

15-853: Algorithms in the Real World

Indexing and Searching III

- Link Analysis
- Near duplicate removal

15-853

Page 1

Indexing and Searching Outline

Introduction: model, query types

Inverted Indices: Compression, Lexicon, Merging

Vector Models: Weights, cosine distance

Latent Semantic Indexing:

➔ **Link Analysis:** PageRank (Google), HITS

Near Duplicate Removal:

15-853

Page 2

Link Analysis

The goal is to rank pages.

We want to take advantage of the link structure to do this.

Two main approaches:

- **Static:** we will use the links to calculate a ranking of the pages offline (Google)
- **Dynamic:** we will use the links in the results of a search to dynamically determine a ranking (Clever - hubs and authorities)

15-853

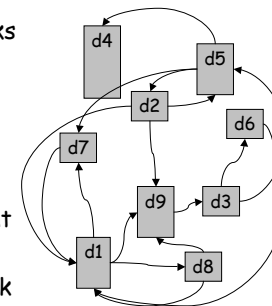
Page 3

The Link Graph

View documents as graph nodes and the hyperlinks between documents as directed edges.

Can give weights on edges (links) based on

1. position in the document
2. weight of anchor term
3. # of occurrences of link
4. ...



15-853

Page 4

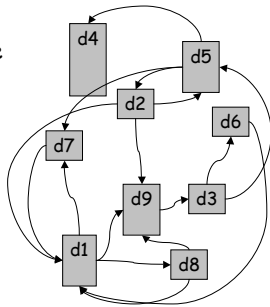
The Link Matrix

An adjacency matrix of the link graph.

Each row corresponds to the Out-Links of a vertex

Each column corresponds to the In-links of a vertex

Often normalize columns to sum to 1 so the dot-product with self is 1.

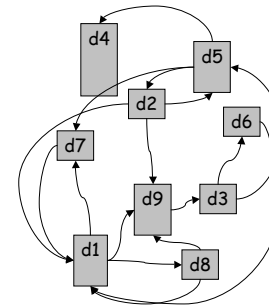


15-853

Page 5

The Link Matrix

0	0	0	0	0	0	1	1	1
1	0	0	0	1	0	0	0	1
0	0	0	0	1	1	0	0	0
0	0	0	0	0	0	0	0	0
0	1	0	1	0	0	1	0	0
1	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	1
0	0	1	0	0	0	0	0	0



15-853

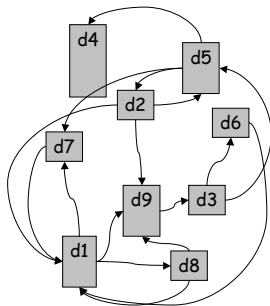
Page 6

Google: Page Rank (1)

Goal: to give a total order (rank) on the "importance" of all pages they index

Possible Method: count the number of in-links.

Problems?



15-853

Page 7

Google: Page Rank (2)

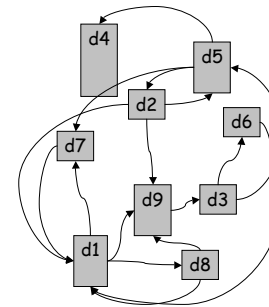
Refined Method: weigh the source by its importance.

Seems like this is a self recursive definition.

Imagine the following:

Jane Surfer, has a lot of time to waste, and randomly selects an outgoing link each time she gets to a page.

She counts how often she visits each page.



Problem?

15-853

Page 8

Side Topic: Markov Chains

A **discrete time stochastic process** is a sequence of random variables $\{X_0, X_1, \dots, X_n, \dots\}$ where the $0, 1, \dots, n, \dots$ are discrete points in time.

A **Markov chain** is a discrete-time stochastic process defined over a finite (or countably infinite) set of states S in terms of a matrix P of **transition probabilities**.

Memorylessness property: for a Markov chain

$$\Pr[X_{t+1} = j \mid X_0 = i_0, X_1 = i_1, \dots, X_t = i] = \Pr[X_{t+1} = j \mid X_t = i]$$

Side Topic: Markov Chains

Let $\pi_i(t)$ be the probability of being in state i at time step t .

Let $\pi(t) = [\pi_0(t), \pi_1(t), \dots]$ be the vector of probabilities at time t .

For an initial probability distribution $\pi(0)$, the probabilities at time t are

$$\pi(t) = \pi(0) P^t$$

A probability distribution π is **stationary** if $\pi = \pi P$ (i.e. is an eigenvector of P with eigenvalue 1)

Side Topic: Markov Chains

A Markov chain is **irreducible** if the underlying graph (of P) consists of a single strong component (i.e. all states are reachable from all other states).

A **period** of a state i is the integer t such that $P^n_{ii} = 0$ for all $n \neq t, 2t, 3t, \dots$ (i.e. only comes back to itself every t steps).

A state is **aperiodic** if it has period 1.

A Markov chain is **aperiodic** if all states are aperiodic.

A state is **recurrent** if upon entering this state the process will definitely return to the state

Side Topic: Markov Chains

Fundamental Theorem of Markov Chains:

Any irreducible, finite, and aperiodic Markov chain has the following properties:

1. All states are **ergodic** (aperiodic and recurrent)
2. There is a **unique stationary distribution** $\pi > 0$ for all states.
3. Let $N(i,t)$ be the number of times the Markov chain visits state i in t steps. Then
$$\lim_{t \rightarrow \infty} \frac{N(i,t)}{t} = \pi_i$$

Google: Page Rank (3)

Stationary probability of a Markov Chain.

Transition probability matrix is simply matrix with equal probability of leaving on each link

Remaining problem: Gives too high of a rank to the Yahoo's of the world, and too low to clusters with low connectivity into the cluster.

Solution: With some probability jump to a random page instead a random outgoing link.

Transition probability matrix is:

$$T = \epsilon U + (1 - \epsilon)A$$

U is the uniform matrix, and A is the normalized link matrix (each column sums to 1).

Google: Page Rank (4)

Want to find π , such that $T\pi = \pi$

This corresponds to finding the principal eigenvector.

Methods:

- **Power method:** iterate $p_{i+1} = Tp_i$ until it settles
pick p_0 randomly,
decomposing p_0 into eigenvectors $a_1 e_1 + a_2 e_2 + \dots$
we have $p_i = \lambda_1^i a_1 e_1 + \lambda_2^i a_2 e_2 + \dots$
- **Multilevel method:** solve on contracted graph, use as initial vector for power method on the expanded graph.
- **Lancoz method:** diagonalize, then iterate

All methods are quite expensive

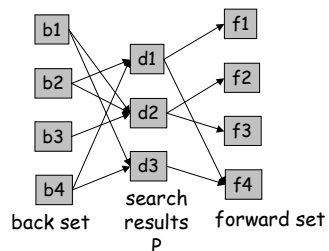
Hubs and Authorities (1)

Goal: To get a ranking for a particular query (instead of the whole web)

Assume: We have a search engine that can give a set of pages P that match a query

Find all pages that point to P (**back set**) and that P point to (**forward set**).

Include b, d and f in a document set D.



Hubs and Authorities (2)

For a particular query we might want to find:

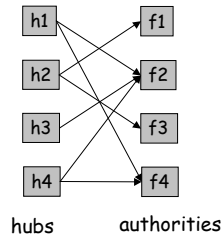
1. The "authorities" on the topic (where the actual information is kept)
2. The "hub" sites that point to the best information on the topic (typically some kind of link site)

Important authorities will have many pages that point to them

Important hubs will point to many pages

But: As with page-rank, just having a lot of pointers does not mean much

Hubs and Authorities (3)



vector h is a weight for each document in D for how good a hub it is

vector a is a weight for each document in D for how good an authority it is

T is the adjacency matrix

Now repeat:

$$h = T^T a$$

$$a = T h$$

until it settles.

15-853

Page 17

Hubs and Authorities (4)

What is this process doing?

$$a_{i+1} = (T T^T) a_i$$

$$h_{i+1} = (T^T T) h_i$$

Power method for eigenvectors of $T T^T$ and $T^T T$

SVD(T) : first column of U and first row of V

15-853

Page 18

Hubs and Authorities (5)

Problem (Topic Drift): Hubs and Authorities will tend to drift to the general, instead of specific. e.g. you search for "Indexing Algorithms" which identifies all pages with those terms in it, and the back and forward pages.

Now you use the SVD to find Hubs + Authorities. These will most likely be pages on Algorithms in general rather than Indexing Algorithms.

Solution: weight the documents with "indexing" prominently in them, more heavily.

$a = T W h$, $h = T^T W a$, where W is a diagonal matrix with weights on the diagonal.

Same as finding SVD of $T' = T W$

15-853

Page 19

Extensions

What about second eigenvectors?

Mix terms and links.

15-853

Page 20

Indexing and Searching Outline

Introduction: model, query types

Inverted Indices: Compression, Lexicon, Merging

Vector Models: Weights, cosine distance

Latent Semantic Indexing:

Link Analysis: PageRank (Google), HITS

➡ **Near Duplicate Removal:**

15-853

Page 21

Syntactic Clustering of the Web

Problem: Many pages on the web are almost the same but not exactly the same.

- Mirror sites with local identifier and/or local links.
- Spam sites that try to improve their "pagerank"
- Plagiarized sites
- Revisions of the same document
- Program documentation converted to HTML
- Legal documents

These near duplicates cause web indexing sites a huge headache.

15-853

Page 22

Syntactic Clustering

Goal: find pages for which either

1. Most of the content is the same
2. One is contained in the other (possibly with changes)

Constraints:

- Need to do it with only a small amount of information per document: a **sketch**
- Comparison of documents should take time linear in the size of the sketch
- Should be reasonably robust against adversary

15-853

Page 23

Syntactic Clustering

Any ideas?

15-853

Page 24

w-Shingling



View each shingle as a value.

$S_w(A)$: the set of w-shingles for document A.

As shorthand, I will often drop the w subscript.

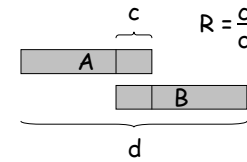
15-853

Page 25

Resemblance and Containment

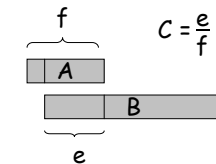
Resemblance:

$$R(A, B) = \frac{|S(A) \cap S(B)|}{|S(A) \cup S(B)|}$$



Containment:

$$C(A, B) = \frac{|S(A) \cap S(B)|}{|S(A)|}$$



15-853

Page 26

Calculating Resemblance/Containment

Method 1:

Keep all shingles and calculate union and intersection of the shingles directly.

Problem: can take $|D|w$ space per document (larger than original document)

Other ideas?

- First 100 shingles?
- Shingles from every k^{th} position?
- Shingles from random positions (preselected)?
e.g. 23, 786, 1121, 1456, ...

15-853

Page 27

Using Min Valued Shingles

How about viewing each shingle as a number and selecting the 100 minimum valued shingles?

- e.g. append the ASCII codes

Possible problem?

15-853

Page 28

Hash

How about a random hash?

i.e. for each shingle x take $h(x)$, then pick the minimum k values.

$$S_k(A) = \min_k\{h(x) : x \in S(A)\}$$

15-853

Page 29

Some Theory: Pairwise independent

Universal Hash functions: (Pairwise independent)

- H : a finite collection (family) of hash functions mapping $U \rightarrow \{0 \dots m-1\}$

- H is **universal** if,

• for $h \in H$ picked uniformly at random,

• and for all $x_1, x_2 \in U, x_1 \neq x_2$

$$\Pr(h(x_1) = h(x_2)) \leq 1/m$$

The class of hash functions

- $h_{ab}(x) = ((a \cdot x + b) \bmod p) \bmod m$

is universal ($p \geq m$ is a prime, $a = \{1 \dots p-1\}$, $b = \{0 \dots p-1\}$)

15-853

Page 30

Some Theory: Minwise independent

Minwise independent permutations:

- S_n : a finite collection (family) of permutations mapping $\{1 \dots n\} \rightarrow \{1 \dots n\}$

- H is **minwise independent** if,

• for $\pi \in S_n$ picked uniformly at random,

• and for $X \subseteq \{1 \dots n\}$, and all $x \in X$

$$\Pr(\min\{\pi(X)\} = \pi(x)) \leq 1/|X|$$

It is actually hard to find a "compact" collection of hash functions that is minwise independent, but we can use an approximation.

15-853

Page 31

Back to similarity

If $\pi \in S_n$ and S_n is minwise independent then:

$$\Pr(\min\{\pi(S(A))\} = \min\{\pi(S(B))\}) = \frac{|S(A) \cap S(B)|}{|S(A) \cup S(B)|} = r(A, B)$$

This suggests we could just keep one minimum value as our "sketch", but our confidence would be low.

What we want for a sketch of size k is either

- use k π 's,

- or keep the k minimum values for one π

15-853

Page 32