# Web Structure & Content
## (<u>not</u> a theory talk)

Filippo Menczer

The University of Iowa

**http://dollar.biz.uiowa.edu/~fil/**
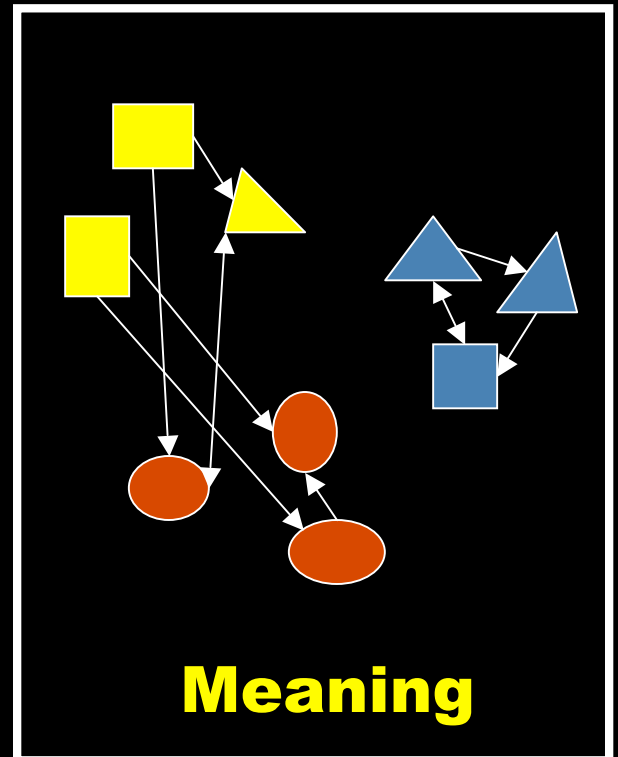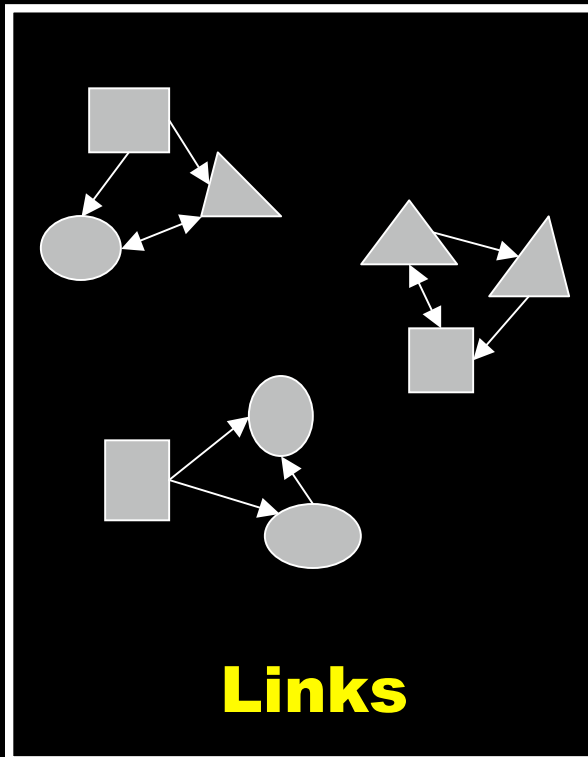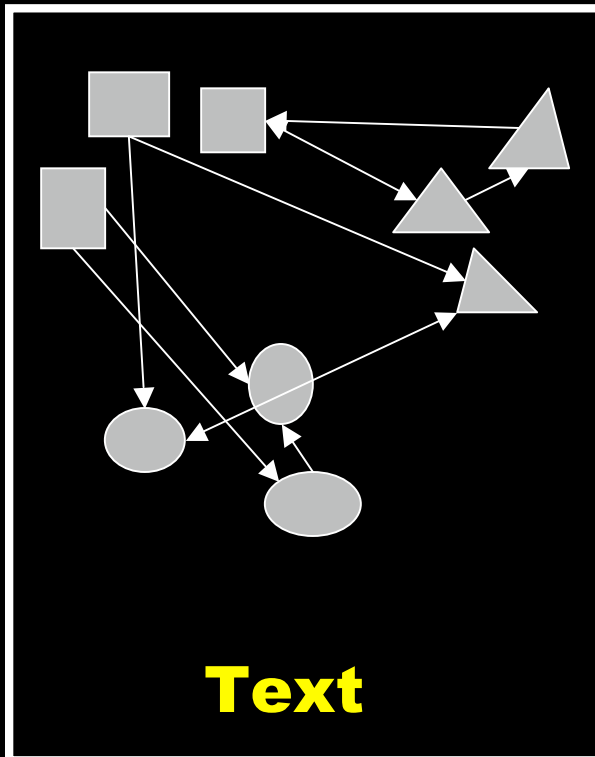
WAW 2002

# Exploiting the Web's text and link cues

- *Pages close in word vector space tend to be related*
  - Cluster hypothesis (van Rijsbergen 1979)
  - The WebCrawler (Pinkerton 1994)
  - The whole first generation of search engines

- *Pages that link to each other tend to be related*
  - Link-cluster conjecture (Menczer 1997)
    - Many formulations: Gibson & al 1998, Bharat & Henzinger 1998, Chakrabarti & al 1998, Dean & Henzinger 1999, Davison 2000, etc
  - Link eigenvalue analysis: HITS, hubs and authorities
    - (Kleinberg & al 1998 segg. @ Almaden etc.)
  - Google's PageRank analysis
    - (Brin & Page 1998)
  - The whole second generation of search engines

# Three topologies

What about the *relationship* between
lexical / link cues and page meaning?
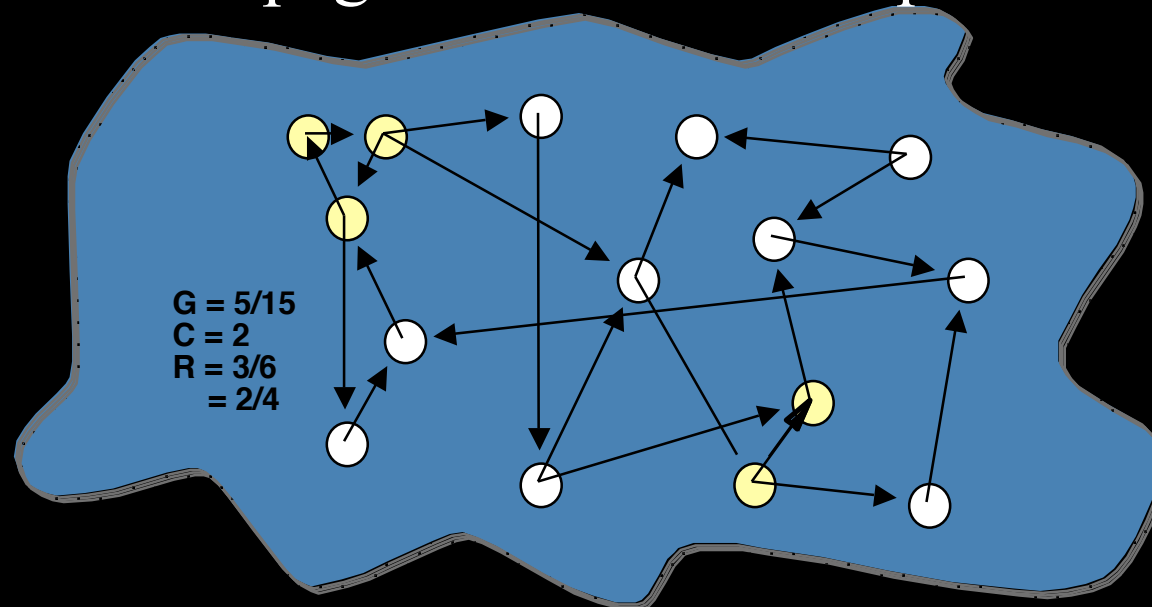


Text

Links

Meaning

# Talk outline

- The topologies of the Web
- Correlations, distributions, projections
- Power laws and Web growth models
- Navigating optimal paths
- Semantic maps (?)

# The "link-cluster" conjecture

- Connection between semantic topology (relevance) and linkage topology (hypertext)
  - $G = Pr[rel(p)] \sim$ fraction of relevant pages (generality)
  - $R = Pr[rel(p) \mid rel(q) \text{ AND } link(q,p)]$
- Related nodes are clustered if $R > G$
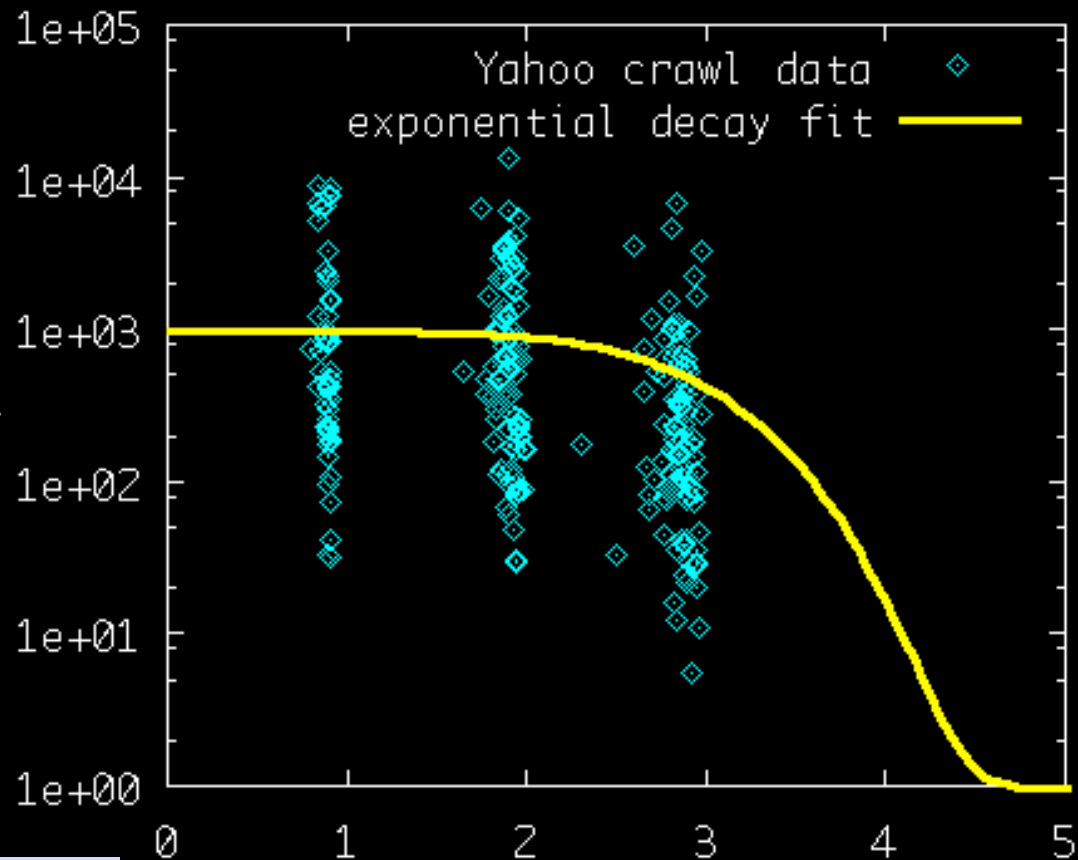  - Necessary and sufficient condition for a random crawler to find pages related to start points



G = 5/15
C = 2
R = 3/6
  = 2/4

$$\frac{R(q,\delta)}{G(q)} \equiv \frac{\Pr[rel(p)|rel(q) \land \|path(q,p)\| \le \delta]}{\Pr[rel(p)]}$$

# Link-cluster conjecture
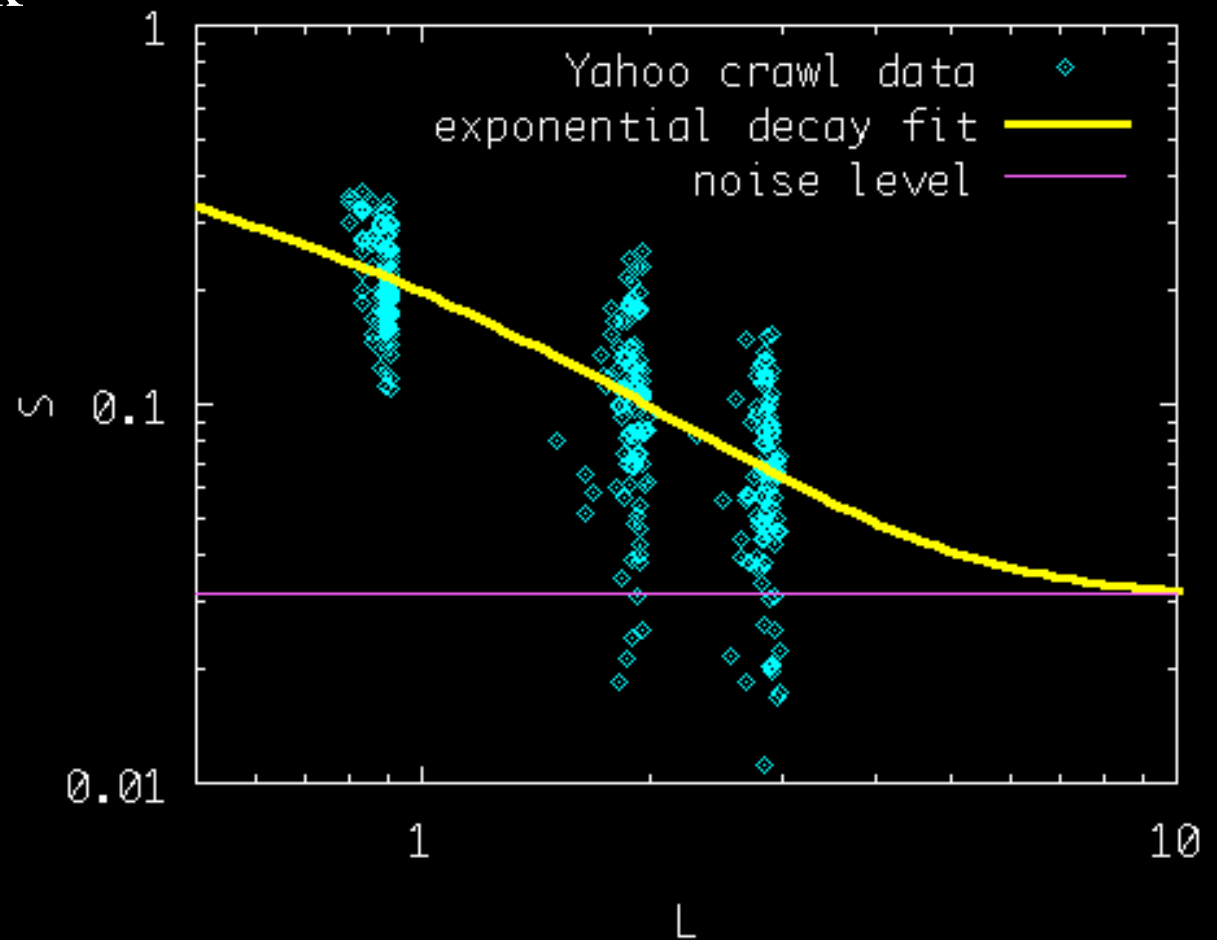
- Preservation of semantics (meaning) across links



$$L(q,\delta) \equiv \frac{\sum\limits_{\{p:\|path(q,p)\|\le\delta\}}\|path(q,p)\|}{\left|\{p:\|path(q,p)\| \le \delta\}\right|}$$

# The "link-content" conjecture
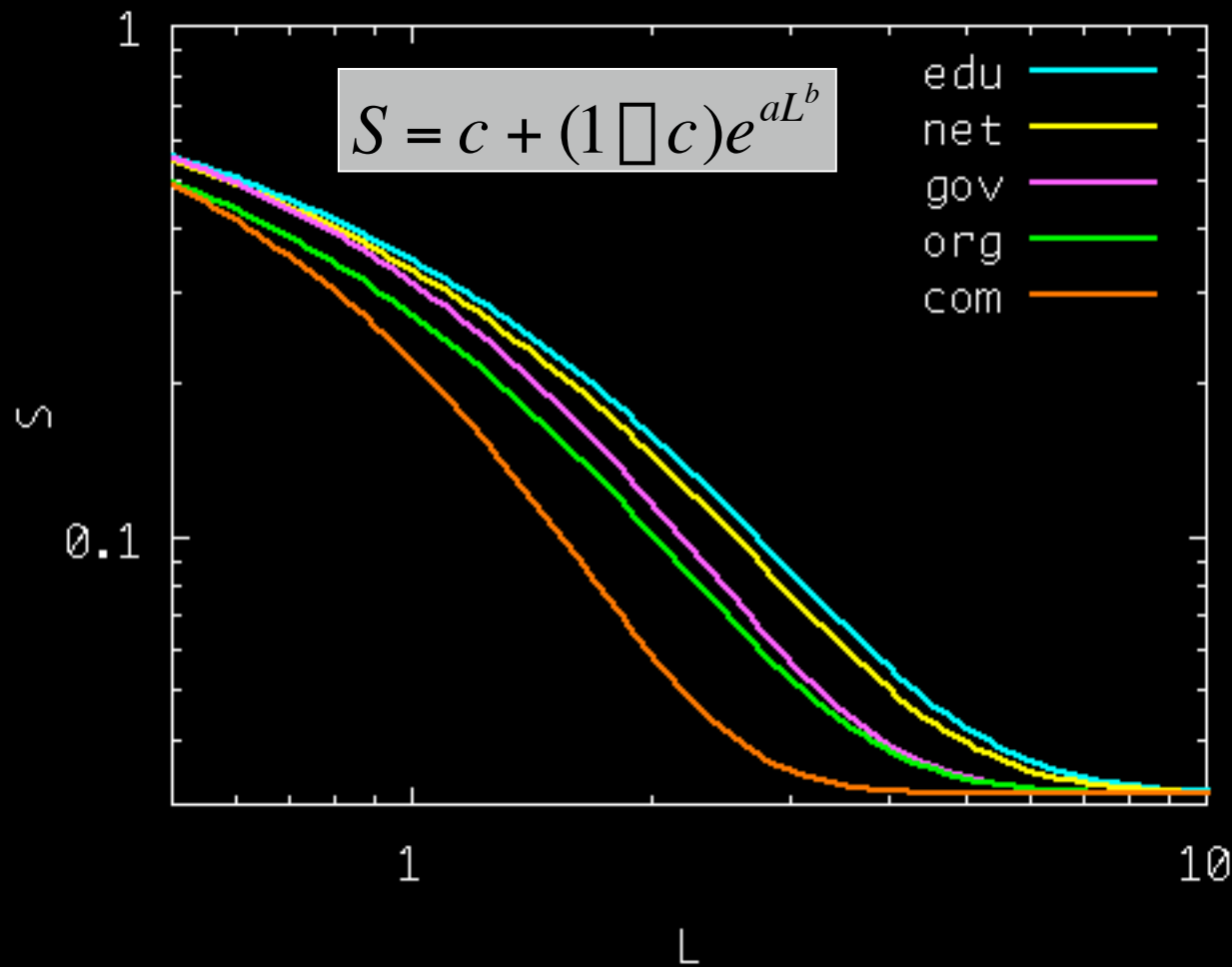
- Correlation of lexical and linkage topology

- **L(δ):** average link distance

- **S(δ):** average similarity to start (topic) page from pages up to distance δ

- Correlation $\rho(L,S) = -0.76$

$$S(q,\delta) \equiv \frac{\displaystyle\sum_{\{p:\|path(q,p)\|\leq\delta\}} sim(q,p)}{\left|\{p:\|path(q,p)\| \leq \delta\}\right|}$$

# Heterogeneity of lexical-linkage correlation



$$S = c + (1-c)e^{aL^b}$$

edu
net
gov
org
com

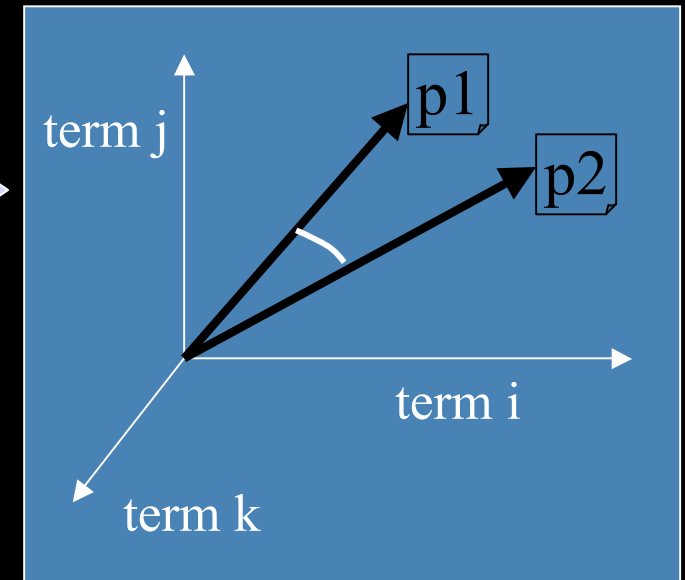signif. diff. *a* only ($\alpha$=0.05)

signif. diff. *a* & *b* ($\alpha$=0.05)

# Mapping the relationship between topologies
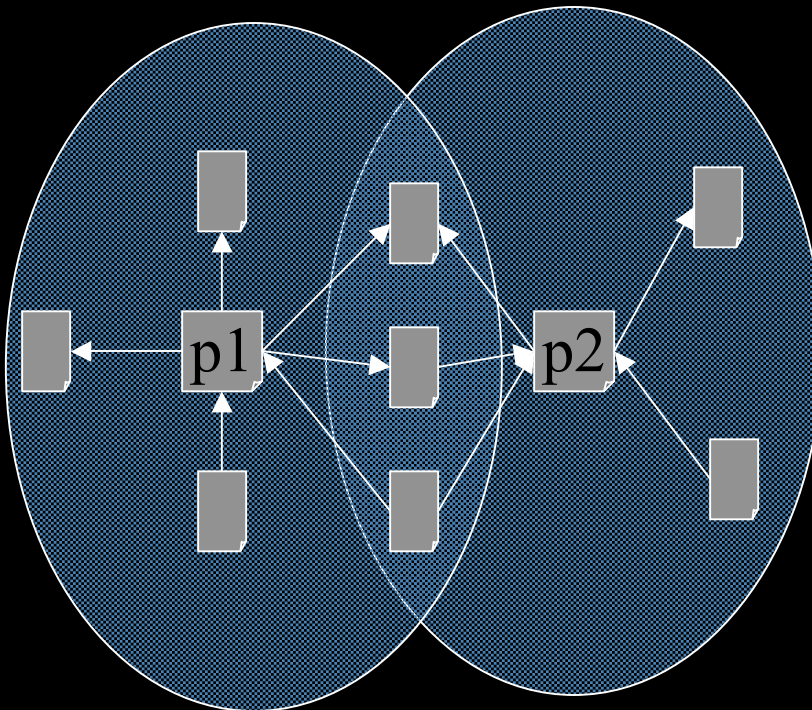
- Any pair of pages rather than linked pages from crawl
- Data: Open Directory Project (dmoz.org)
  - RDF Snapshot: 2002-02-14 04:01:50 GMT
  - After cleanup: 896,233 URLs in 97,614 topics
  - After sampling: 150,000 URLs in 47,174 topics
    - 10,000 from each of 15 top-level branches
- Need 'similarity' or 'proximity' metric for each topology, given a pair of pages:
  - Content: textual/lexical (cosine) similarity
  - Link: co-citation/bibliographic coupling
  - Semantic: relatedness inferred from manual classification

$$\sigma_c(p_1, p_2) = \frac{\sum_{k \in p_1 \cap p_2} f_{kp_1} f_{kp_2}}{\sqrt{\sum_{k \in p_1} f_{kp_1}^2 \sum_{k \in p_2} f_{kp_2}^2}}$$

Content similarity

term j

p1

p2

term i

term k

Link similarity

$$\sigma_l(p_1, p_2) = \frac{\left| U_{p_1} \cap U_{p_2} \right|}{\left| U_{p_1} \cup U_{p_2} \right|}$$

p1

p2

# Semantic similarity

$$\sigma_s(c_1, c_2) = \frac{2 \log \Pr[\mathrm{lca}(c_1, c_2)]}{\log \Pr[c_1] + \log \Pr[c_2]}$$

- Information-theoretic measure based on classification tree (Lin 1998)
- Classic path distance in special case of balanced tree

Correlations between similarities

3.84 x 10^9 pairs

# Joint distribution cube

# Web growth models

- Preferential attachment "BA"
  - At each step $t$ add page $p_t$
  - Create $m$ new links from $p_t$ to $p_{i<t}$

  (Barabasi & Albert 1999, de Solla Price 1976)

  $$\Pr(i) \propto k(i)$$

- Modified BA

  (Bianconi & Barabasi 2001, Adamic & Huberman 2000)

  $$\Pr(i) \propto \eta(i)k(i)$$

- Mixture

  (Pennock & al. 2002,
  Cooper & Frieze 2001, Dorogovtsev & al 2000)

  $$\Pr(i) \propto \psi \cdot k(i) + (1 - \psi) \cdot c$$

- Web copying

  (Kleinberg, Kumar & al 1999, 2000)

  $$\Pr(i) \propto \psi \cdot \Pr(j \rightarrow i) + (1 - \psi) \cdot c$$

- Mixture with Euclidean distance in graphs
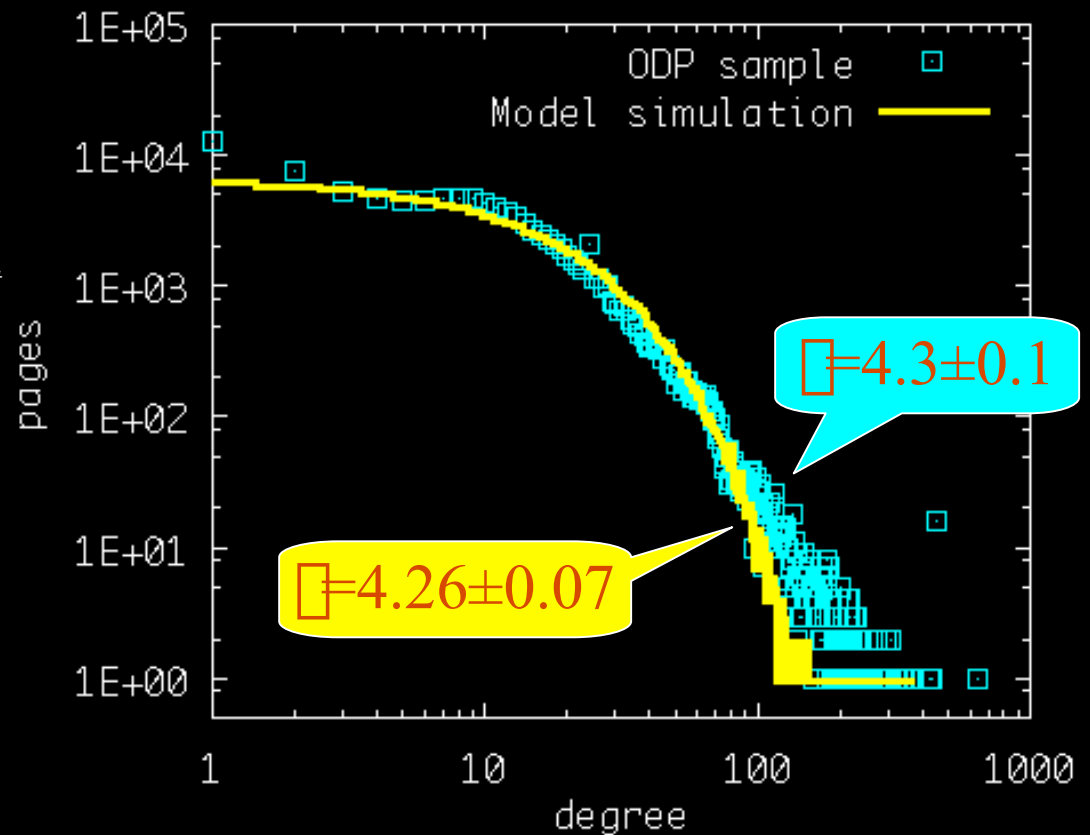
  (Fabrikant, Koutsoupias & Papadimitriou 2002)

  $$i = \arg\min(\phi r_{it} + g_i)$$

# Local content-based growth model

$$\Pr(p_t \to p_{i<t}) = \begin{cases} \dfrac{k(i)}{mt} & \text{if } r(p_i, p_t) < \rho^* \\ c[r(p_i, p_t)]^{-\alpha} & \text{otherwise} \end{cases}$$

- Similar to preferential attachment (BA)
- At each step $t$ add page $p_t$
- Create $m$ new links from $p_t$ to existing pages
- Use degree ($k$) info only for nearby pages
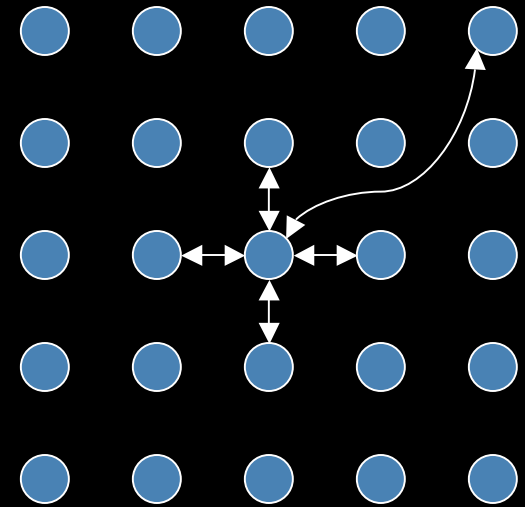  (popularity/importance of similar/related pages)

# Efficient crawling algorithms?

- <u>Theory</u>: since the Web is a small world network, or has a scale free degree distribution, there exist short paths between any two pages:
  - $\sim \log N$ (Barabasi & Albert 1999)
  - $\sim \log N / \log \log N$ (Bollobas 2001)
- <u>Practice</u>: can't find them!
  - Greedy algorithms based on location in geographical small world networks: $\sim$ poly(N) (Kleinberg 2000)
  - Greedy algorithms based on degree in power law networks: $\sim$ N (Adamic, Huberman *et al* 2001)

# Exception # 1

- Geographical networks (Kleinberg 2000)
  - Local links to all lattice neighbors
  - Long-range link probability distribution: power law $\text{Pr} \sim r^{-\alpha}$
    - $r$: lattice (Manhattan) distance
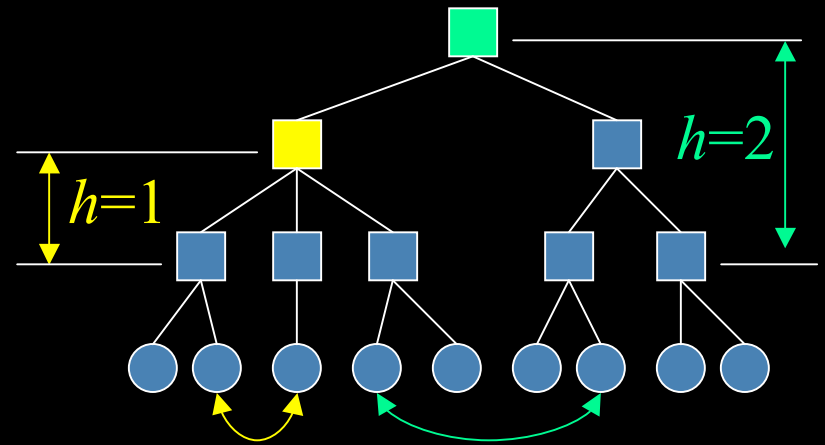    - $\alpha$: constant clustering exponent

$$t \sim \log^2 N \Leftrightarrow \alpha = D$$

# Exception # 2



- Hierarchical networks (Kleinberg 2002, Watts & *al.* 2002)
  - Nodes are classified at the leaves of tree
  - Link probability distribution: exponential tail $\Pr \sim e^{-h}$
    - $h$: tree distance (height of lowest common ancestor)
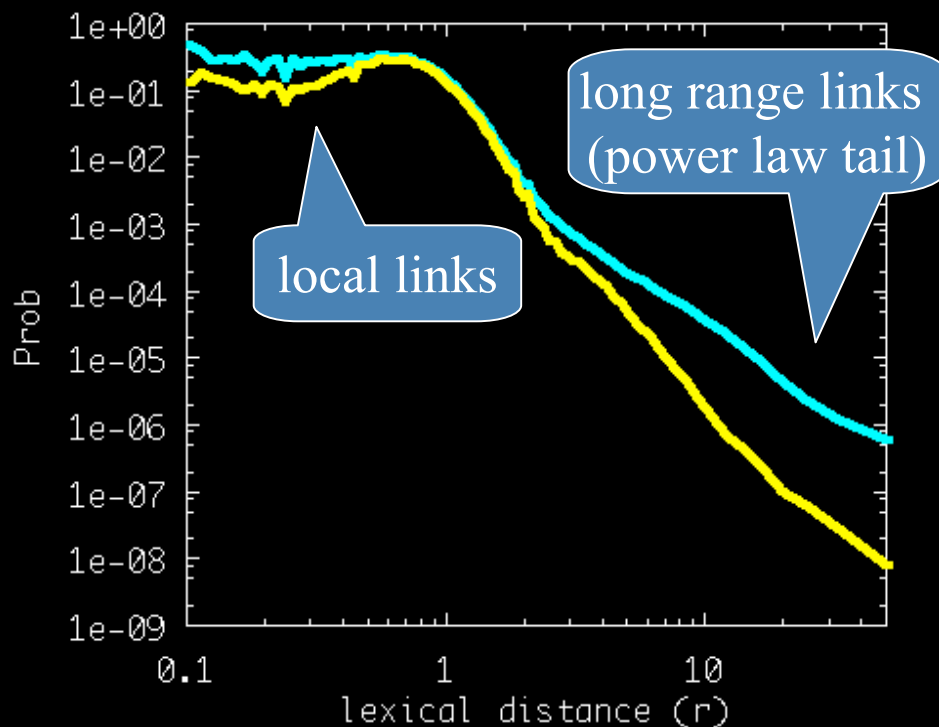
$$t \sim \log^{\varepsilon} N, \varepsilon \geq 1$$

# Is the Web one of these exceptions?

**Geographical model**
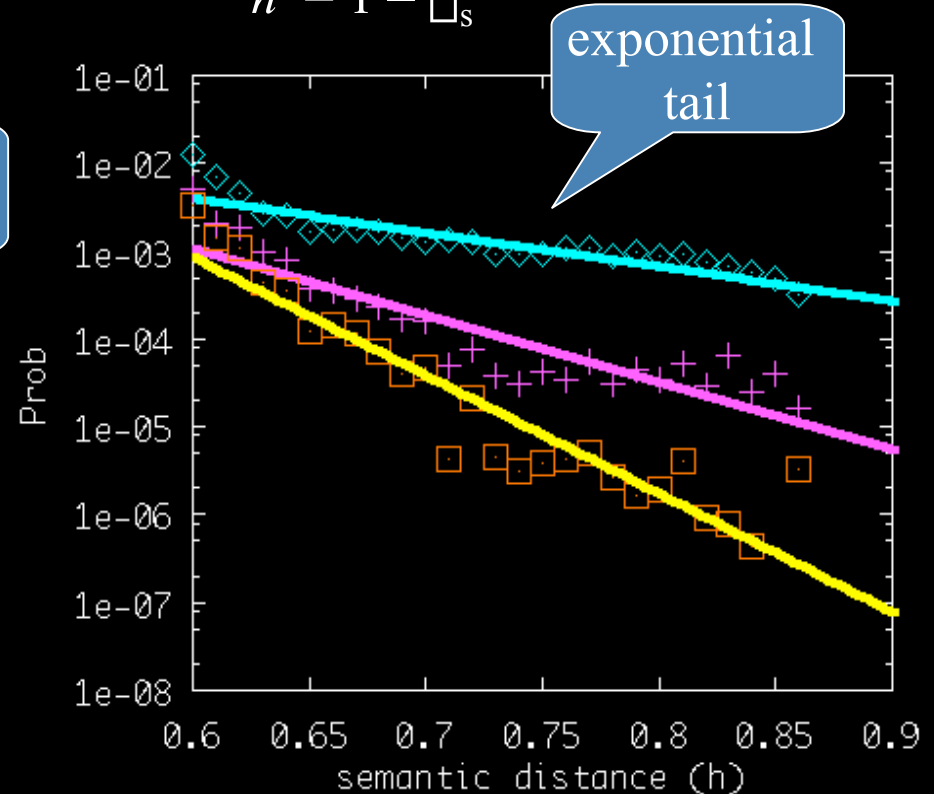- Replace lattice distance by lexical distance

$$r = (1 / \sigma_c) - 1$$

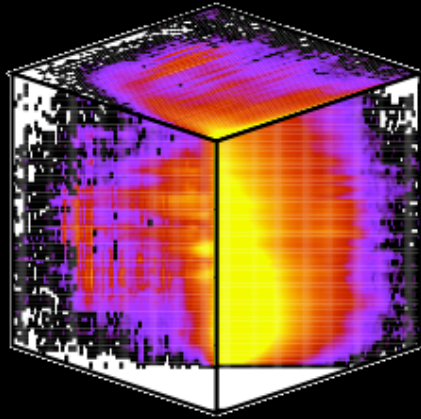**Hierarchical model**
- Replace tree distance by semantic distance

$$h = 1 - \sigma_s$$

# Talk outline

- The topologies of the Web
- Correlations, distributions, projections
- Power laws and Web growth models
- Navigating optimal paths
- **Semantic maps**

# Semantic maps: define "local" Precision and Recall

$$P(s_c, s_l) = \frac{\sum\limits_{\{p,q:\sigma_c = s_c, \sigma_l = s_l\}} \sigma_s(p,q)}{\left|\{p,q : \sigma_c = s_c, \sigma_l = s_l\}\right|}$$

Averaging semantic similarity

$$R(s_c, s_l) = \frac{\sum\limits_{\{p,q:\sigma_c = s_c, \sigma_l = s_l\}} \sigma_s(p,q)}{\sum\limits_{\{p,q\}} \sigma_s(p,q)}$$

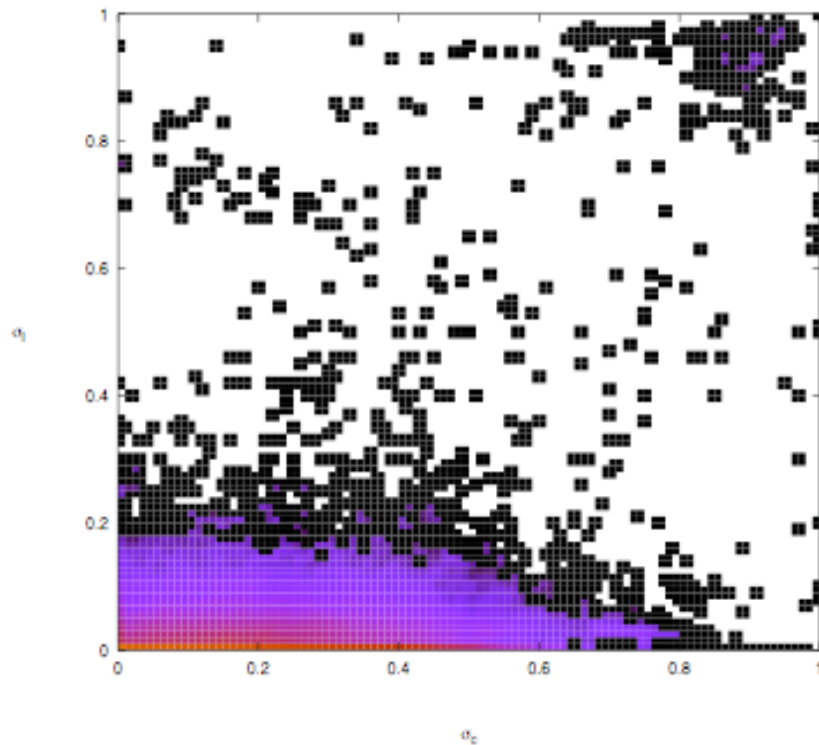Summing semantic similarity

# Semantic maps: *Business*
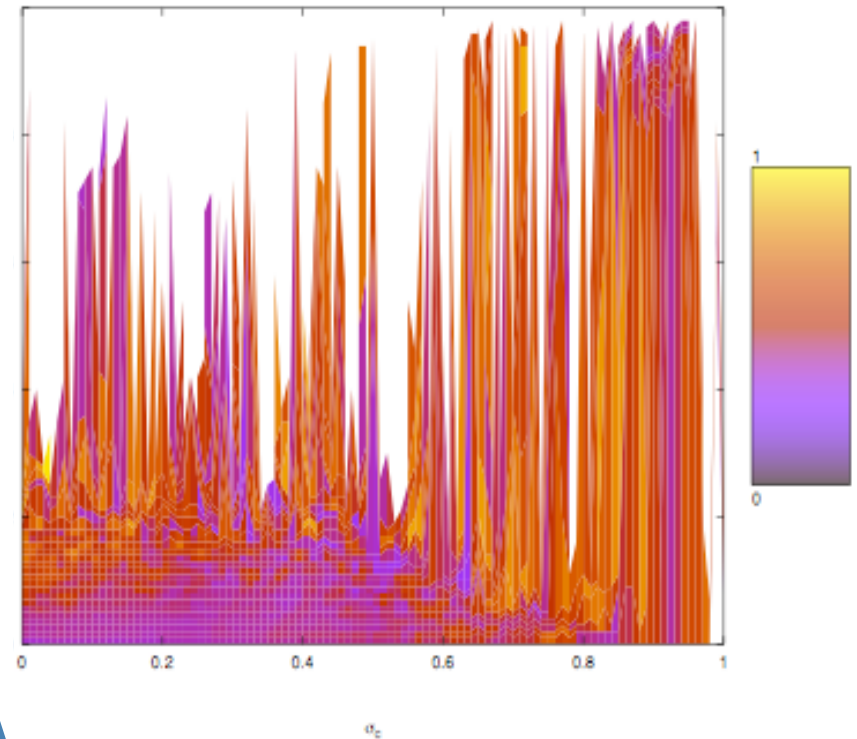


log Recall

Precision

$\sigma_\ell$

$\sigma_c$

# Semantic maps: *Adult*

log Recall

$\sigma_\ell$

Precision



$\sigma_c$

# Semantic maps: *Computers*



log Recall

Precision

$\sigma_\ell$

$\sigma_c$

# Semantic maps: *Home*

log Recall
$\sigma_\ell$
Precision

# Semantic maps: *News*



log Recall

$\sigma_\ell$

Precision

$\sigma_e$

# Semantic maps: all pairs



log Recall    Precision

$\sigma_\ell$

$\sigma_e$

# So what?

- Interpret performance of search engines
- Understand "ranking optimization"
  - vs filtering
  - vs combinations
- Topical signatures for topical/community portals
- Design "better" (and more scalable) crawlers
  - Topic driven
  - Query driven
  - User/community/peer driven
- Competitive intelligence, security applications

# Talk outline

- The topologies of the Web
- Correlations, distributions, projections
- Power laws and Web growth models
- Navigating optimal paths
- Semantic maps
- Questions?

**http://dollar.biz.uiowa.edu/~fil/**