# Search Ads

1

# First generation of search ads: Goto (1996)



(Cost to advertiser: $0.38)

2

# First generation of search ads: Goto (1996)



- o Buddy Blake bid the maximum ($0.38) for this search.
- o He paid $0.38 to Goto every time somebody clicked on the link.
- o Pages were simply ranked according to bid – revenue maximization for Goto.
- o No separation of ads/docs. Only one result list!
- o Upfront and honest. No relevance ranking, . . .
- o . . . but Goto did not pretend there was any.

# Second generation of search ads: Google (2000/2001)

Strict separation of search results and search ads

4

# Two ranked lists: web pages (left) and ads (right)



SogoTrade appears in search results.

SogoTrade appears in ads.

Do search engines rank advertisers higher than non-advertisers?

All major search engines claim no.

5

# QUIZ: PAID RANKING

- Why is it not a good idea for Goto.com to show the amount successfully bid by the advertiser? (name just one good reason.)

# Do ads influence editorial content?

- Similar problem at newspapers / TV channels
- A newspaper is reluctant to publish harsh criticism of its major advertisers.
- The line often gets blurred at newspapers / on TV.
- No known case of this happening with search engines yet?

# How are the ads on the right ranked?

Google   discount broker          Search | Advanced Search
                                           Preferences

Web                              Results 1 - 10 of about 807,000 for discount broker [definition]. (0.12 seconds)

Sponsored Links

**Discount Broker** Reviews
Information on online **discount brokers** emphasizing rates, charges, and customer comments
and complaints.
www.**broker**-reviews.us/ - 94k - Cached - Similar pages

Rated #1 Online **Broker**
No Minimums. No Inactivity Fee
Transfer to Firstrade for Free!
www.firstrade.com

**Discount Broker** Rankings (2008 **Broker** Survey) at SmartMoney.com
**Discount Brokers**. Rank/ **Brokerage**/ Minimum to Open Account, Comments, Standard
Commis- sion*, Reduced Commission, Account Fee Per Year (How to Avoid), Avg. ...
www.smartmoney.com/**brokers**/index.cfm?story=2004-**discount**-table - 121k -
Cached - Similar pages

**Discount Broker**
Commission free trades for 30 days.
No maintenance fees. Sign up now.
TDAMERITRADE.com

Stock Brokers | Discount Brokers | Online Brokers
Most Recommended. Top 5 **Brokers** headlines. 10. Don't Pay Your **Broker** for Free Funds
May 15 at 3:39 PM. 5. Don't **Discount** the Discounters Apr 18 at 2:41 PM ...
www.fool.com/investing/**brokers**/index.aspx - 44k - Cached - Similar pages

TradeKing - Online **Broker**
$4.95 per Trade, Market or Limit
SmartMoney Top **Discount Broker** 2007
www.TradeKing.com

**Discount Broker**
**Discount Broker** - Definition of **Discount Broker** on Investopedia - A stockbroker who carries
out buy and sell orders at a reduced commission compared to a ...
www.investopedia.com/terms/d/**discountbroker**.asp - 31k - Cached - Similar pages

Scottrade Brokerage
$7 Trades, No Share Limit. In-Depth
Research. Start Trading Online Now!
www.Scottrade.com

Discount Brokerage and Online Trading for Smart Stock Market ...
Online stock **broker** SogoTrade offers the best in **discount brokerage** investing. Get stock
market quotes from this internet stock trading company.
www.sogotrade.com/ - 39k - Cached - Similar pages

Stock trades $1.50 - $3
100 free trades, up to $100 back
for transfer costs, $500 minimum
www.sogotrade.com

15 questions to ask **discount brokers** - MSN Money
Jan 11, 2004 ... If you're not big on hand-holding when it comes to investing, a **discount**
**broker** can be an economical way to go. Just be sure to ask these ...
moneycentral.msn.com/content/Investing/Startinvesting/P66171.asp - 34k -
Cached - Similar pages

$3.95 Online Stock Trades
Market/Limit Orders, No Share Limit
and No Inactivity Fees
www.Marsco.com

INGDIRECT | ShareBuilder

8

# How are ads ranked?

- Advertisers bid for keywords – sale by auction.
- Open system: Anybody can participate and bid on keywords.
- Advertisers are only charged when somebody clicks on your ad.
- How does the auction determine an ad's rank and the price paid for the ad?
- Basis is a second price auction, but with twists
- For the bottom line, this is perhaps the most important research area for search engines – computational advertising.
  - Squeezing an additional fraction of a cent from each ad means billions in additional revenue for the search engine.

9

# How are ads ranked?

- First cut: according to bid price **only** `a la Goto
  - Bad idea: open to abuse
  - Example: query [does my wife cheat?] → ad for divorce lawyer
  - We don't want to show nonrelevant ads.
- Instead: rank based on bid price and relevance
- Key measure of ad relevance: clickthrough rate
  - clickthrough rate = CTR = clicks per impressions
- Result: A nonrelevant ad will be ranked low.
  - Even if this decreases search engine revenue short-term
  - Hope: Overall acceptance of the system and overall revenue is maximized if users get useful information.
- Other ranking factors: location, time of day, quality and loading speed of landing page
- The main ranking factor: the query

# Google AdsWords demo

11

# Google's second price auction

| advertiser | bid | CTR | ad rank | rank | paid |
|---|---|---|---|---|---|
| A | $4.00 | 0.01 | 0.04 | 4 | (minimum) |
| B | $3.00 | 0.03 | 0.09 | 2 | $2.68 |
| C | $2.00 | 0.06 | 0.12 | 1 | $1.51 |
| D | $1.00 | 0.08 | 0.08 | 3 | $0.51 |

- bid: maximum bid for a click by advertiser
- CTR: click-through rate: when an ad is displayed, what percentage of time do users click on it? CTR is a measure of relevance.
- ad rank: bid × CTR: this trades off (i) how much money the advertiser is willing to pay against (ii) how relevant the ad is
- rank: rank in auction
- paid: second price auction price paid by advertiser

# Google's second price auction

| advertiser | bid | CTR | ad rank | rank | paid |
|---|---|---|---|---|---|
| A | $4.00 | 0.01 | 0.04 | 4 | (minimum) |
| B | $3.00 | 0.03 | 0.09 | 2 | $2.68 |
| C | $2.00 | 0.06 | 0.12 | 1 | $1.51 |
| D | $1.00 | 0.08 | 0.08 | 3 | $0.51 |

Second price auction: The advertiser pays the minimum amount necessary to maintain their position in the auction (plus 1 cent).

$price_1 \times CTR_1 = bid_2 \times CTR_2$ (this will result in $rank_1 = rank_2$)

$price_1 = bid_2 \times CTR_2 / CTR_1$

$p_1 = bid_2 \times CTR_2/CTR_1 = 3.00 \times 0.03/0.06 = 1.50$
$p_2 = bid_3 \times CTR_3/CTR_2 = 1.00 \times 0.08/0.03 = 2.67$
$p_3 = bid_4 \times CTR_4/CTR_3 = 4.00 \times 0.01/0.08 = 0.50$

Notice 2nd guy pays more than 1st guy

13

# Keywords with high bids

According to http://www.cwire.org/highest-paying-search-terms/

$69.1   mesothelioma treatment options

$65.9   personal injury lawyer michigan

$62.6   student loans consolidation

$61.4   car accident attorney los angeles

$59.4   online car insurance quotes

$59.4   arizona dui lawyer

$46.4   asbestos cancer

$40.1   home equity line of credit

$39.8   life insurance quotes

$39.2   refinancing

$38.7   equity line of credit

$38.0   lasik eye surgery new york city

$37.0   2nd mortgage

$35.9   free car insurance quote

14

# Search ads: A win-win-win?

- The search engine company gets revenue every time somebody clicks on an ad.

- The user only clicks on an ad if they are interested in the ad.

    - Search engines punish misleading and nonrelevant ads.

    - As a result, users are often satisfied with what they find after clicking on an ad.

- The advertiser finds new customers in a cost-effective way.

15

# Quiz: Search Ads

- Why is web search potentially more attractive for advertisers than TV spots, newspaper ads or radio spots? (name just one reason.)

16

# Not a win-win-win: Keyword arbitrage

- Buy a keyword on Google
- Then redirect traffic to a third party that is paying much more than you are paying Google.
  - E.g., redirect to a page full of ads
- This rarely makes sense for the user.
- Ad spammers keep inventing new tricks.
- The search engines need time to catch up with them.

17

# Not a win-win-win: Violation of trademarks

- Example: geico
- During part of 2005: The search term "geico" on Google was bought by competitors.
- Geico lost this case in the United States.
- Louis Vuitton lost similar case in Europe.
- See https://www.cnet.com/news/geico-sues-google-overture-over-trademarks/
- It's potentially misleading to users to trigger an ad of a trademark if the user can't buy the product on the site.

18

# SPAM
## (SEARCH ENGINE OPTIMIZATION)

19

# THE TROUBLE WITH PAID SEARCH ADS …

- It costs money.  What's the alternative?
- *Search Engine Optimization:*
  - "Tuning" your web page to rank highly in the algorithmic search results for select keywords
  - Alternative to paying for placement
  - Thus, intrinsically a marketing function
- Performed by companies, webmasters and consultants ("Search engine optimizers") for their clients
- Some perfectly legitimate, some very shady

20

# SEARCH ENGINE OPTIMIZATION (SPAM)

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

# Simplest forms

- First generation engines relied heavily on *tf/idf*
  - The top-ranked pages for the query `maui resort` were the ones containing the most `maui'`s and `resort'`s
- SEOs responded with dense repetitions of chosen terms
  - e.g., `maui resort maui resort maui resort`
  - Often, the repetitions would be in the same color as the background of the web page
    - Repeated terms got indexed by crawlers
    - But not visible to humans on browsers

Pure word density cannot be trusted as an IR signal

# VARIANTS OF KEYWORD STUFFING

- Misleading meta-tags, excessive repetition
- Hidden text with colors, style sheet tricks, etc.

**Meta-Tags** =
"… London hotels, hotel, holiday inn, hilton, discount, booking, reservation, sex, mp3, britney spears, viagra, …"

23

# CLOAKING

- Serve fake content to search engine spider
- DNS cloaking: Switch IP address, impersonate.

# More spam techniques

- **Doorway pages**
  - Pages optimized for a single keyword that re-direct to the real target page
- **Link spamming**
  - Mutual admiration societies, hidden links, awards – more on these later
  - *Domain flooding:* numerous domains that point or re-direct to a target page
- **Robots**
  - Fake query stream – rank checking programs
    - "Curve-fit" ranking programs of search engines
  - Millions of submissions via Add-Url

# THE WAR AGAINST SPAM

- Quality signals - Prefer authoritative 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 by machine learning
  - Training set based on known spam
- Family friendly filters
  - Linguistic analysis, general classification techniques, etc.
  - For images: flesh tone detectors, source text analysis, etc.
- 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
  - https://www.bing.com/toolbox/webmaster/
  - http://www.google.com/intl/en/webmasters/
- Adversarial IR: the unending (technical) battle between SEO's and web search engines
- Research  http://airweb.cse.lehigh.edu/

# SIZE OF THE WEB

28

# WHAT IS THE SIZE OF THE WEB ?

- Issues
  - The web is really infinite
    - Dynamic content, e.g., calendars
    - Soft 404: www.yahoo.com/<anything> is a valid page
  - Static web contains syntactic duplication, mostly due to mirroring (~30%)
  - Some servers are seldom connected
- Who cares?
  - Media, and consequently the user
  - Engine design
  - Engine crawl policy. Impact on recall.

29

# WHAT CAN WE ATTEMPT TO MEASURE?

- The relative sizes of search engines
  - The notion of a page being indexed is still *reasonably* well defined.
  - Already there are problems
    - Document extension: e.g., engines index pages not yet crawled, by indexing anchor text.
    - Document restriction: All engines restrict what is indexed (first *n* words, only relevant words, etc.)

30

Sec. 10.5

# NEW DEFINITION?

- The statically indexable web is whatever search engines index.
  - IQ is whatever the IQ tests measure.
- Different engines have different preferences
  - max url depth, max count/host, anti-spam rules, priority rules, etc.
- Different engines index different things under the same URL:
  - frames, meta-keywords, document restrictions, document extensions, …

31

# Relative Size from Overlap Given two engines A and B

**Sample** URLs randomly from A

**Check** if contained in B and vice versa

```
A ∩ B =  (1/2) * Size A
A ∩ B =  (1/6) * Size B

(1/2)*Size A = (1/6)*Size B

∴ Size A / Size B =
            (1/6)/(1/2) = 1/3
```

32

**Each test involves:** (i) <u>Sampling</u> (ii) Checking

# Sampling URLs

- Ideal strategy: Generate a random URL and check for containment in each index.

-  Problem: Random URLs are hard to find! Enough to generate a random URL contained in a given Engine.

- Approach 1: Pick a random URL contained in a given engine
  - Suffices for the estimation of relative size
- Approach 2: Random walks / IP addresses
  - In theory: might give us a true estimate of the size of the web (as opposed to just relative sizes of indexes)

33

# STATISTICAL METHODS

- Approach 1
  - Random queries
  - Random searches
- Approach 2
  - Random IP addresses
  - Random walks

34

# RANDOM URLS FROM RANDOM QUERIES

- Generate <u>random query</u>: how?

  Not an English dictionary

  - **Lexicon:** 400,000+ words from a web crawl

  - **Conjunctive Queries:** $w_1$ and $w_2$

    *e.g., vocalists AND rsi*

- Get 100 result URLs from engine A

- Choose a random URL as the candidate to check for presence in engine B

- This distribution induces a probability weight W(p) for each page.

35

# QUERY BASED CHECKING

- *Strong Query* to check whether an engine *B* has a document *D*:
  - Download *D*. Get list of words.
  - Use 8 low frequency words as AND query to *B*
  - Check if *D* is present in result set.
- Problems:
  - Near duplicates
  - Frames
  - Redirects (to docs not on engine *B*)
  - Engine time-outs
  - Is 8-word query good enough?

36

# ADVANTAGES & DISADVANTAGES

- Statistically sound under the "induced weight".
- Biases induced by random query
  - Query Bias: Favors content-rich pages in the language(s) of the lexicon
  - Ranking Bias: *Solution:* Use conjunctive queries & fetch all
  - Checking Bias: Duplicates, impoverished pages omitted
  - Document or query restriction bias: engine might not deal properly with 8 words conjunctive query
  - Malicious Bias: Sabotage by engine
  - Operational Problems: Time-outs, failures, engine inconsistencies, index modification.

37

# Random searches

- Choose random searches extracted from a local log [Lawrence & Giles 97] or build "random searches" [Notess]
  - Use only queries with small result sets.
  - Count normalized URLs in result sets.
  - Use ratio statistics

38

# ADVANTAGES & DISADVANTAGES

- Advantage
  - Might be a better reflection of the human perception of coverage (because it covers all the human searches)
- Issues
  - Samples are correlated with source of log
  - Duplicates
  - Technical statistical problems (must have non-zero results, ratio average not statistically sound)

# RANDOM SEARCHES

- 575 & 1050 queries from the NEC RI employee logs
- 6 Engines in 1998, 11 in 1999
- Implementation:
  - Restricted to queries with < 600 results in total
  - Counted URLs from each engine after verifying query match
  - Computed size ratio & overlap for individual queries
  - Estimated index size ratio & overlap by averaging over all queries

40

# QUIZ: QUERIES FROM NEC STUDY

- *adaptive access control*
- *neighborhood preservation topographic*
- *hamiltonian structures*
- *right linear grammar*
- *pulse width modulation neural*
- *unbalanced prior probabilities*
- *ranked assignment method*
- *internet explorer favourites importing*
- *karvel thornber*
- *zili liu*

- *softmax activation function*
- *bose multidimensional system theory*
- *gamma mlp*
- *dvi2pdf*
- *john oliensis*
- *rieke spikes exploring neural*
- *video watermarking*
- *counterpropagation network*
- *fat shattering dimension*
- *abelson amorphous computing*

41

## What's the problem with these queries?

# Random IP addresses

- Generate random IP addresses
- Find a web server at the given address
  - If there's one
- Collect all pages from server
  - From this, choose a page at random

# RANDOM IP ADDRESSES

- HTTP requests to random IP addresses
  - Ignored: empty or authorization required or excluded
  - [Lawr99] Estimated 2.8 million IP addresses running crawlable web servers (16 million total) from observing 2500 servers.
  - OCLC using IP sampling found 8.7 M hosts in 2001
    - Netcraft [Netc02] accessed 37.2 million hosts in July 2002
- [Lawr99] exhaustively crawled 2500 servers and extrapolated
  - Estimated size of the web to be 800 million pages
  - Estimated use of metadata descriptors:
    - Meta tags (keywords, description) in 34% of home pages, Dublin core metadata in 0.3%

43

# ADVANTAGES & DISADVANTAGES

- Advantages
  - Clean statistics
  - Independent of crawling strategies
- Disadvantages
  - Doesn't deal with duplication
  - Many hosts might share one IP, or not accept requests
  - No guarantee all pages are linked to root page.
    - E.g.: employee home pages
  - Power law for # pages/hosts generates bias towards sites with few pages.
    - But bias can be accurately quantified IF underlying distribution understood
  - Potentially influenced by spamming (multiple IP's for same server to avoid IP block)

44

# RANDOM WALKS

- View the Web as a directed graph
- Build a random walk on this graph
  - Includes various "jump" rules back to visited sites
    - Does not get stuck in spider traps!
    - Can follow all links!
  - Converges to a stationary distribution
    - Must assume graph is finite and independent of the walk.
    - Conditions are not satisfied (cookie crumbs, flooding)
    - Time to convergence not really known
  - Sample from stationary distribution of walk
  - Use the "strong query" method to check coverage by search engine

# ADVANTAGES & DISADVANTAGES

- Advantages
  - "Statistically clean" method, at least in theory!
  - Could work even for infinite web (assuming convergence) under certain metrics.
- Disadvantages
  - List of seeds is a problem.
  - Practical approximation might not be valid.
  - Non-uniform distribution
    - Subject to link spamming

46

# CONCLUSIONS

- No sampling solution is perfect.
- Lots of new ideas ...
- ....but the problem is getting harder
- Quantitative studies are fascinating and a good research problem

# DUPLICATE DETECTION

48

# DUPLICATE DOCUMENTS

- The web is full of duplicated content
- Strict duplicate detection = exact match
  - Not as common
- But many, many cases of near duplicates
  - E.g., last-modified date the only difference between two copies of a page

# DUPLICATE/NEAR-DUPLICATE DETECTION

- *Duplication*: Exact match  can be detected with fingerprints

- *Near-Duplication*: Approximate match

  - Overview
    - Compute syntactic similarity with an edit-distance measure
    - Use similarity threshold to detect near-duplicates
      - E.g.,  Similarity > 80% => Documents are "near duplicates"
      - Not transitive though sometimes used transitively

50

# COMPUTING SIMILARITY

- Features:
  - Segments of a document (natural or artificial breakpoints)
  - **Shingles** (Word N-Grams)
  - ***a rose is a rose is a rose*** →

    a_rose_is_a

    rose_is_a_rose

    is_a_rose_is

    a_rose_is_a

- Similarity Measure between two docs (= <u>sets of shingles</u>)
  - Jaccard coefficient: Size_of_Intersection / Size_of_Union

# SHINGLES + SET INTERSECTION

- Computing <u>exact</u> set intersection of shingles between <u>all</u> pairs of documents is expensive/intractable
  - Approximate using a cleverly chosen subset of shingles from each (a *sketch*)
- Estimate (size_of_intersection / size_of_union) based on a short sketch



Doc A → Shingle set A → Sketch A

Doc B → Shingle set B → Sketch B

Jaccard

52

# SKETCH OF A DOCUMENT

- Create a "sketch vector" (of size ~200) for each document

  - Documents that share $\geq t$ (say 80%) corresponding vector elements are near duplicates

  - For doc $D$, sketch$_D[\,i\,]$ is as follows:
    - Let f map all shingles in the universe to $0..2^m{-}1$ (e.g., f = fingerprinting)
    - Let $\pi_i$ be a *random permutation* on $0..2^m{-}1$
    - Pick MIN $\{\pi_i(f(s))\}$ over all shingles $s$ in $D$

53

# COMPUTING SKETCH[I] FOR DOC1

**Document 1**



$2^{64}$  **Start with 64-bit $f$(shingles)**

$2^{64}$  **Permute on the number line with $\pi_i$**

$2^{64}$

$2^{64}$  **Pick the min value**

# TEST IF DOC1.SKETCH[I] = DOC2.SKETCH[I]



**Document 1**

**Document 2**

Apply the same perm. On Doc 2.

$2^{64}$

$2^{64}$

$2^{64}$

$2^{64}$

$2^{64}$

$2^{64}$

$2^{64}$

$2^{64}$

A

B

Are these equal?

Test for **200** random permutations: $\pi_1, \pi_2, \ldots \pi_{200}$

55

# HOWEVER…



**Document 1**   **Document 2**

$2^{64}$   $2^{64}$

$2^{64}$   $2^{64}$

A   B

$2^{64}$   $2^{64}$

$2^{64}$   $2^{64}$

Why?

Theorem:
   Jaccard (D1, D2) = Prob(A = B)

56

# SET SIMILARITY OF SETS $C_I$ , $C_J$

$$\text{Jaccard}(C_i, C_j) = \frac{|C_i \cap C_j|}{|C_i \cup C_j|}$$

- View sets as columns of a matrix A; one row for each element in the universe. $a_{ij} = 1$ indicates presence of item i in set j
- Example

**C₁  C₂**

| | |
|---|---|
| 0 | 1 |
| 1 | 0 |
| 1 | 1 |
| 0 | 0 |
| 1 | 1 |
| 0 | 1 |

Jaccard**(C₁,C₂) = 2/5 = 0.4**

| C1 | C2 | C3 |
|----|----|----|
| 1  | 0  | 1  |
| 1  | 0  | 0  |
| 0  | 0  | 0  |
| 1  | 1  | 1  |
| 0  | 1  | 1  |
| 0  | 0  | 1  |
| 1  | 1  | 1  |
| 0  | 1  | 0  |
| 1  | 0  | 1  |

- By Jaccard, which one is more similar to C1: is it C2 or C3? Why?

58

# KEY OBSERVATION

- For columns $C_i$, $C_j$, four types of rows

|   | $C_i$ | $C_j$ |
|---|-------|-------|
| **A** | 1 | 1 |
| **B** | 1 | 0 |
| **C** | 0 | 1 |
| **D** | 0 | 0 |

- Overload notation: A = # of rows of type A
- **Claim**

$$\text{Jaccard}(C_i, C_j) = \frac{A}{A + B + C}$$

# "MIN" HASHING

- Randomly permute rows
- Hash $h(C_i)$ = index of first row with 1 in column $C_i$
- Surprising Property

$$P\left(h(C_i) = h(C_j)\right) = \text{Jaccard}\left(C_i, C_j\right)$$

- Why?
  - Both are A/(A+B+C)
  - Look down columns $C_i$, $C_j$ until first non-Type-D row
  - $h(C_i) = h(C_j) \longleftrightarrow$ type A row

# MIN-HASH SKETCHES

- Pick *P* random row permutations
- MinHash sketch

$\text{Sketch}_D$ = list of *P* indexes of first rows with 1 in column C

- Similarity of signatures
  - Let $\text{sim}[\text{sketch}(C_i), \text{sketch}(C_j)]$ = fraction of permutations where MinHash values agree
  - Observe $E[\text{sim}(\text{sketch}(C_i), \text{sketch}(C_j))]$ = $\text{Jaccard}(C_i, C_j)$

61

# EXAMPLE

**Signatures**

|  | $S_1$ | $S_2$ | $S_3$ |
|---|---|---|---|
| Perm 1 = (12345) | 1 | 2 | 1 |
| Perm 2 = (54321) | 4 | 5 | 4 |
| Perm 3 = (34512) | 3 | 5 | 4 |

|  | $C_1$ | $C_2$ | $C_3$ |
|---|---|---|---|
| $R_1$ | 1 | 0 | 1 |
| $R_2$ | 0 | 1 | 1 |
| $R_3$ | 1 | 0 | 0 |
| $R_4$ | 1 | 0 | 1 |
| $R_5$ | 0 | 1 | 0 |

**Similarities**

|  | 1-2 | 1-3 | 2-3 |
|---|---|---|---|
| Col-Col | 0.00 | 0.50 | 0.25 |
| Sig-Sig | 0.00 | 0.67 | 0.00 |

62

# ALL SIGNATURE PAIRS

- Now we have an extremely efficient method for estimating a Jaccard coefficient for a single pair of documents.

- But we still have to estimate $N^2$ Jaccard coefficients where $N$ is the number of web pages.
  - Still slow

- One solution: locality sensitive hashing (LSH)

- Another solution: sorting (Henzinger 2006)

# More resources

- IIR Chapter 19