### **Network Structure of the Web**

Chapter 21 (plus additional references)

#### The Web as a Graph

- Nodes: static Web pages
- Links: directed hyperlinks
- Notion of "in-links", "out-links", "in-degree", "out-degree".
- Questions:
  - What is the structure of this graph?
  - How did this structure come about?
  - How to best search, given this structure?
  - How to best crawl, given this structure?
  - How to mine the structure to find communities, themes?





















# The Science of Networks

Are there properties common to all complex networks?

If so, why?

#### **Observed common properties:**

- Small world property
- Scale-free structure
- Clustering and community structure
- Robustness to random node failure
- Vulnerability to targeted hub attacks
- Vulnerability to cascading failures





































### Other examples of power-laws in nature

- Magnitude vs. frequency of earthquakes
- Magnitude vs. frequency of stock market crashes
- Income vs. frequency (of people with that income)
- Populations of cities vs. frequency (of cities with that population)
- Word rank vs. frequency in English text

"More normal than 'normal'?

How are scale-free networks created?

Barabàsi and Albert: Preferential attachment

Netlogo demo













# Ergodicity in Markov Chains Markov chain is ergodic if there is some non-zero probability of visiting every state. Theorem: For any ergodic Markov chain with transition matrix P, there is a unique steady-state probability vector π that is the principal left eigenvector of P, such that if η(*i*,*i*) is the number of visits to state *i* in *t* steps, then $\lim_{t \to \infty} \frac{\eta(i, t)}{t} = \pi(i),$ where π(i) > 0 is the steady-state probability for state *i*. Corollary: The "random walk with teleporting" process described above results in a unique distribution of steady state probabilities over the states of the Markov chain.

#### **Definition**:

PageRank(i) =  $\pi(i)$ .







## Google and PageRank

- **Google:** "PageRank relies on the uniquely democratic nature of the web by using its vast link structure as an indicator of an individual page's value. In essence, Google interprets a link from page A to page B as a vote, by page A, for page B. But, Google looks at more than the sheer volume of votes, or links a page receives; it also analyzes the page that casts the vote. Votes cast by pages that are themselves 'important' weigh more heavily and help to make other pages 'important'."
- Google recalculates PageRank each time it crawls the web and rebuilds its index.
- PageRank tends to favor older pages (more links)
- http://directory.google.com/

From http://www.geek.com/articles/chips/googles-pagerank-algorithm-tracedback-to-the-1940s-20100217/

Earlier forerunner to PageRank in the work of the Harvard economist Wassily Leontief:

"In 1941, Leontief published a paper in which he divides a country's economy into sectors that both supply and receive resources from each other, although not in equal measure. One important question is: what is the value of each sector when they are so tightly integrated? Leontief's answer was to develop an iterative method of valuing each sector based on the importance of the sectors that supply it. Sound familiar? In 1973, Leontief was awarded the Nobel Prize in economics for this work."

#### Other Uses for Page Rank

- http://www.eigenfactor.org (for journal impact)
- Ranking doctoral programs (network: departments are nodes, one node links to another if it hires faculty from that dept.)
- Food webs species that are essential to an ecosystem

# Googling Food Webs: Can an Eigenvector Measure Species' Importance for Coextinctions?

#### Stefano Allesina<sup>1</sup>\*, Mercedes Pascual<sup>2,3,4</sup>

1 National Center for Ecological Analysis and Synthesis, Santa Barbara, California, United States of America, 2 Department of Ecology and Evolutionary Biology, Universi of Michigan, Ann Arbor, Michigan, United States of America, 3 Santa Fe Institute, Santa Fe, New Mexico, United States of America, 4 Howard Hughes Medical Institu

#### Abstract

A major challenge in ecology is forecasting the effects of species' extinctions, a pressing problem given current human impacts on the planet. Consequences of species losses such as secondary extinctions are difficult to forecast because species are not isolated, but interact instead in a complex network of ecological relationships. Because of their mutual dependence, the loss of a single species can cascade in multiple coextinctions. Here we show that an algorithm adapted from the one Google uses to rank web-pages can order species according to their importance for coextinctions, providing the sequence of losses that results in the fastest collapse of the network. Moreover, we use the algorithm to bridge the gap between qualitative (who eats whom) and quantitative (at what rate) descriptions of food webs. We show that our simple algorithm finds the best possible solution for the problem of assigning importance from the perspective of secondary extinctions in all analyzed networks. Our approach relies on network structure, but applies regardless of the specific dynamical model of species' interactions, because it identifies the subset of coextinctions common to all possible models, those that will happen with certainty given the complete loss of prey of a given predator. Results show that previous measures of importance based on the concept of "hubs" or number of connections, as well as centrality measures, do not identify the most effective extinction sequence. The proposed algorithm provides a basis for further developments in the analysis of extinction risk in ecosystems.

#### Googling Food Webs: Can an Eigenvector Measure Species' Importance for Coextinctions?

#### Stefano Allesina<sup>1</sup>\*, Mercedes Pascual<sup>2,3,4</sup>

1 National Center for Ecological Analysis and Synthesis, Santa Barbara, California, United States of America, 2 Department of Ecology and Evolutionary Biology, Universi of Michigan, Ann Arbor, Michigan, United States of America, 3 Santa Fe Institute, Santa Fe, New Mexico, United States of America, 4 Howard Hughes Medical Institu

#### Abstract

A major challenge in ecology is forecasting the effects of species' extinctions, a pressing problem given current human impacts on the planet. Consequences of species losses such as secondary extinctions are difficult to forecast because species are not isolated, but interact instead in a complex network of ecological relationships. Because of their mutual dependence, the loss of a single species can cascade in multiple coextinctions. Here we show that an algorithm adapted oviding

the gap simple

specific

nodels,

revious

do not s in the

the se between the second and the se

algorit one Google uses to rank web-pages can order species dynam according to their importance for coextinctions,

those identify analys fastest collapse of the network."



# Hubs and Authorities (HITS Algorithm)

- Proposed by Jon Kleinberg (Cornell) at same time Brin and Page were developing PageRank
- HITS: Hyperlinked-induced topic search
- Supposedly used by Teoma and Ask.com

#### Hubs and Authorities Main ideas

Each node has a hub score and an authority score

Hub: Web site that points to a lot of good authorities

Authority: Web site that is pointed to by a lot of good hubs

Circular definition  $\implies$  iterative computation

$$h(v) \leftarrow \sum_{v \mapsto y} a(y)$$
$$a(v) \leftarrow \sum_{y \mapsto v} h(y)$$

where  $a \mapsto b$  means a links to b

#### **Hubs and Authorities** How to compute scores

Let  $\mathbf{h}$  = vector of hub scores of web pages in collection Let  $\mathbf{a}$  = vector of authority scores of web pages in collection

Let A denote the adjacency matrix of the web pages in collection

Then:

 $\mathbf{h} \leftarrow \mathbf{A}\mathbf{A}^T\mathbf{h}$ h ← Aa  $\mathbf{a} \leftarrow \mathbf{A}^T \mathbf{A} \mathbf{a}$  $\mathbf{a} \leftarrow \mathbf{A}^T \mathbf{h}$ 

and

Can turn these into eigenvector equations:

$$\lambda_h \mathbf{h} = (\mathbf{A}\mathbf{A}^T)\mathbf{h}$$
$$\lambda_h \mathbf{a} = (\mathbf{A}^T\mathbf{A})\mathbf{a}$$

# **Hubs and Authorities** How to compute scores, continued • HITS Algorithm: - Assemble the target subset of web pages from a given query, form the graph induced by their hyperlinks, and compute $AA^{T}$ and $A^{T}A$ – Compute the principal eigenvectors of $AA^{T}$ and $A^{T}A$ to form the vector of hub scores **h** and the vector of authority scores a. - Output the top-scoring hubs and top-scoring authorities

#### Choosing the subset of the web (One proposed method, described in textbook)

- 1. Given a query , use a text index to get all pages containing the terms of the query. Call this the *root* set of pages.
- 2. Build the *base* set of pages to include the root set as well as any page that either links to a page in the root set, or is linked to by a page in the root set.

Use the base set to compute hub and authority scores.

#### Differences between PageRank and HITS

- PageRank computers one score per document. HITS computers two.
- PageRank executed at indexing time, HITS executed at query time. Hub and authority scores are query-specific, whereas PageRank scores are query-independent.
- PageRank is assigned to all documents in collection. HITS scores are assigned only to relevant subset.

# **In-Class Exercises**

Draw a picture of a Web structure containing nodes with high hub score and high authority score.

Draw two different pictures of Web structures containing nodes with high PageRank.