

Mathematics of Image and Data Analysis
Math 5467

PageRank

Instructor: Jeff Calder
Email: jcalder@umn.edu

<http://www-users.math.umn.edu/~jwcalder/5467>

Last time

- Spectral Clustering

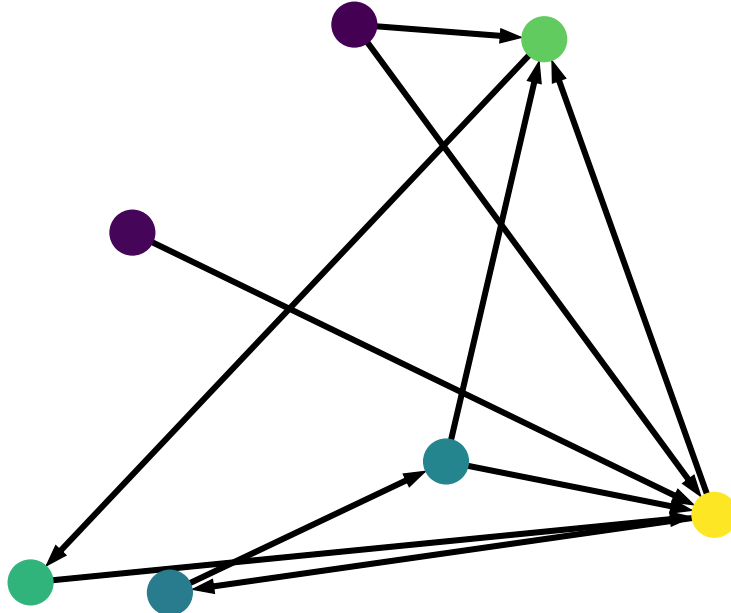
Today

- PageRank

PageRank

The PageRank algorithm ranks websites based on the link structure of the internet. It was used to sort Google search results until 2006, and has been used in

- Biology (GeneRank), chemistry, ecology, neuroscience, physics, sports, and computer systems...



PageRank

Main Idea: Take a random walk on the internet for T steps.

$$\text{Rank of site } i = \lim_{T \rightarrow \infty} \frac{1}{T} (\text{Number of times site } i \text{ is visited}).$$

Problem: Random walks can get stuck in disconnected components of the internet, and may never visit a given site i .

Solution: Every so often, the random walker teleports to a random site on the internet. The walker is called a **random surfer**.

Code demo

Mathematics of PageRank

To describe PageRank mathematically, we start with an adjacency matrix W

$$W(i, j) = \begin{cases} 1, & \text{if site } i \text{ links to site } j \\ 0, & \text{otherwise.} \end{cases}$$

We also have a probability transition matrix P for the random walk:

$$P(i, j) = \text{Probability of stepping from } j \text{ to } i.$$

Both P and W are $n \times n$ matrices, n = number of webpages.

Mathematics of PageRank

Clicking on a link at random from webpage j leads to the transition probabilities

$$P(i, j) = \frac{W(j, i)}{\underbrace{\sum_{k=1}^n W(j, k)}_{D(j, j)}}$$

$$D = \begin{pmatrix} D(1,1) & & 0 \\ & \ddots & \\ 0 & & D(n,n) \end{pmatrix}$$

Exercise 1. Show that $P = W^T D^{-1}$, where D is the diagonal matrix with diagonal entries $D(i, i) = \sum_{j=1}^n W(i, j)$. △

Random surfer

Let $\alpha \in [0, 1)$ be the random walk probability, and let $v \in \mathbb{R}^n$ be the teleportation probability distribution. That is, $v(i) \geq 0$ for all i , and $\sum_i v(i) = 1$.

Random surfer dynamics: When at website j , the random surfer chooses the next site as follows:

1. With probability α the surfer clicks an outgoing link at random, that is, the surfer navigates to website i with probability $P(i, j)$.
2. With probability $1 - \alpha$ the surfer teleports to website i with probability $v(i)$.

Teleportation

Teleportation distribution: Common choices are

- $v(i) = 1/n$ for all i (jump to a site uniformly at random).
- (Localized PageRank) $v(i) = \delta_{ij}$ (always jump back to site j).

Localized PageRank ranks all sites based on their similarity to site j .

The PageRank vector

For $k \geq 0$ define

$x_k(i)$ = Probability that the random surfer is at page i on step k .

Definition 2. The PageRank vector x is

$$x(i) = \lim_{k \rightarrow \infty} x_k(i),$$

provided the limit exists.

Transition probabilities

To see how x_k transitions to x_{k+1} requires some probability. We condition on the location of the surfer at step k , and on the outcome of the coin flip, to obtain

$$x_{k+1}(i) = (1 - \alpha)v(i) + \alpha \sum_{j=1}^n P(i, j)x_k(j).$$

We can write this in matrix/vector form as

$$(1) \quad x_{k+1} = (1 - \alpha)v + \alpha Px_k.$$

If x_k converges to a vector x as $k \rightarrow \infty$, then x should satisfy

$$x = (1 - \alpha)v + \alpha Px.$$

Question: Does x_k converge as $k \rightarrow \infty$, and if so, how quickly does it converge?

Analysis of PageRank

We consider the PageRank equation

$$(2) \quad x = (1 - \alpha)v + \alpha Px.$$

Lemma 3. *Let $v \in \mathbb{R}^n$ and $0 \leq \alpha < 1$. Then there is a unique vector $x \in \mathbb{R}^n$ solving the PageRank equation (2). Furthermore, the following hold.*

(i) *We have $\sum_{i=1}^n x(i) = \sum_{i=1}^n v(i)$.*

(ii) *If $v(i) \geq 0$ for all i , then $x(i) \geq 0$ for all i .*

The ℓ_1 -norm

$$\|x\|_2 = \sqrt{x^{(1)^2} + \dots + x^{(n)^2}}$$

It will be more convenient to work in the ℓ_1 -norm $\|\cdot\|_1$ defined by

$$\|x\|_1 = \sum_{i=1}^n |x^{(i)}|.$$

In the ℓ_1 -norm, the transition matrix P is non-expansive.

Proposition 4. We have $\|Px\|_1 \leq \|x\|_1$.

Proof:

$$\begin{aligned} \|Px\|_1 &= \sum_{i=1}^n |(Px)_i| \\ &= \sum_{i=1}^n \left| \sum_{j=1}^n P(i,j) x^{(j)} \right| \\ &\leq \sum_{i=1}^n \sum_{j=1}^n |P(i,j)| |x^{(j)}| \end{aligned}$$

$$P(i,j) = \frac{w(j,i)}{\sum_{k=1}^n w(j,k)}$$

$$= \sum_{j=1}^n |x(j)| \underbrace{\left(\sum_{i=1}^n P(i,j) \right)}_{=1} = \|x\|_1 \quad \square$$

$$\sum_{i=1}^n P(i,j) = \frac{\sum_{i=1}^n w(j,i)}{\sum_{k=1}^n w(j,k)} = 1$$

Proof of Lemma 3 $x = (1 - \alpha)v + \alpha Px.$

$$x - \alpha Px = (1 - \alpha)v$$

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

A

$$Ax = v$$

Look at $\ker(A)$. Claim $\ker(A) = \{0\}$.

Let $z \in \ker(A)$, so $Az = 0$. Need to show $z = 0$.

$$A = (1-\alpha)^{-1}(I - \alpha P)$$

$$Az = 0 \iff z - \alpha Pz = 0$$

$$\iff z = \alpha Pz$$

$$\|z\|_1 = \|\alpha Pz\|_1$$

$$= \alpha \|Pz\|_1 \leq \alpha \|z\|_1$$

$$(1-\alpha)\|z\|_1 \leq 0, \quad \alpha < 1$$

$$\Rightarrow 1 - \alpha = 0 \quad \text{or} \quad \|z\|_2 = 0$$

$$\alpha < 1 \Rightarrow \|z\|_2 = 0 \quad \square$$

(i) We have $\sum_{i=1}^n x(i) = \sum_{i=1}^n v(i)$.

(ii) If $v(i) \geq 0$ for all i , then $x(i) \geq 0$ for all i .

$$(i) \quad x = (1 - \alpha)v + \alpha Px.$$

$$x(i) = (1 - \alpha)v(i) + \alpha \sum_{j=1}^n P(i,j)x(j)$$

$$\sum_{i=1}^n x(i) = (1 - \alpha) \sum_{i=1}^n v(i) + \alpha \sum_{i=1}^n \sum_{j=1}^n P(i,j)x(j)$$

$$= (1-\alpha) \sum_{i=1}^{\infty} v(i) + \alpha \sum_{j=1}^{\infty} \left[\sum_{i=1}^{\infty} P(i,j) \right] x(j)$$

$$= (1-\alpha) \sum_{i=1}^{\infty} v(i) + \alpha \sum_{j=1}^{\infty} x(j)$$

$$\cancel{(1-\alpha)} \sum_{i=1}^{\infty} x(i) = \cancel{(1-\alpha)} \sum_{i=1}^{\infty} v(i)$$

$$\alpha < 1$$

\square

(i) We have $\sum_{i=1}^n x(i) = \sum_{i=1}^n v(i)$.

(ii) If $v(i) \geq 0$ for all i , then $x(i) \geq 0$ for all i .

$$x(i) = (1-\alpha)v(i) + \alpha \sum_{j=1}^n P(i,j)x(j)$$

$$|x(i)| = \left| (1-\alpha)v(i) + \alpha \sum_{j=1}^n P(i,j)x(j) \right|$$

$v(i) \geq 0 \downarrow$

$$\leq (1-\alpha)v(i) + \alpha \sum_{j=1}^n P(i,j)|x(j)|$$

$$\sum_{i=1}^n |x(i)| \leq (1-\alpha) \sum_{i=1}^n v(i) + \alpha \sum_{j=1}^n |x(j)|$$

$$= (1-\alpha) \sum_{i=1}^n x(i) + \alpha \sum_{j=1}^n |x(j)|$$

$$\cancel{(1-\alpha)} \sum_{i=1}^n |x(i)| \leq \cancel{(1-\alpha)} \sum_{i=1}^n x(i)$$

$$\Rightarrow \sum_{i=1}^n |x(i)| = \sum_{i=1}^n x(i)$$

$$\sum_{i=1}^n \underbrace{(|x(i)| - x(i))}_{=0} = 0$$

$$= 0$$

\Rightarrow

$$x(i) = |x(i)|$$

$$\geq 0$$



Eigenvector problem

Remark 5. When v is a probability distribution, it is common to re-write the PageRank problem (2) as an eigenvector problem

$$P_\alpha x = x$$

where

$$P_\alpha := (1 - \alpha)v\mathbf{1}^T + \alpha P.$$

$$x - \alpha Px = (1 - \alpha)v$$

$$\alpha Px + (1 - \alpha)v\mathbf{1}^T x = x$$

$$\underbrace{(\alpha P + (1 - \alpha)v\mathbf{1}^T)}_{P_\alpha} x = x$$

P_α

$$P_\alpha x = x$$

$$x(i) \geq 0$$

$$v(i) \geq 0$$

$$\sum x(i) = 1$$

$$\sum v(i) = 1$$

$$\mathbf{1}^T x = 1$$

$$\mathbf{1}^T v = 1$$

Convergence of the PageRank iteration

Let $v \in \mathbb{R}^n$ and $0 \leq \alpha < 1$. Let x_k satisfy the PageRank iteration

$$x_{k+1} = (1 - \alpha)v + \alpha Px_k,$$



and let x be the unique solution of the PageRank problem

$$x = (1 - \alpha)v + \alpha Px.$$



Theorem 6. *We have*

$$(3) \quad \|x_k - x\|_1 \leq \alpha^k \|x_0 - x\|_1.$$



Since $0 \leq \alpha < 1$, this is convergence of $x_k \rightarrow x$ with a **linear** convergence rate of α .

Proof:

$$x_{k+1} = (1 - \alpha)v + \alpha Px_k,$$

$$x = (1 - \alpha)v + \alpha Px.$$

$$\begin{aligned} x_{k+1} - x &= \alpha Px_k - \alpha Px \\ &= \alpha P(x_k - x) \end{aligned}$$

$$\begin{aligned} \|x_{k+1} - x\|_2 &= \|\alpha P(x_k - x)\|_2 \\ &= \alpha \|P(x_k - x)\|_2 \end{aligned}$$

$$\|Px\|_2 \leq \|x\|_2 \quad \leq \alpha \|x_k - x\|_2$$


Induction

$$\|x_{k+1} - x\|_2 \leq \alpha^{k+1} \|x_0 - x\|_2$$



Power iteration

Remark 7. In the eigenvector formulation discussed above, the PageRank iteration $x_{k+1} = P_\alpha x_k$ is basically the power iteration to find the largest eigenvector of P . The normalization step is not needed since $\|x_k\|_1 = 1$ for all k .



Personalized PageRank for image retrieval ([.ipynb](#))