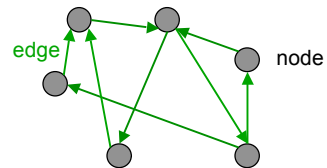


Link-based ranking Part 2

1

Goal

- **Intuition:** when Web page **points** to another Web page, it **confers status/authority/popularity** to that page
- Find a measure that **captures intuition**

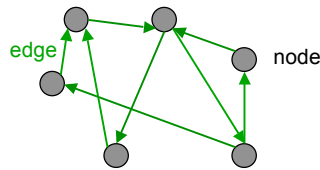


- **Not just web linking**
 - Citations in books, articles
 - Doctors referring to other doctors

2

Review: first measure PageRank

- Given a directed graph with n nodes
- Assign each node a score that represents its importance in structure
 - Call score **PageRank**: $pr(\text{node})$



3

Conferring importance

Core ideas:

- A node should **confer** some of its importance **to the nodes to which it points**
 - If a node is important, the nodes it links to should be important
- A node should **not transfer more** importance **than it has**
- Address **problems** with:
 - **Sinks** (nodes with no edges out)
 - **Cyclic** behavior

4

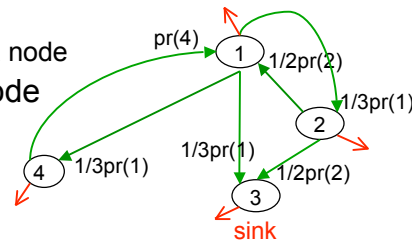
Random walk model (review)

1. **Move** from node to linked neighbor with probability $1/\text{outdegree}$

Outdegree of a node
= number of edges out of a node

2. **Randomly jump** to any node

- Break cycles
- Escape from sinks



Captured with:

$$\text{pr}_{\text{new}}(\mathbf{k}) = \alpha/n + (1-\alpha) \sum_{i \text{ with edge from } i \text{ to } k} (\text{pr}(i) / t_i)$$

- α parameter chosen empirically
- t_i outdegree of node i

Steady state probability of being at a node = $\text{pr}(\text{node})$

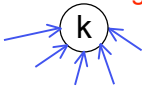
5

Normalized?

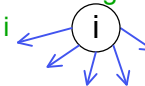
- Would like $\sum_{1 \leq k \leq n} (\text{pr}(k)) = 1$
- Consider $\sum_{1 \leq k \leq n} (\text{pr}_{\text{new}}(k))$

$$\begin{aligned} &= \sum_{1 \leq k \leq n} (\alpha/n + (1-\alpha) \sum_{i \text{ with edge from } i \text{ to } k} (\text{pr}(i) / t_i)) \\ &= \sum_{1 \leq k \leq n} (\alpha/n) + \sum_{1 \leq k \leq n} ((1-\alpha) \sum_{i \text{ with edge from } i \text{ to } k} (\text{pr}(i) / t_i)) * \\ &= \alpha + (1-\alpha) \sum_{1 \leq k \leq n} \sum_{i \text{ with edge from } i \text{ to } k} (\text{pr}(i) / t_i) \\ &= \alpha + (1-\alpha) \sum_{1 \leq i \leq n} \sum_{k \text{ with edge from } i \text{ to } k} (\text{pr}(i) / t_i) * \\ &= \alpha + (1-\alpha) \sum_{i \text{ with edge from } i} \text{pr}(i) \end{aligned}$$

*inner sum \sum_i over incoming edges for one k



*inner sum \sum_k over outgoing edges for one i



6

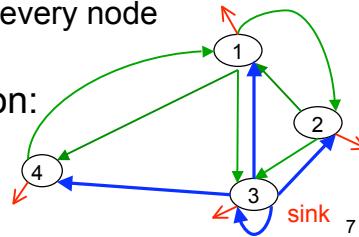
Problem for desired normalization

- Have

$$\sum_{1 \leq k \leq n} (\mathbf{pr}_{\text{new}}(\mathbf{k})) = \alpha + (1-\alpha) \sum_{i \text{ with edge from } i} \mathbf{pr}(i)$$
- **Missing $\mathbf{pr}(i)$** for nodes with no edges from them
 - sinks!
- **Solution:** add n edges out of every sink
 - Edge to every node including self
 - Gives $1/n$ contribution to every node

Gives desired normalization:

If $\sum_{1 \leq k \leq n} (\mathbf{pr}_{\text{initial}}(k)) = 1$
 then $\sum_{1 \leq k \leq n} (\mathbf{pr}(k)) = 1$



Matrix formulation

- Let E be the n by n adjacency matrix

$$E(i,k) = 1 \text{ if there is an edge from node } i \text{ to node } k$$

$$= 0 \text{ otherwise}$$
- Define **new matrix L** :
 - For each row i of E ($1 \leq i \leq n$)
 - If row i contains $t_i > 0$ ones, $L(i,k) = (1/t_i) E(i,k)$, $1 \leq k \leq n$
 - If row i contains 0 ones, $L(i,k) = 1/n$, $1 \leq k \leq n$
- Vector **\mathbf{pr}** of PageRank values defined by

$$\mathbf{pr} = (\alpha/n, \alpha/n, \dots, \alpha/n)^T + (1-\alpha) L^T \mathbf{pr}$$
- has a solution representing the **steady-state values $\mathbf{pr}(k)$**

8

Calculation

- Choose α
 - No single best value
 - Page and Brin originally used $\alpha=.15$
- Simple iterative calculation
 - Initialize $pr_{\text{initial}}(k) = 1/n$ for each node k
 - so $\sum_{1 \leq k \leq n} (pr_{\text{initial}}(k)) = 1$
 - $pr_{\text{new}}(k) = \alpha/n + (1-\alpha) \sum_{1 \leq i \leq n} L(i,k)pr(i)$
- Converges
 - Has necessary mathematical properties
 - In practice, choose convergence criterion
 - Stops iteration

9

PageRank Observations

- PageRank can be calculated for *any* graph
- Google calculates on entire Web graph
- Huge calculation for Web graph
 - precomputed
 - 1998 Google:
 - 52 iterations for 322 million links
 - 45 iterations for 161 million links
- PageRank must be combined with query-based scoring for final ranking
 - Many variations
 - What Google exactly does secret
 - Can make some guesses by results

10

HITS

Hyperlink Induced Topic Search

- Second well-known algorithm
- By Jon Kleinberg while at IBM Almaden Research Center
- Same general goal as PageRank
- Distinguishes **2 kinds of nodes**
 - **Hubs**: resource pages
 - **Point to many authorities**
 - **Authorities**: good information pages
 - **Point to many hubs**

11

Mutual reinforcement

- Authority weight node j : $a(j)$
 - Vector of weights \mathbf{a}
- Hub weight node j : $h(j)$
 - Vector of weights \mathbf{h}
- Update:

$$a_{\text{new}}(k) = \sum_{i \text{ with edge from } i \text{ to } k} (h(i))$$

$$h_{\text{new}}(k) = \sum_{j \text{ with edge from } k \text{ to } j} (a(j))$$

12

Matrix formulation

Steady state:

$$\mathbf{a} = \mathbf{E}^T \mathbf{h}$$

$$\mathbf{h} = \mathbf{E} \mathbf{a}$$

$$\mathbf{a} = \mathbf{E}^T \mathbf{E} \mathbf{a}$$

$$\mathbf{h} = \mathbf{E} \mathbf{E}^T \mathbf{h}$$

Interpretation:

- $\mathbf{E}^T \mathbf{E}(i,j)$: number nodes **point to** both node i and node j
 - “Co-citation”
- $\mathbf{E} \mathbf{E}^T(i,j)$: number nodes **pointed to by** both node i and node j
 - “Bibliographic coupling”

13

Iterative Calculation

$$\mathbf{a} = \mathbf{h} = (1, \dots, 1)^T$$

While (not converged) {

$$\mathbf{a}_{\text{new}} = \mathbf{E}^T \mathbf{h}$$

$$\mathbf{h}_{\text{new}} = \mathbf{E} \mathbf{a}$$

$$\mathbf{a} = \mathbf{a}_{\text{new}} / \|\mathbf{a}_{\text{new}}\| \quad \text{normalize to unit vector}$$

$$\mathbf{h} = \mathbf{h}_{\text{new}} / \|\mathbf{h}_{\text{new}}\| \quad \text{normalize to unit vector}$$

}

14

Convergence

- Linear algebra - eigenvalues
- Kleinberg uses slightly different iteration and slightly different proof than in *Intro to IR* book

– Normalization important

$$\mathbf{a}_0 = \mathbf{h}_0 = (1, \dots, 1)^T$$

For k^{th} iteration {

$$\mathbf{a}_k = \text{normalized } (E^T \mathbf{h}_{k-1})$$

$$\mathbf{h}_k = \text{normalized } (E \mathbf{a}_k) \quad \text{uses new value of } \mathbf{a}$$

}

$$\text{Then } \mathbf{a}_k = \text{normalized } ((E^T E)^{k-1} E^T \mathbf{a}_0)$$

$$\text{Then } \mathbf{h}_k = \text{normalized } ((E E^T)^k \mathbf{h}_0)$$

15

General Theorem:

If M is a symmetric n by n matrix and \mathbf{v} is a vector not orthogonal to the principal eigenvector \mathbf{w}_1 of M ,

then the unit vector in the direction of $M^k \mathbf{v}$ converges to \mathbf{w}_1 as k goes to infinity.

Application:

Since $\mathbf{h}_0 = (1, 1, \dots, 1)^T$, \mathbf{h}_0 is not orthogonal to the principal eigenvector of $E E^T$

$\Rightarrow \mathbf{h}_k$ converges

\mathbf{a}_k similar but little more work because first vector $E^T \mathbf{a}_0$

16

Use of HITS

- Actual use of HITS by IBM people was **after** find Web pages satisfying query:
 1. Retrieve documents satisfy query and **rank by term-based** techniques
 2. Keep **top c documents**: root set of nodes
 - c a chosen constant - tunable
 3. Make base set:
 1. Root set
 2. *Plus nodes pointed to by* nodes of **root set**
 3. *Plus nodes pointing to* nodes of **root set**
 4. Make base graph: base set plus edges from Web graph between these nodes
 5. Apply HITS to base graph

17

Results using HITS

- Documents ranked by authority score $a(\text{doc})$ and hub score $h(\text{doc})$
 1. Authority score primary score for search results
- Heuristics:
 - delete all links between pages in same domain
 - Keep only pre-determined number of pages linking into root set (~200)
- Findings (original paper)
 - Number iterations in original tests ~50
 - most authoritative pages **do not** contain initial query terms
 - Compare LSI “concepts”

18

Observations

- HITS can be applied to any graph
- Base graph **much smaller** than Web graph
- Kleinberg identified bad phenomena
 - Topic diffusion: generalizes topic when expand root graph to base graph
 - Want *compilers* - generalized to *programming*

19

HITS and clustering

- Non-principal eigenvectors of EE^T and E^TE have positive and negative component values
 - Denote a_{e_2}, a_{e_3}, \dots
matching h_{e_2}, h_{e_3}, \dots
- For a matched pair of eigenvectors \mathbf{a}_{e_j} and \mathbf{h}_{e_j}
 - Denote k^{th} component of j^{th} eigenvector $\mathbf{a}_{e_j}(k)$ and $\mathbf{h}_{e_j}(k)$
 - Make a “community” of size c (a chosen constant):
 - Choose c pages with most positive $\mathbf{h}_{e_j}(k)$ - hubs
 - Choose c pages with most positive $\mathbf{a}_{e_j}(k)$ - authorities
 - Make another “community” of size c :
 - Choose c pages with most negative $\mathbf{h}_{e_j}(k)$ - hubs
 - Choose c pages with most negative $\mathbf{a}_{e_j}(k)$ - authorities
- Compare LSI

20

Eigenvalues and clustering

General class of techniques for clustering a graph using eigenvectors of adjacency matrix (or similar matrix) called

Spectral clustering

First described in 1973

More later, maybe ...

21