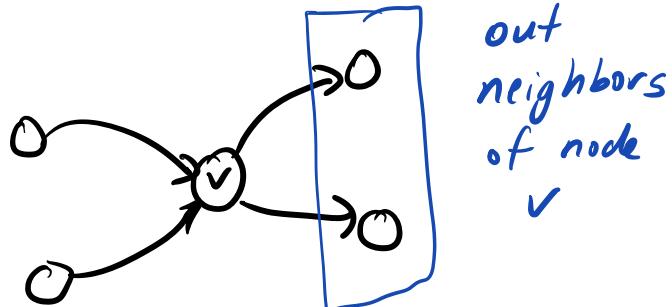


Random Walks and Centrality Scores in Graphs

Definition: Given a directed graph $G = (V, E)$ and a starting node $s \in V$, the standard random walk of length K is a sequence of random nodes

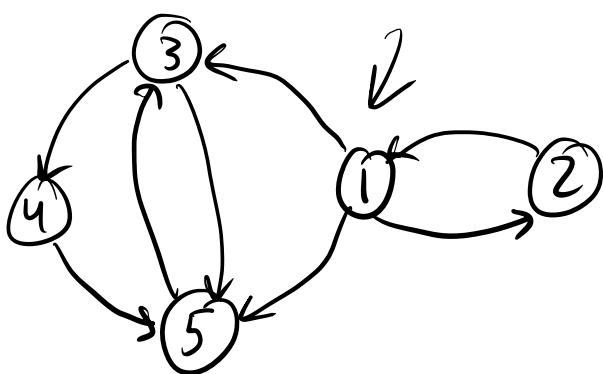
$$S = v_0, v_1, v_2, \dots, v_K$$

such that node v_i is chosen uniformly at random from among the out neighbors of v_{i-1} .



We will use d_v to denote the out-degree of node v (# of out neighbors).

Math behind random walks



$$d = d_{\text{out}} = \begin{bmatrix} 3 \\ 1 \\ 2 \\ 1 \\ 1 \end{bmatrix}$$

D is the degree matrix, $D_{ii} = d_i$

$$A = \begin{array}{ccccc|c} & 1 & 2 & 3 & 4 & 5 \\ \hline 0 & 1 & 1 & 0 & 1 & \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & \\ 0 & 0 & 0 & 0 & 1 & \\ 0 & 0 & 1 & 0 & 0 & \end{array} \quad D = \begin{bmatrix} 3 & & & & \\ & 1 & & & \\ & & 2 & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix}$$

$$e_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$A^T e_1$ gives me the set of nodes I can travel to if I start at node 1.

$$P = A^T D^{-1} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{3} & & & & \\ & 1 & & & \\ & & \frac{1}{2} & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix}$$

=  $\begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ \frac{1}{3} & 0 & 0 & 0 & 0 \\ \frac{1}{3} & 0 & 0 & 0 & 1 \\ 0 & 0 & \frac{1}{2} & 0 & 0 \\ \frac{1}{3} & 0 & \frac{1}{2} & 1 & 0 \end{bmatrix} \quad P_{31}$

So P_{ij} gives me the probability of taking a step to node i , if I start at node j . P is the "random walk matrix".

$P^K x_s$ is the probability distribution for where I can be after K steps.

$$\text{where } x_s(i) = \begin{cases} 1 & \text{if } i=s \\ 0 & \text{otherwise} \end{cases}$$

If $e^T x = 1$, then $y = Px$ satisfies

$$y^T e = 1.$$

A matrix with columns that sum to 1
(and has nonzero entries) is called a
column stochastic matrix.

Observation: If $A = A^T$ then we

have

$$\hat{d} = \frac{d}{e^T d}$$

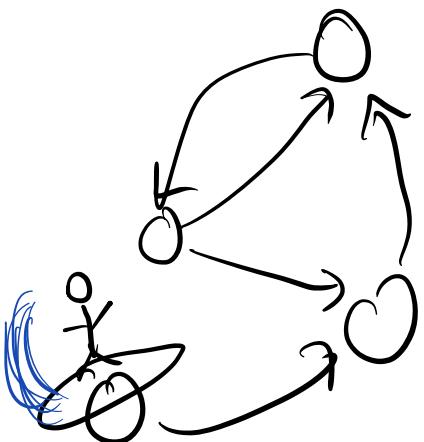
$$\begin{aligned} P \hat{d} &= A^T D^{-1} \hat{d} \\ &= A \underline{D^{-1} \hat{d}} = \underline{\hat{d}} \end{aligned}$$

So \hat{d} is the "steady state" vector
for the graph.

PageRank

Notes follow paper "PageRank Beyond the Web" D. Bleich
(SIAM Review)

Model : random surfer



At each step

- with probability $\alpha \in (0, 1)$ surfer follows a uniform random out link
- with probability $(1 - \alpha)$ teleport to a different node based on some distribution given by a vector $v \geq 0$ satisfying $v^T e = 1$

$v(i)$ = probability I teleport to node i .

Simple case $v = \frac{e}{n}$

Page Rank centrality (informal)

A webpage (node) is more "important" or "central" if the surfer spends a higher proportion of time there.

I.e. we want to compute the average amount of time the surfer spends at node i .

Def: A matrix P is column stochastic if $P \geq 0$ and $e^T P = e^T$.

If is sub-stochastic if

$$e^T P \leq e^T$$

(columns sum to ≤ 1).

Page Rank centrality (formal)

Let $G = (V, E)$ be a graph and let v be a column stochastic vector defining starting probabilities.

e.g. $v = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\underline{v} = \begin{bmatrix} \frac{1}{5} \\ \frac{1}{5} \\ \frac{1}{5} \\ \frac{1}{5} \\ \frac{1}{5} \end{bmatrix} = \frac{e}{n}$ ($n=5$)

with probability α we transition based on the standard random walk matrix

$$\boxed{P} = A^T D^{-1}$$

with probability $1-\alpha$ the surfer teleports to a node based on \mathbf{v} , i.e. teleportation matrix is

$$\underline{\mathbf{v}}\mathbf{e}^T = \begin{bmatrix} 1 & 2 & 3 \\ v & v & v & \cdots & v \end{bmatrix}$$

So the PageRank transition matrix is

$$\underline{M = \alpha P + (1-\alpha)\underline{\mathbf{v}}\mathbf{e}^T}$$

Let $x^{(0)} = \mathbf{v}$ be the starting vector, then the probability distribution at step $i+1$ is given by

$$\begin{aligned} \underline{x}^{(i+1)} &= \underline{M} \underline{x}^{(i)} \\ &= \alpha P \underline{x}^{(i)} + (1-\alpha) \underline{\mathbf{v}} \mathbf{e}^T \underline{x}^{(i)} \end{aligned}$$

$$\underline{x} = \alpha P \underline{x}^{(i)} + (1-\alpha) v$$

If this procedure converges, that means $\underline{x}^{(K)} \rightarrow \underline{x}$

$$\underline{x} = \alpha P \underline{x} + (1-\alpha) v$$

$$\Rightarrow \underline{x} - \alpha P \underline{x} = (1-\alpha) v$$

$$\Rightarrow \boxed{(\underline{I} - \alpha P) \underline{x} = (1-\alpha) v}$$

So the steady state vector \underline{x} is the solution to a linear system.

This is also the solution to an eigenvector problem

$$\boxed{\underline{x} = M \underline{x}}.$$

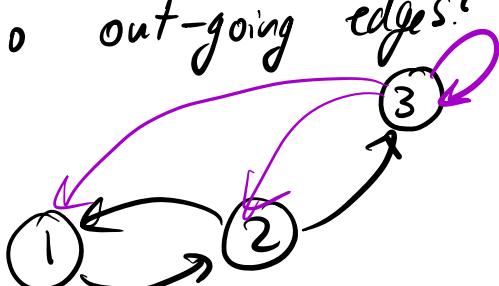
Definition: Let P be a column stochastic matrix. Assume $v \in \mathbb{R}^n$ satisfies $v^T e = 1$ and $v \geq 0$, and $\alpha \in (0, 1)$. The PageRank linear system is given by

$$(I - \alpha P)x = (1 - \alpha)v \quad (*)$$

where x is called the PageRank vector.

Fact: As long as $\alpha \in (0, 1)$ the solution x to $(*)$ exists, is unique and satisfies $x \geq 0$ and $e^T x = 1$.

Issue: What do we do if node 1 has no out-going edges?



$$D = \begin{bmatrix} 1 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\underline{P} = "A^T D^{-1}"$$

$$D^{-1} = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & ? \end{bmatrix}$$

Ideas: If I get to a dangling node:

- teleport based on v .
 - add self-loop
 - link it to every other node
(and add self loop?)
-

The Pseudo-PageRank System

Definition: Let \hat{P} be a column substochastic matrix ($e^T \hat{P} \leq e^T$, $\hat{P} \geq 0$)

Let f be a nonnegative vector and $\alpha \in (0, 1)$. The pseudo-PageRank problem is to solve the linear system

$$(I - \alpha \hat{P}) y = f \quad (\star)$$

and y is called the Pseudo PageRank vector.

Theorem: Let y be the solution to (\star)

and define $v = \frac{f}{e^T f}$ and

$$c = e - \hat{P}^T e \geq 0 \quad \text{and} \quad x = \frac{y}{e^T y}.$$

Then x is the PageRank vector for the system

$$(I - \alpha P)x = (1 - \alpha)v$$

where $P = \hat{P} + \nu c^T$.

Proof: Check that v is a stochastic vector. ✓
 Check that P is a column stochastic matrix:

$$\begin{aligned} e^T P &= e^T (\hat{P} + \nu c^T) \\ &= \underline{e^T \hat{P}} + \cancel{\nu c^T} = \boxed{e^T} \end{aligned}$$

because $c = e - \hat{P}^T e$

$$\begin{aligned} \Rightarrow e &= c + \hat{P}^T e \\ \Rightarrow \underline{e^T} &= \underline{e^T \hat{P} + c^T} \end{aligned}$$

Note from $(I - \alpha P)x = (1-\alpha)v$ that

$$\begin{aligned} x &= \alpha \underline{Px} + (1-\alpha)v \\ &= \alpha \hat{P}x + \alpha \underline{\nu c^T x} + (1-\alpha)v \\ &= \alpha \hat{P}x + [\alpha c^T x + (1-\alpha)]v \\ &= \alpha \hat{P}x + \underline{[\alpha c^T x + (1-\alpha)]f} \end{aligned}$$

$e^T f$

call this γ

$$= \alpha \hat{P}x + \gamma f$$

$$\Rightarrow (I - \alpha \hat{P})x = \gamma f$$

$$\Rightarrow (I - \alpha \hat{P}) \begin{bmatrix} x \\ \gamma \end{bmatrix} = f$$

\downarrow

$$\Rightarrow x = \gamma y$$

and since $x^T e = 1$

$$\Rightarrow \gamma y^T e = 1 \Rightarrow \gamma = \frac{1}{y^T e}$$

$$\Rightarrow x = \frac{y}{y^T e}.$$



Message: Every Pseudo-PageRank system corresponds to a standard PageRank system.

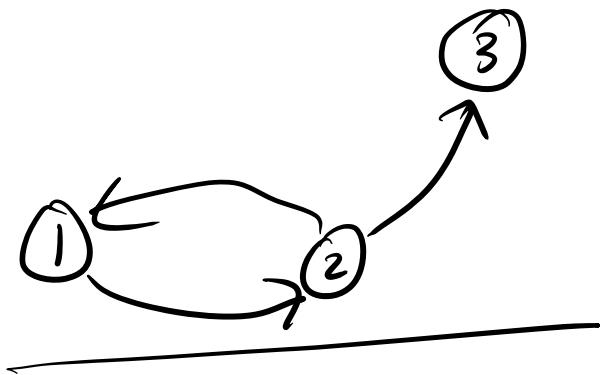
Let's fix the dangling node issue

Define D^+ to be the diagonal matrix

$$D_{ii}^+ = \begin{cases} \frac{1}{d_i} & \text{if } d_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

To define PageRank we use the matrix

$$\hat{P} = A^T D^+ \quad e^T \hat{P} \leq e^T$$



$$D^+ = \begin{bmatrix} 1 & & \\ & \frac{1}{2} & \\ & & 0 \end{bmatrix}$$

$$\hat{P} = \begin{bmatrix} 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 \end{bmatrix}$$

We define a correction vector by

$$c = e - \hat{P}^T e$$

in the example:

$$c = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 & 1 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$c[i] = \begin{cases} 1 & \text{if } i \text{ is dangling} \\ 0 & \text{otherwise} \end{cases}$$

Options for solving dangling node issue

- ① Solve pseudo-PR system

$$(I - \alpha \hat{P}) y = f \quad \text{where}$$

$f = v$. By our theorem, this will correspond to the Page Rank system

$$(I - \alpha P) x = (1 - \alpha)v$$

where $P = \hat{P} + vc^T$

e.g. $vc^T = v[0 \ 0 \ 1]$

$$= \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

αP

This is called "strongly-preferential"

Page Rank

② "weakly preferential"

$$P = \hat{P} + c u^T \quad \text{where } u = \frac{e}{n}$$

③ "sink preferential"

$$P = \hat{P} + \text{diag}(c)$$

$$P = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Personalized PageRank System

Define $v_R = \begin{cases} \frac{1}{|R|} & \text{if } i \in R \\ 0 & \text{otherwise} \end{cases}$

where $R \subseteq V$ is a subset of nodes in the graph. Then

$$(I - \alpha P)x = (1-\alpha)v_R$$

is a Personalized PageRank system and x is the PPR vector.

often $R = \{v\}$ for one node $v \in V$.