# SEO and spam

Dan Smith

**UEA**
University of East Anglia

---

# Overview

- Google
- Search engine optimisation
- Spam and adversarial IR

3

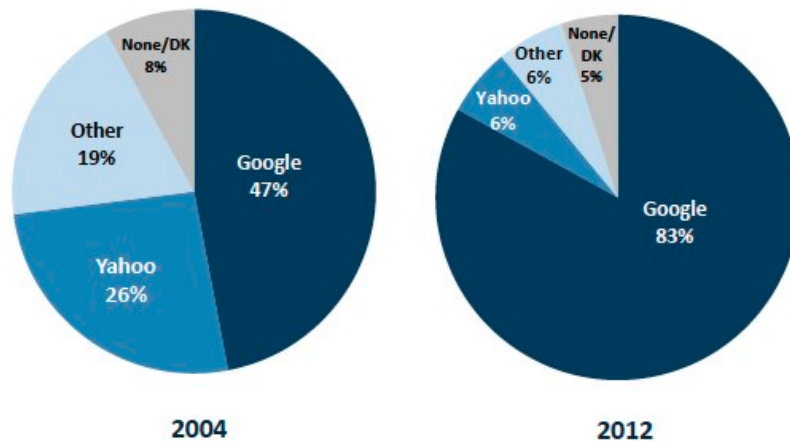**GOOGLE**

4

# Web search engines

- Google has approx. 90% of the world's web search traffic
  - UK 88%
  - USA 76%
- Other search engines
  - Bing the biggest competitor in USA, globally 3.37%
  - Baidu important in China, globally 0.79%
  - Yahoo! declining, globally 3.43%

http://www.statista.com/statistics/216573/worldwide-market-share-of-search-engines/

# The rise of Google

*% of search users who answered the question: Which search engine do you use MOST OFTEN?*

None/DK
8%

Other
19%

Google
47%

Yahoo
26%

2004

None/DK
5%

Other
6%

Yahoo
6%

Google
83%

2012

# Google now: Knowledge Graph

- Major development aimed to give better answers to conversational and question searches
- Based on a large graph of (mostly) named entities with links between them
- Graph is built from search history data
  - e.g. people who search for Impressionists look at pages about Claude Monet

# Google ranking factors



How to make search engines like your pages

# SEARCH ENGINE OPTIMISATION

8

# SEO objectives

- Make the site easy for search engines to
  - crawl
  - index
  - understand
- First, make the site easy for users to use...
  - they are the consumers
  - the search engine helps them find it
  - ... but this lecture's not about UX or web design

# SEO: what's important?

"Google has a lot of very clever people working on a lot of very advanced algorithms. Rather than trying to keep up and discover ways of cheating the system, we should spend more time looking at how a human sees our sites and engages with them, ..."

Gareth Owen, Oct. 2010,
  http://searchenginewatch.com/3641474

11

# Title and Meta tags

- Use brief, descriptive titles
  - each page (or small group of pages) should have a separate title
  - it's what people see in the search results
- Use the description meta tag to provide a concise summary of the page
  - it may be used as the snippet for the page

12

# URLs

- Crawlers use URLs to navigate the site
- Users use URLs to infer the site structure
- Other sites may use URLs as the displayed link text
- URLs are displayed in search results

- Search engines prefer shorter URLs with recognisable words
  - dislike long numeric or mixed identifiers, session identifiers, ...
  - use a single URL per page, as multiple URLs dilute the rank

13

# Navigation

- Important to understand how crawlers and users navigate the site
- Use a logical structure
- Users often truncate URLs to move to more general content
- A site map is a page showing the hierarchical structure of the site for humans
- An XML Sitemap is a structure to help search engines (Google)

14

# Links

- Search engines like
  - text links
  - text links with informative text
  - customised missing link (404) pages
- They dislike
  - broken links
  - uninformative or non-existent 404 messages
- Don't leak rank to spammers
  - add the rel="nofollow" attribute to links in comments, unfavourable reviews, …

15

# Content is paramount

- Good content gets linked to
- Stale content is ignored
    - keep updating the content
- Users will search using different vocabularies
    - according to experience, background, ...
    - search engines provide tools to help identify which terms are used to access a site, synonyms and variants for keywords, ...
- Create content for humans, not search engines

16

# Images

- Don't bury text in images
    - crawlers can't see it
- Always provide alt text
    - it helps crawlers (and humans)
- Use (short) descriptive filenames

# Page structure

- Use heading tags (h1-h6) appropriately
  - to reflect the structure and importance of text
- Use headings sparingly
  - avoid lots of very short sections

- Use the site's robots.txt file to keep crawlers away from content you don't want indexed
  - but not as a security measure

---

# SEO SPAMMING

19

# Search engine optimization (Spam)

- Motives
  - Commercial, political, religious, lobbies
  - Promotion funded by advertising budget
- Operators
  - Contractors (Search Engine Optimizers) for lobbies, companies
  - Web masters
  - Hosting services
- Forums
  - e.g., Webmaster world ( www.webmasterworld.com )
    - Search engine specific tricks
    - Discussions about academic papers ☺

20

# Cloaking

- Serve different content to crawlers and humans
  - sites suspected of cloaking get removed from search results
- DNS cloaking: Switch IP address

# The war against spam

- Quality signals – prefer pages based on:
  - Votes from authors (linkage signals)
  - Votes from users (usage signals)
- Policing of URL submissions
  - Anti robot test
- Limits on meta-keywords
- Robust link analysis
  - Ignore statistically implausible linkage (or text)
  - Use link analysis to detect spammers (guilt by association)

- Spam recognition
  - Training set based on known spam
- Family friendly filters
  - Linguistic analysis, general classification techniques, ...
  - For images: flesh tone detectors, source text analysis, ...
- Editorial intervention
  - Blacklists
  - Top queries audited
  - Complaints addressed
  - Suspect pattern detection

# More on spam

- Web search engines have policies on SEO practices they tolerate/block
  - http://help.yahoo.com/help/us/ysearch/index.html
  - http://www.google.com/intl/en/webmasters/
- Adversarial IR: the (technical) battle between SEOs and web search engines
- Research   http://airweb.cse.lehigh.edu/

23

# References

- Intro. to IR, Ch 19, 21
- Garfield E. (1972) Citation as a tool in journal evaluation, *Science* 178, 471–479
- Small H. (1973) Co-Citation in the Scientific Literature: A New Measure of the Relationship between Publications, *J. Am. Soc. Info. Sci.*
- Kleinberg J. (1998) Authoritative sources in a hyperlinked environment, *Proc. ACM-SIAM Symposium on Discrete Algorithms*
- Brin S. and Page L. (1998) The anatomy of a large-scale hypertextual Web-search engine, *Proc. 7th International World Wide Web Conference*