How does the Web really grow?

Filippo Menczer

School of Informatics Department of Computer Science Indiana University

http://informatics.indiana.edu/fil/



Research supported by NSF CAREER Award IIS-0133124

Conclusion

To explain the emergent topology of the Web, we model the cognitive processes behind Web dynamics - content matters!

Characterizing the Web's link structure

- The Web is a small-world: small diameter

 (Barabasi, Albert & al @ ND,
 Huberman, Adamic & al @ PARC/HP)
- The Web is a bow tie

 (Broder & al @ AltaVista/Compag/IBM)
- The Web is scale-free: power laws (all of the above)
- Subsets of the Web are not scale free (Pennock, Lawrence, Giles & al @ NECI)





Web growth models classes

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)

Modified BA / fitness / hidden variable

(Bianconi & Barabasi 2001, Adamic & Huberman 2000...)

• Mixture

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

Web copying

(Kleinberg, Kumar & al 1999, 2000)

Mixture with metric distance

(Fabrikant, Koutsoupias & Papadimitriou 2002, Aldous 2003)

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

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

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

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

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

Mapping the relationship between topologies

- Given any pair of pages, need 'similarity' or 'proximity' metric for each topology:
 - Content: textual/lexical (cosine) similarity
 - Link: co-citation/bibliographic coupling
- 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

$$\sigma_c(\vec{p_1}, \vec{p_2}) = \overline{1}$$

$$\sigma_{c}\left(\overrightarrow{p_{1}},\overrightarrow{p_{2}}\right) = \frac{\overrightarrow{p_{1}}\cdot\overrightarrow{p_{2}}}{\left\|\overrightarrow{p_{1}}\right\|\cdot\left\|\overrightarrow{p_{2}}\right\|}$$

Content similarity





Link similarity

$$\sigma_{l}(p_{1}, p_{2}) = \frac{\left|U_{p_{1}} \cap U_{p_{2}}\right|}{\left|U_{p_{1}} \cup U_{p_{2}}\right|}$$

Þ

Individual metric distributions







Link probability vs lexical distance

$$r = 1/\sigma_c - 1$$
$$\Pr(\lambda \mid \rho) = \frac{|(p,q): r = \rho \land \sigma_l > \lambda|}{|(p,q): r = \rho|}$$



Local content-based growth model

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

- 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)



So, many models get the degree distribution right

- Which is "right"?
- Need an independent observation (other than degree) to validate models
- Distribution of content similarity across linked pairs
 - Across all pairs: $\Pr[\sigma_c] \sim 10^{-7\sigma_c}$

None of these models is right!



Back to the mixture model

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

- Pennock & al. 2002
- Cooper & Frieze 2001
- Dorogovtsev & al 2000
- Extra free parameter
 - Can tune popularity with "random"



 Bias choice by content similarity instead of uniform distribution

Mixture model class $Pr(i) = \psi \frac{k(i)}{2mt} + (1 - \psi)\hat{Pr}(i)$

Degree-uniform mixture:

$$\hat{\Pr}(i) = rac{1}{t}$$
 $\psi = 0.3$

Degree-similarity mixture:

$$\hat{\Pr}(i) \propto \left(rac{1}{\sigma_{m{c}}(i,t)}-1
ight)^{-lpha}$$

 $\psi=0.2, \ lpha=1.7$



Both mixture models get the degree distribution right...



...but the degree-similarity mixture model predicts the similarity distribution much better



Conclusion

- To explain the emergent topology of the Web, we model the cognitive processes behind Web dynamics - content matters!
- What about paper citation networks?





Questions?

http://informatics.indiana.edu/fil/



Research supported by NSF CAREER Award IIS-0133124