







8.1 Ranking Web Pages

Page Rank [Brin & Page 98]

• Idea: the more high ranked pages link to a web page, the higher its rank.

$$PageRank(v) = p + (1-p) \sum_{u \to v} \frac{PageRank(u)}{OutDegree(u)}$$

• Interpretation by random walk:

PageRank is the probability that a "random surfer" visits a page

- » Parameter p is probability that the surfer gets bored and starts on a new random page.
- » (1-p) is the probability that the random surfer follows a link on current page.

• PageRanks correspond to principal eigenvector of the normalized link matrix.

• Can be calculated using an efficient iterative algorithm.

SFU, CMPT 740, 03-3, Martin Ester



8.1 Ranking Web Pages

HITS

Method

- Collect seed set of pages S (e.g., returned by search engine).
- Expand seed set to contain pages that point to or are pointed to by pages in seed set.
- Initialize all hub/authority weights to 1.
- Iteratively update hub weight h(p) and authority weight a(p) for each page:

$$a(p) = \sum_{q \to p} h(q)$$
 $h(p) = \sum_{p \to q} a(q)$

• Stop, when hub/authority weights converge.

SFU, CMPT 740, 03-3, Martin Ester































8.3 Text Clustering	
Suffix Tree Clustering [Zamir & Etzioni 1998]	
Forming Clusters	
Not by similar feature vectors	
But by common terms	
Strengths of Suffix Tree Clustering (STC)	
Efficiency: runtime $O(n)$ for <i>n</i> text documents	
Overlapping clusters	
Method	
1. Identification of "basic clusters"	
2. Combination of basic clusters	
SFU, CMPT 740, 03-3, Martin Ester	398



8.3 Text Clustering

Combination of Basic Clusters

- Basic clusters are highly overlapping
- Merge basic clusters having too much overlap
- *Basic clusters graph*: nodes represent basic clusters Edge between A and B iff $|A \cap B| / |A| > 0.5$ and $|A \cap B| / |B| > 0.5$
- Composite cluster:

a component of the basic clusters graph

• Drawback of this approach:

Distant members of the same component need not be similar No evaluation on standard test data

400

SFU, CMPT 740, 03-3, Martin Ester

8.3 Text Clustering Example from the Grouper System Groups Grouper All results (388) • load balancing (15) polyserve server clustering and load balancing (7 window nt (18) clustering algorithm (32) numerical taxonomy (6) clustering Search data set (12) international conference on (7) Results from each engine: 200 💌 Search for All of these words 💌 Show results initially as: © Mixed+Grouper I C Mixed C Index C Clusters C Combined C Ranked List hierarchical clustering (16) technical report (6) . department of computer science (7) Show debug output? 🗆 Yes cluster analysis (11) unsupervised learning (6) • Use new DF ranking scheme? 🗖 Yes research center (6) linux clustering (11) availability clustering (11) . server clustering (13) Default Search 10 seconds search for good results; gets results from all engines clustering ru (6) cluster (92) clustering technology (9) Quality Search 30 seconds search for **best** results, downloads documents from Web • nt clustering (10) SFU, CMPT 740, 03-3, Martin Ester 401

















8.4 Text Classification

Bag of Words Model

Problem

- Term t_i appears in no training document of class c_i
- *t_i* appears in a document d to be classified
- Document d also contains terms which strongly indicate class *c_i*

$$P(d_j | c) = 0$$
 and $P(d | c) = 0$

Solution

• Smoothing of the relative frequencies in the training documents

SFU, CMPT 740, 03-3, Martin Ester















References

Beil F., Ester M., Xu X.: "Frequent Term-Based Text Clustering", Proc. KDD '02, 2002. Brin S., Page L.: "The anatomy of a large-scale hypertextual Web search engine", Proc. WWW7, 1998.

Chakrabarti S., Dom B., Indyk P.: "Enhanced hypertext categorization using hyperlinks", Proc. ACM SIGMOD 98, 1998.

Craven M., DiPasquo D., Freitag D., McCallum A., Mitchell T., Nigam K., Slattery S.: "Learning to Extract Symbolic Knowledge from the World Wide Web", Proc. AAAI 98, 1998.

Deerwater S., Dumais S. T., Furnas G. W., Landauer T. K., Harshman R.: "Indexing by latent semantic analysis", Journal of the Society for Information Science, 41(6), 1990.

Kleinberg J.: "Authoritative sources in a hyperlinked environment", Proc. SODA, 1998. Nigam K., McCallum A., Thrun S., Mitchell T.: "Text Classification from Labeled and

Unlabeled Documents using EM", Machine Learning, 39(2/3). pp. 103-134. 2000.

Steinbach M., Karypis G., Kumar V.: "A comparison of document clustering techniques", Proc. KDD Workshop on Text Mining, 2000.

Zamir O., Etzioni O.: "Web Document Clustering: A Feasibility Demonstration", Proc. SIGIR 1998.

SFU, CMPT 740, 03-3, Martin Ester