

# Authoritative Sources in a Hyperlinked Environment

Dimitris Sakavalas

June 15, 2010

# Outline

## Introduction

- Problem description

- Authorities and Broad-topic queries

- Algorithm outline

## Construction of a Focused Subgraph of the www

- Expanding  $R_\sigma$  into the base set  $S_\sigma$

- Heuristics improving base set  $S_\sigma$

## Computing Hubs and Authorities

- Iterative Algorithm

- Elements of Linear Algebra

- Convergence of the Iterative Algorithm

## Similar page queries

## Multiple Sets of Hubs and Authorities

- Non-Principal Eigenvectors and Clustering of  $G_\sigma$

# Introduction

## Problem description

### *www*-search

- Discover pages relevant to a given **query string**
- Discover the most **authoritative** pages  
(subjective, requires human evaluation)

# Introduction

## Problem description

### *www*-search

- Discover pages relevant to a given **query string**
- Discover the most **authoritative** pages  
(subjective, requires human evaluation)

### Queries

- **Specific** queries → *Scarcity Problem*
- **Broad-topic** queries → *Abundance Problem*
- **Similar-page** queries

# Introduction

Authorities and Broad-topic queries

In search of a definition of authority

# Introduction

## Authorities and Broad-topic queries

### In search of a definition of authority

- No endogenous measure of a page to assess authority (ex. *Text-based* ranking functions for queries: “Harvard”, “search engines”, “automobile manufacturers”, images)

# Introduction

## Authorities and Broad-topic queries

### In search of a definition of authority

- No endogenous measure of a page to assess authority (ex. *Text-based* ranking functions for queries: “Harvard”, “search engines”, “automobile manufacturers”, images)
- Solution: *Analysis of the link structure*



# Introduction

## Authorities and Broad-topic queries

### In search of a definition of authority

- No endogenous measure of a page to assess authority (ex. *Text-based* ranking functions for queries: “Harvard”, “search engines”, “automobile manufacturers”, images)
- Solution: *Analysis of the link structure*



### Pitfalls:

- Advertising links
- Navigational links
- Universal Popularity links



# Introduction

## Authorities and Broad-topic queries

### In search of a definition of authority

- No endogenous measure of a page to assess authority (ex. *Text-based* ranking functions for queries: “Harvard”, “search engines”, “automobile manufacturers”, images)
- Solution: *Analysis of the link structure*



### Pitfalls:

- Advertising links
- Navigational links
- Universal Popularity links

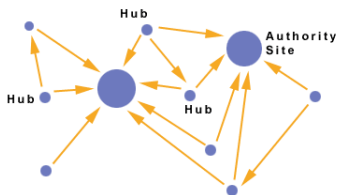
**Question.** How reliable is page  $p$ ?

# Introduction

Algorithm outline

## Hubs

Pages that link to many Authorities

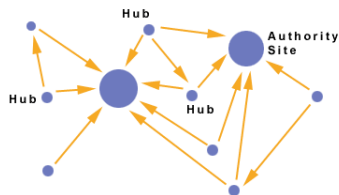


# Introduction

## Algorithm outline

### Hubs

Pages that link to many Authorities



## Outline of the algorithm

- 1 Create a **focused subgraph** of the www from the output of a text-based engine (Small collection of pages likely to contain the most authoritative pages)
- 2 Identify *hubs* and *authorities*

## Construction of a Focused Subgraph of the www

**Directed graph  $G = (V, E)$  of the www.**

$V$  the collection of hyperlinked pages of the www

- $V$ : (nodes correspond to the pages)
- $(p, q) \in E \Leftrightarrow$  There is a link from  $p$  to  $q$

**$G[W]$ .** Subgraph of  $G$  induced on  $W \subseteq V$

# Construction of a Focused Subgraph of the www

**Directed graph  $G = (V, E)$  of the www.**

$V$  the collection of hyperlinked pages of the www

- $V$ : (nodes correspond to the pages)
- $(p, q) \in E \Leftrightarrow$  There is a link from  $p$  to  $q$

**$G[W]$ .** Subgraph of  $G$  induced on  $W \subseteq V$

Goal

Construct  $G_\sigma = G[S_\sigma]$  such that

- ①  $|S_\sigma|$  is relatively small  $\rightarrow$  (computational cost)
- ②  $S_\sigma$  rich in relevant pages  $\rightarrow$  (search quality)
- ③  $S_\sigma$  contains most of the *strongest authorities*

# Construction of a Focused Subgraph of the www

**Directed graph  $G = (V, E)$  of the www.**

$V$  the collection of hyperlinked pages of the www

- $V$ : (nodes correspond to the pages)
- $(p, q) \in E \Leftrightarrow$  There is a link from  $p$  to  $q$

**$G[W]$ .** Subgraph of  $G$  induced on  $W \subseteq V$

Goal

Construct  $G_\sigma = G[S_\sigma]$  such that

- ①  $|S_\sigma|$  is relatively small  $\rightarrow$  (computational cost)
- ②  $S_\sigma$  rich in relevant pages  $\rightarrow$  (search quality)
- ③  $S_\sigma$  contains most of the *strongest authorities*

Observation

$S_\sigma$  rich in relevant pages  $\rightarrow$  many links to authorities

# Construction of a Focused Subgraph of the www

## Notation

- $\Gamma^+(p)$ : the set of all pages  $p$  points to
- $\Gamma^-(p)$ : the set of all pages pointing to  $p$
- $\mathcal{E}$ : a *text-based search engine*

# Construction of a Focused Subgraph of the www

## Notation

- $\Gamma^+(p)$ : the set of all pages  $p$  points to
- $\Gamma^-(p)$ : the set of all pages pointing to  $p$
- $\mathcal{E}$ : a *text-based search engine*

## 1st Step

Root set  $R_\sigma = \{t \text{ highest ranked pages for the query } \sigma \text{ from } \mathcal{E}\}$

- $R_\sigma$  satisfies (i), (ii) .
- Generally doesn't satisfy (iii) and is *structureless* (few links)



# Construction of a Focused Subgraph of the www

## Notation

- $\Gamma^+(p)$ : the set of all pages  $p$  points to
- $\Gamma^-(p)$ : the set of all pages pointing to  $p$
- $\mathcal{E}$ : a *text-based search engine*

## 1st Step

Root set  $R_\sigma = \{t \text{ highest ranked pages for the query } \sigma \text{ from } \mathcal{E}\}$

- $R_\sigma$  satisfies (i), (ii) .
- Generally doesn't satisfy (iii) and is *structureless* (few links)

## Observation

*If  $a \notin R_\sigma$ , strong authority, there is a high probability that  $\exists p \in R_\sigma$  such that  $a \in \Gamma^+(p)$*

## Expanding $R_\sigma$ into the base set $S_\sigma$

### 2nd Step

**Idea:** Expand  $R_\sigma$  into the Base set  $S_\sigma$  by adding strong authorities and relevant pages

## Expanding $R_\sigma$ into the base set $S_\sigma$

### 2nd Step

**Idea:** Expand  $R_\sigma$  into the Base set  $S_\sigma$  by adding strong authorities and relevant pages

### Subgraph( $\sigma, \mathcal{E}, t, d$ )

$t, d \in \mathbb{N}$

Set  $S_\sigma := R_\sigma$

For each page  $p \in R_\sigma$

    Add all pages in  $\Gamma^+(p)$  to  $S_\sigma$

    If  $|\Gamma^-(p)| \leq d$  then,

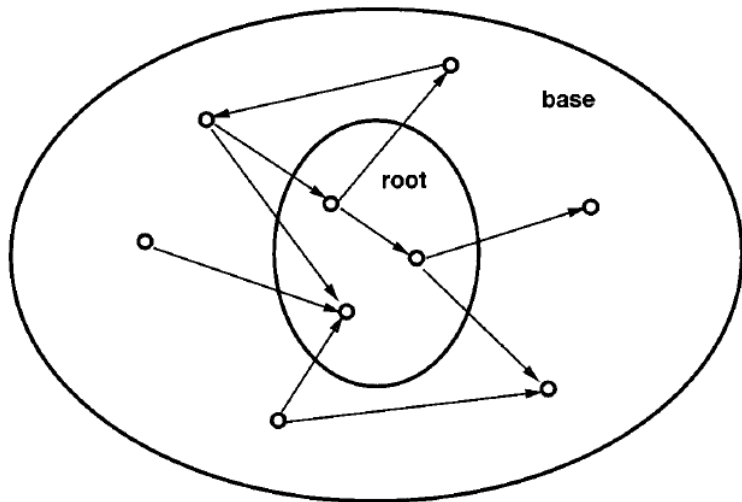
        Add all pages in  $\Gamma^-(p)$  to  $S_\sigma$

    Else

        Add an arbitrary set of  $d$  pages from  $\Gamma^-(p)$  to  $S_\sigma$

End.

# Expanding $R_\sigma$ into the base set $S_\sigma$



# Heuristics improving $S_\sigma$

## Links

- *Transverse*: link between pages with different domain name
- *Intrinsic*: link between pages with the same domain name (navigational)

# Heuristics improving $S_\sigma$

## Links

- *Transverse*: link between pages with different domain name
- *Intrinsic*: link between pages with the same domain name (navigational)

## 3rd Step: Heuristics

- 1 Delete all *Intrinsic* links from the graph  $G_\sigma$
- 2 Allow up to a small number  $m$  pages from a single domain to point to any given page  $p$

# Heuristics improving $S_\sigma$

## Links

- *Transverse*: link between pages with different domain name
- *Intrinsic*: link between pages with the same domain name (navigational)

## 3rd Step: Heuristics

- 1 Delete all *Intrinsic* links from the graph  $G_\sigma$
- 2 Allow up to a small number  $m$  pages from a single domain to point to any given page  $p$

Finally we obtain a small subgraph  $G_\sigma$ , relatively focused on the query  $\sigma$ , containing many relevant pages and strong authorities.

# Hubs and Authorities

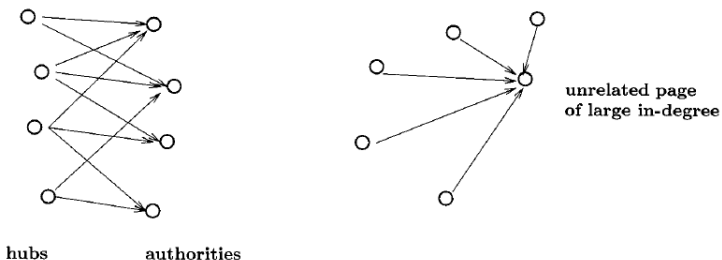
Purely In-degree ordering of Authorities



# Hubs and Authorities

## Purely In-degree ordering of Authorities

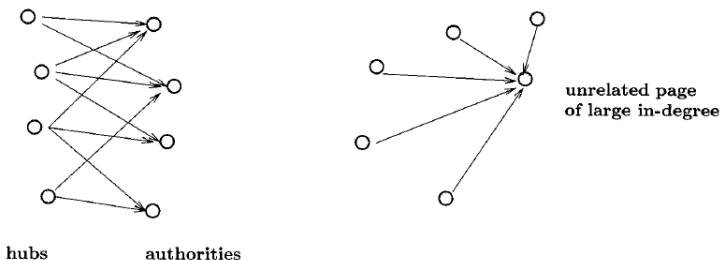
**Problem.** “Universally popular” pages, large in-degree regardless of the underlying query topic (ex. Amazon Books)



# Hubs and Authorities

## Purely In-degree ordering of Authorities

**Problem.** “Universally popular” pages, large in-degree regardless of the underlying query topic (ex. Amazon Books)



## Observation

**Authorities:** *high in-degree, considerable overlap in the sets of pages that point to them(hubs)*

# Hubs and Authorities

Observation (mutually reinforcing relationship)

- **Good Hub:** *a page that points to many good authorities*
- **Good Authority:** *a page that is pointed by many good hubs*

# Hubs and Authorities

Observation (mutually reinforcing relationship)

- **Good Hub:** *a page that points to many good authorities*
- **Good Authority:** *a page that is pointed by many good hubs*

Definitions

$\forall p \in S_\sigma = \{1, 2, \dots, n\}$  assign:

- authority weight  $x^{<p>} \geq 0$
- hub weight  $y^{<p>} \geq 0$

# Hubs and Authorities

## Observation (mutually reinforcing relationship)

- **Good Hub**: a page that points to many good authorities
- **Good Authority**: a page that is pointed by many good hubs

## Definitions

$\forall p \in S_\sigma = \{1, 2, \dots, n\}$  assign:

- authority weight  $x^{<p>} \geq 0$
- hub weight  $y^{<p>} \geq 0$

- corresponding vectors  $x = \begin{bmatrix} x^{<1>} \\ x^{<2>} \\ \vdots \\ x^{<n>} \end{bmatrix}$  and  $y = \begin{bmatrix} y^{<1>} \\ y^{<2>} \\ \vdots \\ y^{<n>} \end{bmatrix}$

Normalized such that  $\|x\| = \|y\| = 1$

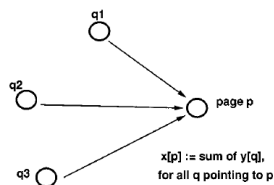
**example.** For norm  $\|\cdot\|_1$ ,  $\sum_{i=1}^n x^{<i>} = \sum_{i=1}^n y^{<i>} = 1$

# Hubs and Authorities

## Definitions

- $\mathcal{I}(x, y)$  ( $x$ -weight update):

$$\forall x^{<p>}, x^{<p>} \leftarrow \sum_{q:(q,p) \in E} y^{<q>}$$



# Hubs and Authorities

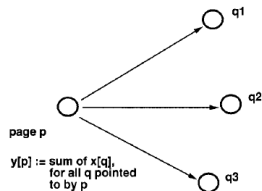
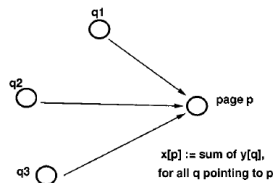
## Definitions

- $\mathcal{I}(x, y)$  ( $x$ -weight update):

$$\forall x^{<p>}, x^{<p>} \leftarrow \sum_{q:(q,p) \in E} y^{<q>}$$

- $\mathcal{O}(x, y)$  ( $y$ -weight update):

$$\forall y^{<p>}, y^{<p>} \leftarrow \sum_{q:(p,q) \in E} x^{<q>},$$



# Hubs and Authorities

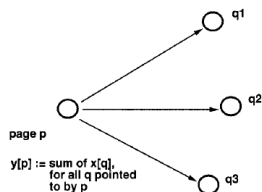
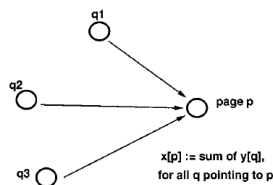
## Definitions

- $\mathcal{I}(x, y)$  ( $x$ -weight update):

$$\forall x^{<p>}, x^{<p>} \leftarrow \sum_{q:(q,p) \in E} y^{<q>}$$

- $\mathcal{O}(x, y)$  ( $y$ -weight update):

$$\forall y^{<p>}, y^{<p>} \leftarrow \sum_{q:(p,q) \in E} x^{<q>},$$



## Idea

Find the desired “*equilibrium*” values for the weights by applying  $\mathcal{I}, \mathcal{O}$  operations iteratively.



# Hubs and Authorities

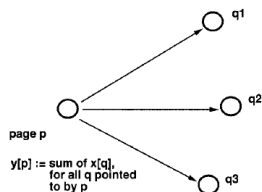
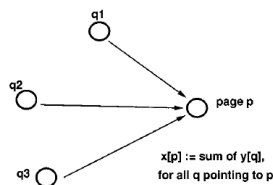
## Definitions

- $\mathcal{I}(x, y)$  ( $x$ -weight update):

$$\forall x^{<p>}, x^{<p>} \leftarrow \sum_{q:(q,p) \in E} y^{<q>}$$

- $\mathcal{O}(x, y)$  ( $y$ -weight update):

$$\forall y^{<p>}, y^{<p>} \leftarrow \sum_{q:(p,q) \in E} x^{<q>},$$



## Idea

Find the desired “*equilibrium*” values for the weights by applying  $\mathcal{I}, \mathcal{O}$  operations iteratively.

**Question.** Does this “*equilibrium*” exist?

# Hubs and Authorities

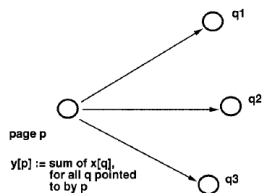
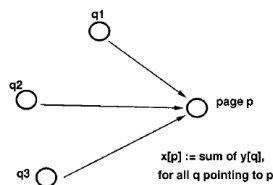
## Definitions

- $\mathcal{I}(x, y)$  ( $x$ -weight update):

$$\forall x^{<p>}, x^{<p>} \leftarrow \sum_{q:(q,p) \in E} y^{<q>}$$

- $\mathcal{O}(x, y)$  ( $y$ -weight update):

$$\forall y^{<p>}, y^{<p>} \leftarrow \sum_{q:(p,q) \in E} x^{<q>},$$



## Idea

Find the desired “*equilibrium*” values for the weights by applying  $\mathcal{I}, \mathcal{O}$  operations iteratively.

**Question.** Does this “*equilibrium*” exist? → YES (linear algebra)

# Iterative algorithm

- **Iterate**( $G_\sigma, k$ )

$z := (1, 1, \dots, 1) \in \mathbb{R}^n$

Set  $x_0 := z$

Set  $y_0 := z$

For  $i = 1, 2, \dots, k$

    Apply the  $\mathcal{I}(x_{i-1}, y_{i-1})$  operation to obtain new  $x$ -weights  $x'_i$

    Apply the  $\mathcal{O}(x'_i, y_{i-1})$  operation to obtain new  $y$ -weights  $y'_i$

    Normalize  $x'_i$ , obtaining  $x_i$

    Normalize  $y'_i$ , obtaining  $y_i$

End

Return  $(x_k, y_k)$

# Iterative algorithm

- **Iterate**( $G_\sigma, k$ )

$z := (1, 1, \dots, 1) \in \mathbb{R}^n$

Set  $x_0 := z$

Set  $y_0 := z$

For  $i = 1, 2, \dots, k$

    Apply the  $\mathcal{I}(x_{i-1}, y_{i-1})$  operation to obtain new  $x$ -weights  $x'_i$

    Apply the  $\mathcal{O}(x'_i, y_{i-1})$  operation to obtain new  $y$ -weights  $y'_i$

    Normalize  $x'_i$ , obtaining  $x_i$

    Normalize  $y'_i$ , obtaining  $y_i$

End

Return  $(x_k, y_k)$

- **Filter**( $G_\sigma, k, c$ )

$(x_k, y_k) := \text{Iterate}(G_\sigma, k)$

Report the pages with the  $c$  largest coordinates in  $x_k$  as authorities

Report the pages with the  $c$  largest coordinates in  $y_k$  as hubs

# Elements of Linear Algebra

## Notation

Matrix  $M \in \mathbb{C}^{n \times n}$

- *Eigenvalues*  $\lambda_i(M)$
- *Eigenvectors*  $\omega_i(M)$
- *eigenspace*  $V_\lambda$  subspace of  $\mathbb{C}^n$
- *multiplicity of  $\lambda$*  :  $m_\lambda = \dim(V_\lambda)$

# Elements of Linear Algebra

## Notation

Matrix  $M \in \mathbb{C}^{n \times n}$

- *Eigenvalues*  $\lambda_i(M)$
- *Eigenvectors*  $\omega_i(M)$
- *eigenspace*  $V_\lambda$  subspace of  $\mathbb{C}^n$
- *multiplicity of  $\lambda$*  :  $m_\lambda = \dim(V_\lambda)$

## Definitions

- *Positive matrix*  $M$  : If  $M_{i,j} > 0, \forall i, j$
- *Primitive matrix*  $M$  : If  $\exists m \in \mathbb{N}, M^m$  : positive

# Elements of Linear Algebra

Let  $M \in \mathbb{R}^{n \times n}$  *symmetric, nonnegative matrix* ( $M_{i,j} \geq 0, \forall i, j$ )

## Theorems

- 1  $M$  has at most  $n$  distinct eigenvalues and  $\sum m_{\lambda_i} = n$

# Elements of Linear Algebra

Let  $M \in \mathbb{R}^{n \times n}$  *symmetric, nonnegative matrix* ( $M_{i,j} \geq 0, \forall i, j$ )

## Theorems

- ①  $M$  has at most  $n$  distinct eigenvalues and  $\sum m_{\lambda_i} = n$
- ②  $M$  has *only real eigenvalues*  $\lambda_i$   
Denote them  $|\lambda_1(M)| \geq |\lambda_2(M)| \geq \dots \geq |\lambda_n(M)|$



# Elements of Linear Algebra

Let  $M \in \mathbb{R}^{n \times n}$  symmetric, nonnegative matrix ( $M_{i,j} \geq 0, \forall i, j$ )

## Theorems

①  $M$  has at most  $n$  distinct eigenvalues and  $\sum m_{\lambda_i} = n$

②  $M$  has only real eigenvalues  $\lambda_i$

Denote them  $|\lambda_1(M)| \geq |\lambda_2(M)| \geq \dots \geq |\lambda_n(M)|$

③ (*Perron-Frobenius*) If  $M$  is primitive then:

- i. Largest eigenvalue  $\lambda_1(M) > 0$  and  $m_\lambda = 1$
- ii.  $\lambda_1 > |\lambda_i| \forall i \neq 1$
- iii.  $\lambda_1$  has a corresponding eigenvector  $\omega_1(M)$  with all entries positive (**the principal eigenvector**)

Consider  $\omega_1(M)$  normalized,  $\|\omega_1(M)\| = 1$

# Elements of Linear Algebra

Let  $M \in \mathbb{R}^{n \times n}$  symmetric, nonnegative matrix ( $M_{i,j} \geq 0, \forall i, j$ )

## Theorems

- 1  $M$  has at most  $n$  distinct eigenvalues and  $\sum m_{\lambda_i} = n$
- 2  $M$  has only real eigenvalues  $\lambda_i$   
Denote them  $|\lambda_1(M)| \geq |\lambda_2(M)| \geq \dots \geq |\lambda_n(M)|$
- 3 (*Perron-Frobenius*) If  $M$  is primitive then:
  - i. Largest eigenvalue  $\lambda_1(M) > 0$  and  $m_\lambda = 1$
  - ii.  $\lambda_1 > |\lambda_i| \forall i \neq 1$
  - iii.  $\lambda_1$  has a corresponding eigenvector  $\omega_1(M)$  with all entries positive (**the principal eigenvector**)  
Consider  $\omega_1(M)$  normalized,  $\|\omega_1(M)\| = 1$
- 4 If  $M$  is primitive and vector  $v$  not orthogonal to the principal eigenvector  $\omega_1(M)$ , let  $v_k$  be the unit vector in the direction of  $M^k v$ . Then  $v_k \xrightarrow{k \rightarrow \infty} \omega_1(M)$

# Elements of Linear Algebra

## Observation

*Note that we can consider  $\omega_1(M)$  as a vector of the orthonormal base of  $V_{\lambda_1}$ . Moreover  $\dim V_{\lambda_1} = m_{\lambda_1} = 1$  and all entries in  $\omega_1(M)$  are positive. Hence the principal eigenvector  $\omega_1(M)$  is unique.*

## Connection of Iterative algorithm to linear Algebra

Let  $A$  be the adjacency matrix of  $G_\sigma$

Observation

Operations  $\mathcal{I}(x, y)$  and  $\mathcal{O}(x, y)$  can be written as

$$x \leftarrow A^\top y \text{ and } y \leftarrow Ax \text{ respectively}$$

# Connection of Iterative algorithm to linear Algebra

Let  $A$  be the adjacency matrix of  $G_\sigma$

Observation

Operations  $\mathcal{I}(x, y)$  and  $\mathcal{O}(x, y)$  can be written as

$$x \leftarrow A^\top y \text{ and } y \leftarrow Ax \text{ respectively}$$

Iterate Algorithm

The algorithm initializes  $x_0 = y_0 = z = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$ . In the  $i$ -th step of the algorithm we have the following steps:

- $\mathcal{I}(x_i, y_{i-1}) \rightarrow x_i = A^\top y_{i-1}$
- $\mathcal{O}(x_i, y_i) \rightarrow y_i = Ax_i$

# Connection of Iterative algorithm to linear Algebra

Let  $A$  be the adjacency matrix of  $G_\sigma$

Observation

Operations  $\mathcal{I}(x, y)$  and  $\mathcal{O}(x, y)$  can be written as

$$x \leftarrow A^\top y \text{ and } y \leftarrow Ax \text{ respectively}$$

Iterate Algorithm

The algorithm initializes  $x_0 = y_0 = z = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$ . In the  $i$ -th step of the algorithm we have the following steps:

- $\mathcal{I}(x_i, y_{i-1}) \rightarrow x_i = A^\top y_{i-1}$
- $\mathcal{O}(x_i, y_i) \rightarrow y_i = Ax_i$

$$\Rightarrow \left. \begin{array}{l} x_i = A^\top Ax_{i-1} = (A^\top A)^i \cdot x_0 \\ y_i = AA^\top y_{i-1} = (AA^\top)^i \cdot y_0 \end{array} \right\} \Rightarrow \begin{array}{l} x_i = (A^\top A)^i \cdot z \\ y_i = (AA^\top)^i \cdot z \end{array} \quad (*)$$

# Convergence of the Iterative Algorithm

## Observation

- *We consider  $x_i, y_i$  normalized in each step as in Iterate algorithm. So in every step  $\|x_i\| = \|y_i\| = 1$*
- *For simplicity, matrices  $AA^\top$  and  $A^\top A$  are considered primitive*
- *Obviously  $z \cdot c \neq 0, \forall c \in \mathbb{R}_{\geq 0}^n$*

# Convergence of the Iterative Algorithm

## Observation

- *We consider  $x_i, y_i$  normalized in each step as in Iterate algorithm. So in every step  $\|x_i\| = \|y_i\| = 1$*
- *For simplicity, matrices  $AA^\top$  and  $A^\top A$  are considered primitive*
- *Obviously  $z \cdot c \neq 0, \forall c \in \mathbb{R}_{\geq 0}^n$*

## Theorem

*The sequences  $\{x_i\}$  and  $\{y_i\}$  produced by Iterate Algorithm converge to limits  $x^*, y^*$  respectively*



# Convergence of the Iterative Algorithm

## Observation

- We consider  $x_i, y_i$  normalized in each step as in Iterate algorithm. So in every step  $\|x_i\| = \|y_i\| = 1$
- For simplicity, matrices  $AA^\top$  and  $A^\top A$  are considered primitive
- Obviously  $z \cdot c \neq 0, \forall c \in \mathbb{R}_{\geq 0}^n$

## Theorem

*The sequences  $\{x_i\}$  and  $\{y_i\}$  produced by Iterate Algorithm converge to limits  $x^*, y^*$  respectively*

Proof.

$A^\top A$  primitive,  $z$  not orthogonal to the principal eigenvector  $\omega_1(A^\top A) = x^*$ ,  $x_i$  the unit vector in the direction  $(A^\top A)^i \cdot z$  Hence from *Theorem 4*  $\Rightarrow \{x_i\} \rightarrow x^*$ , similarly  $\{y_i\} \rightarrow y^* = \omega_1(AA^\top)$

□

# Conclusion

## Conclusion

The process of iterated  $\mathcal{I}/\mathcal{O}$  operations will converge. The desired equilibrium values for the  $x, y$  weights, are the values  $x^*, y^*$ .

# Conclusion

## Conclusion

The process of iterated  $\mathcal{I}/\mathcal{O}$  operations will converge. The desired equilibrium values for the  $x, y$  weights, are the values  $x^*, y^*$ .

## Observation

In practice the convergence of *Iterate* is quite rapid, so one can compute  $x, y$  weights by starting from any initial vectors  $x_0, y_0$ , and performing a fixed number of  $\mathcal{I}, \mathcal{O}$  operations

# Conclusion

## Conclusion

The process of iterated  $\mathcal{I}/\mathcal{O}$  operations will converge. The desired equilibrium values for the  $x, y$  weights, are the values  $x^*, y^*$ .

## Observation

In practice the convergence of *Iterate* is quite rapid, so one can compute  $x, y$  weights by starting from any initial vectors  $x_0, y_0$ , and performing a fixed number of  $\mathcal{I}, \mathcal{O}$  operations

*Hyperlink-Induced Topic Search (HITS) algorithm (Ask.com)*

# Similar-Page Queries

**Query.** Find pages similar to  $p$

Idea

- Initiate search with the page  $p$  instead of a query string  $\sigma$
- Use the link structure to infer “similarity among pages”
- Similar pages to  $p \rightarrow$  Strongest authorities in local region of  $p$

# Similar-Page Queries

**Query.** Find pages similar to  $p$

Idea

- Initiate search with the page  $p$  instead of a query string  $\sigma$
- Use the link structure to infer “similarity among pages”
- Similar pages to  $p \rightarrow$  Strongest authorities in local region of  $p$

Adaptation of the broad-topic query method

- ① Find  $t$  pages pointing to  $p \rightarrow$  Assemble *root set*  $R_p$
- ② Expand  $R_p$  to *base set*  $S_p$  as before
- ③ Search for authorities and hubs in the focused subgraph  $G_p$

# Multiple Sets of Hubs and Authorities

**Situation.** For a query string  $\sigma$ , relevant pages grouped into *clusters*. Reasons

- Different meanings of  $\sigma$ . ex. “jaguar”
- String  $\sigma$  arises as term in multiple technical communities. ex. “randomized algorithms”
- String  $\sigma$  refers to a highly polarized issue. ex. “abortion”

# Multiple Sets of Hubs and Authorities

**Situation.** For a query string  $\sigma$ , relevant pages grouped into *clusters*. Reasons

- Different meanings of  $\sigma$ . ex. “jaguar”
- String  $\sigma$  arises as term in multiple technical communities. ex. “randomized algorithms”
- String  $\sigma$  refers to a highly polarized issue. ex. “abortion”

**Problem.** Identify the *clusters*



# Multiple Sets of Hubs and Authorities

**Situation.** For a query string  $\sigma$ , relevant pages grouped into *clusters*. Reasons

- Different meanings of  $\sigma$ . ex. “jaguar”
- String  $\sigma$  arises as term in multiple technical communities. ex. “randomized algorithms”
- String  $\sigma$  refers to a highly polarized issue. ex. “abortion”

**Problem.** Identify the *clusters*

Idea

- Clusters represented in focused subgraph  $G_\sigma$  as densely linked subgraphs(collection of hubs and authorities)
- Extract densely linked collections of hubs and authorities through *non-principal* eigenvectors of  $AA^\top$  and  $A^\top A$

# Non-Principal Eigenvectors and Clustering of $G_\sigma$

## Proposition

*$AA^\top$  and  $A^\top A$  have the same multiset of eigenvalues, and their eigenvectors can be chosen so that  $\omega_i(AA^\top) = A\omega_i(A^\top A)$*

# Non-Principal Eigenvectors and Clustering of $G_\sigma$

## Proposition

*$AA^\top$  and  $A^\top A$  have the same multiset of eigenvalues, and their eigenvectors can be chosen so that  $\omega_i(AA^\top) = A\omega_i(A^\top A)$*

## Observation

- *For each pair of eigenvectors  $x_i^* = \omega_i(A^\top A)$ ,  $y_i^* = \omega_i(AA^\top)$  operation  $\mathcal{I}(x_i^*, y_i^*)$  keeps  $x$ -weights parallel to  $x_i^*$ :  
 $x = A^\top \omega_i(AA^\top) = A^\top A \omega_i(A^\top A) = \lambda_i x_i^*$   
Similarly operation  $\mathcal{O}(x_i^*, y_i^*)$  keeps  $x$ -weights parallel to  $y_i^*$*

# Non-Principal Eigenvectors and Clustering of $G_\sigma$

## Proposition

*$AA^\top$  and  $A^\top A$  have the same multiset of eigenvalues, and their eigenvectors can be chosen so that  $\omega_i(AA^\top) = A\omega_i(A^\top A)$*

## Observation

- *For each pair of eigenvectors  $x_i^* = \omega_i(A^\top A)$ ,  $y_i^* = \omega_i(AA^\top)$  operation  $\mathcal{I}(x_i^*, y_i^*)$  keeps  $x$ -weights parallel to  $x_i^*$ :  
 $x = A^\top \omega_i(AA^\top) = A^\top A \omega_i(A^\top A) = \lambda_i x_i^*$   
Similarly operation  $\mathcal{O}(x_i^*, y_i^*)$  keeps  $x$ -weights parallel to  $y_i^*$*
- *Non-principal eigenvectors have both positive and negative entries. Hence each pair of  $(x_i^*, y_i^*)$  provide us with two densely connected set of hubs and authorities*
  - *Pages that correspond to the  $c$  coordinates with the most positive values*
  - *Pages that correspond to the  $c$  coordinates with the most negative values*

# Non-Principal Eigenvectors and Clustering of $G_\sigma$

## Description of the method

- For each of the first few non-principal eigenvectors  $(x_i^*, y_i^*)$  find the two densely connected set of hubs and authorities through an algorithm similar (but less clean conceptually) to the *Iterate algorithm*

# Non-Principal Eigenvectors and Clustering of $G_\sigma$

## Description of the method

- For each of the first few non-principal eigenvectors  $(x_i^*, y_i^*)$  find the two densely connected set of hubs and authorities through an algorithm similar (but less clean conceptually) to the *Iterate algorithm*
- The pages with large coordinates in the first few nonprincipal eigenvectors tend to recur, so that essentially the same collection of hubs and authorities will often be generated by several of the strongest nonprincipal eigenvectors

## Similar Concepts in Other Areas

- *Social networks* → “*Standing*”
- *Bibliometrics* → “*Impact*”

## Similar Concepts in Other Areas

- *Social networks* → “*Standing*”
- *Bibliometrics* → “*Impact*”
- *Economics*  
(Wassily Leontief 1941-Nobel prize 1973)