

"Link farms" → to help page x, make lots
of other pages that will link to it.



even though these pages aren't linked
to x guarantees them some pagerank
which they can transfer to x.

Q: How can google fight this?

A: Change G so that when the "random
surfer" restarts it always jumps to "safe" pages.
(e.g. google, com, etc.)

(e.g. google, com, etc.)

$$\text{i.e. } G = \alpha P + (1-\alpha) \begin{bmatrix} \frac{1}{n_{safe}} & 0 & \dots & \frac{1}{n_{safe}} & \dots \\ \ddots & \ddots & \ddots & \ddots & \ddots \end{bmatrix}$$

... But you can't do this completely otherwise new pages that aren't linked to can't get positive page rank...

Last thought:

In real life G changes, but only by a small amount.

Q: Can we take advantage of this?

A: It seems so ... this should make the iterations converge faster if we start from the old pageranks

→ This seems natural, but unfortunately it doesn't usually work

Now that we're finished w/ Page rank, let's go back to the big picture

... go back to ppt