
TEXT RETRIEVAL AND SEARCH ENGINES

The basic concepts, principles, and the major techniques in text retrieval,
which is the underlying science of search engines.

Course author:

ChengXiang Zhai



*University of Illinois at Urbana-Champaign
&
Coursera*

2015

Contents

1 Natural Language Content Analysis	5
1.1 An Example of NLP	5
1.2 The State of the Art	5
1.3 Recommended reading	5
2 Text Access	6
2.1 Two Modes of Text Access: Pull vs. Push	6
2.2 Pull Mode: Querying vs. Browsing	6
2.3 Recommended reading	6
3 Text Retrieval Problem	7
3.1 What Is Text Retrieval?	7
3.2 Formal Formulation of TR	7
3.3 How to Compute $R'(q)$	7
3.4 Theoretical Justification for Ranking	8
3.5 Recommended reading	8
4 Overview of Text Retrieval Methods	9
4.1 How to Design a Ranking Function	9
4.2 Retrieval Models	9
4.3 Common Ideas in State of the Art Retrieval Models	9
4.4 Which Model Works the Best?	10
4.5 Recommended reading	10
5 Vector Space Retrieval Model	11
5.1 Basic Idea	11
5.1.1 Vector Space Model (VSM): Illustration	11
5.1.2 VSM Is a Framework	11
5.1.3 What VSM Doesn't Say	11
5.2 Simplest VSM = Bit-Vector + Dot-Product + BOW	12
5.3 Improved Instantiation	12
5.3.1 Improved VSM with Term Frequency (TF) Weighting	12
5.3.2 IDF Weighting: Penalizing Popular Terms	13
5.3.3 Adding Inverse Document Frequency (IDF)	13
5.4 TF Transformation	13
5.4.1 Ranking Function with TF-IDF Weighting	13
5.4.2 TF Transformation: BM25 Transformation	14
5.4.3 Summary	14
5.5 Doc Length Normalization	14
5.5.1 Pivoted Length Normalization	14
5.5.2 State of the Art VSM Ranking Functions	15
5.6 Further Improvement of VSM?	15
5.7 Further Improvement of BM25	15
5.8 Summary of Vector Space Model	15
5.9 Recommended reading	17

6 System Implementation	18
6.1 Implementation of TR Systems	18
6.1.1 Typical TR System Architecture	18
6.1.2 Tokenization	18
6.1.3 Inverted Index	18
6.1.4 Empirical Distribution of Words	18
6.1.5 Zipf's Law	19
6.1.6 Data Structures for Inverted Index	19
6.2 Inverted Index Construction	19
6.2.1 Inverted Index Compression	20
6.2.2 Integer Compression Methods	20
6.3 Fast Search	20
6.3.1 General Form of Scoring Function	20
6.3.2 A General Algorithm for Ranking Documents	20
6.3.3 Further Improving Efficiency	21
6.4 Some Text Retrieval Toolkits	21
6.5 Summary of System Implementation	21
6.6 Recommended reading	21
7 Evaluation of Text Retrieval Systems	22
7.1 The Cranfield Evaluation Methodology	22
7.2 Basic Measures	22
7.2.1 Evaluating a Set of Retrieved Docs	22
7.2.2 Combine Precision and Recall: F-Measure	22
7.3 Evaluating a Ranked List	23
7.3.1 Evaluating Ranking: Precision-Recall (PR) Curve	23
7.3.2 How to Summarize a Ranking	23
7.3.3 Mean Average Precision (MAP)	23
7.3.4 Summary	24
7.4 Multi-level Relevance Judgments	24
7.5 Statistical Significance Tests	25
7.5.1 Sign test	25
7.5.2 Wilcoxon signed-rank test	26
7.6 Pooling: Avoid Judging all Documents	27
7.7 Summary of TR Evaluation	27
7.8 Recommended reading	27
8 Probabilistic Model	28
8.1 Basic Idea of Probabilistic Model	28
8.2 Language Model	28
8.2.1 The Simplest Language Model: Unigram LM	28
8.2.2 Recommended reading	29
8.3 Ranking based on Query Likelihood	29
8.3.1 How to Estimate $p(w d)$	29
8.3.2 Rewriting the Ranking Function with Smoothing	30
8.3.3 Smoothing Methods	30
8.3.4 Summary of Query Likelihood Probabilistic Model	31
8.3.5 Recommended reading	31

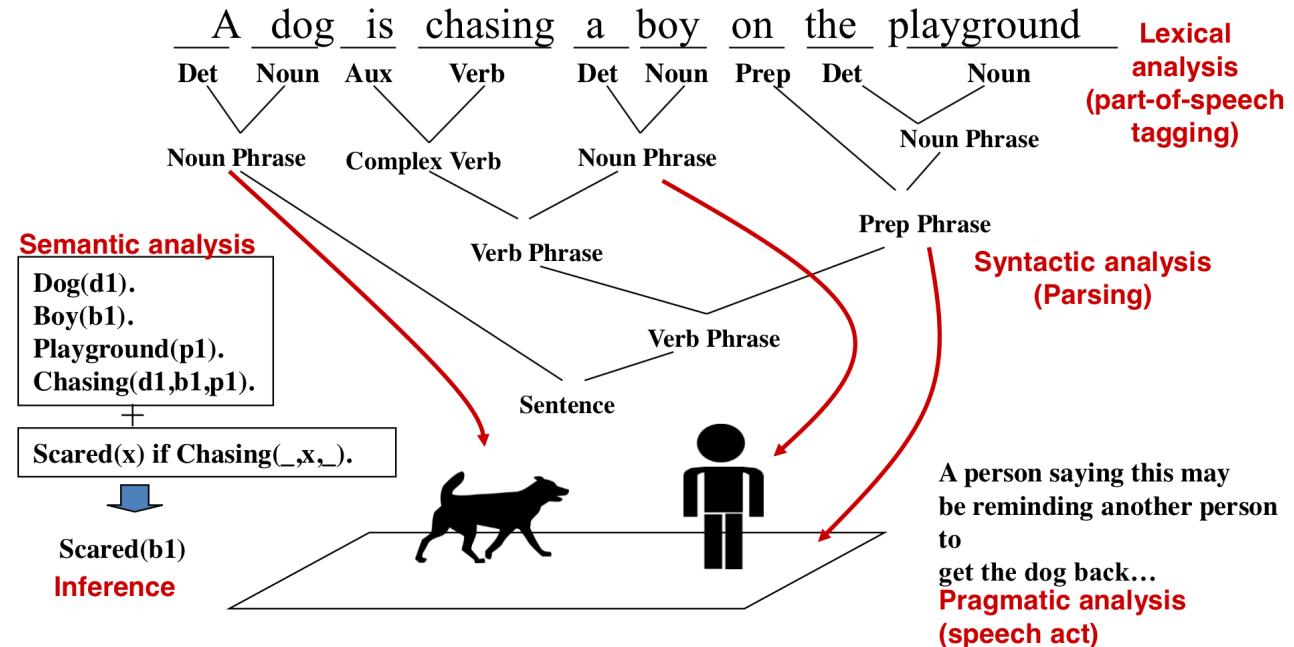
9 Feedback in Text Retrieval	32
9.1 Types of Feedback	32
9.2 Feedback in Vector Space Model	32
9.2.1 Rocchio Feedback	32
9.2.2 Rocchio in Practice	33
9.3 Feedback in Language Model	33
9.3.1 The KL-divergence measure	33
9.3.2 Kullback-Leibler (KL) Divergence Retrieval Model	33
9.3.3 Feedback as Model Interpolation	34
9.3.4 Generative Mixture Model	34
9.4 Recommended reading	35
10 Web Search	36
10.1 Web Search: Challenges & Opportunities	36
10.2 Basic Search Engine Technologies	36
10.3 Crawler/Spider/Robot	36
10.3.1 Major Crawling Strategies	37
10.4 Web Index	37
10.4.1 GFS Architecture	38
10.4.2 MapReduce: A Framework for Parallel Programming	38
10.4.3 MapReduce: Computation Pipeline	38
10.4.4 Inverted Indexing with MapReduce	39
10.5 Link Analysis	40
10.5.1 Ranking Algorithms for Web Search	40
10.5.2 Exploiting Inter-Document Links	40
10.5.3 PageRank: Capturing Page «Popularity»	40
10.5.4 The PageRank Algorithm	41
10.5.5 PageRank in Practice	42
10.5.6 HITS: Capturing Authorities & Hubs	42
10.5.7 The HITS Algorithm	42
10.6 Learning to Rank	43
10.6.1 How Can We Combine Many Features?	43
10.6.2 Regression-Based Approaches	43
10.6.3 More Advanced Learning Algorithms	43
10.6.4 Recommended reading	44
10.7 Future of Web Search	44
10.7.1 Next Generation Search Engines	44
10.7.2 Future Intelligent Information Systems	44
11 Recommender Systems	45
11.1 Recommender ≈ Filtering System	45
11.1.1 Basic Filtering Question	45
11.2 Content-Based Filtering	45
11.2.1 A Typical Content-Based Filtering System	45
11.2.2 Three Basic Problems in Content-Based Filtering	45
11.2.3 Extend a Retrieval System for Information Filtering	46
11.2.4 A General Vector-Space Approach	46
11.2.5 Difficulties in Threshold Learning	46
11.2.6 Empirical Utility Optimization	46
11.2.7 Beta-Gamma Threshold Learning	47

11.3	Collaborative Filtering	47
11.3.1	What is Collaborative Filtering (CF)?	47
11.3.2	Memory-based Approaches to Collaboration Filtering Problem	48
11.3.3	User Similarity Measures	48
11.3.4	Improving User Similarity Measures	48
11.3.5	Summary	48
11.3.6	Recommended reading	49
12	Course Summary	50
12.1	Recommended reading	50
12.2	Main Techniques for Harnessing Big Text Data: Text Retrieval + Text Mining	50

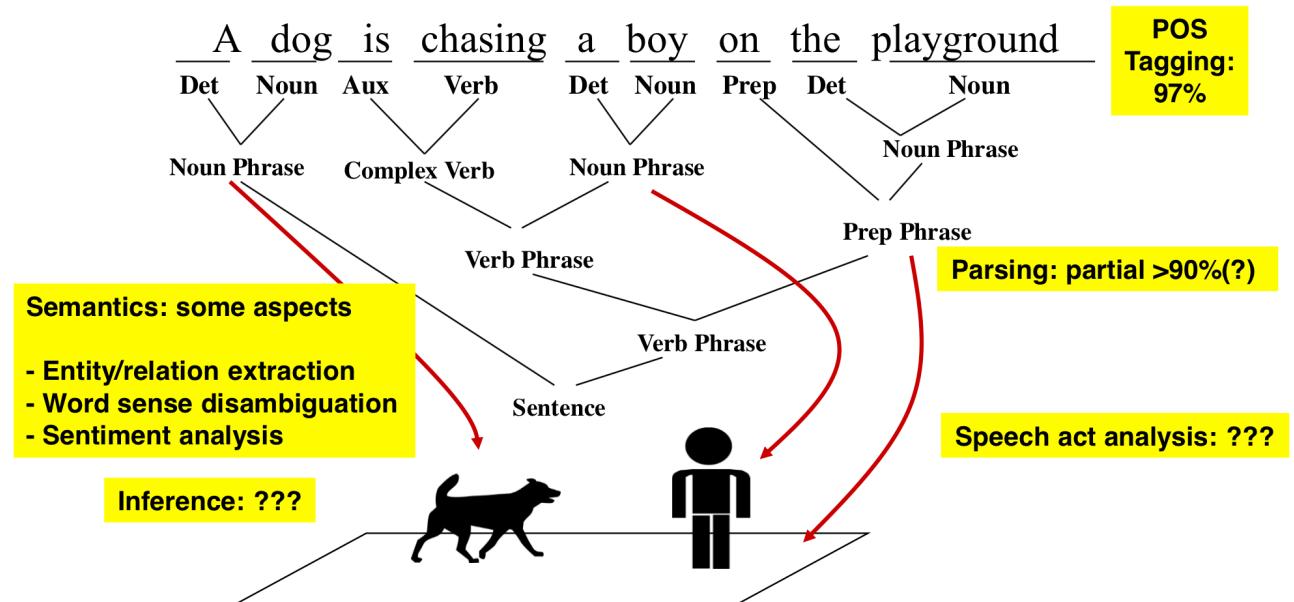
1 Natural Language Content Analysis

NLP = Natural Language Processing

1.1 An Example of NLP



1.2 The State of the Art



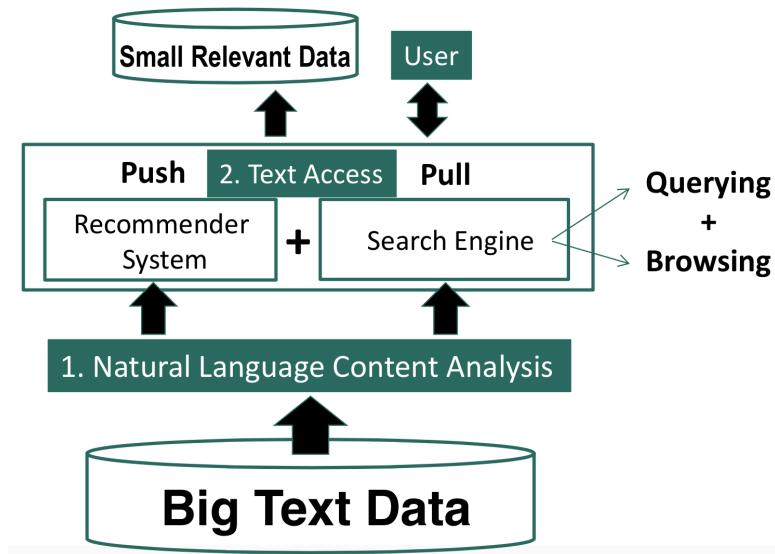
1.3 Recommended reading

- Chris Manning and Hinrich Schütze, «Foundations of Statistical Natural Language Processing», MIT Press. Cambridge, MA: May 1999.

2 Text Access

2.1 Two Modes of Text Access: Pull vs. Push

- Pull Mode (search engines) – Users take initiative
 - Ad hoc information need
- Push Mode (recommender systems)
 - Systems take initiative
 - Stable information need or system has good knowledge about a user's need



2.2 Pull Mode: Querying vs. Browsing

- Querying
 - User enters a (keyword) query
 - System returns relevant documents
 - Works well when the user knows what keywords to use
- Browsing
 - User navigates into relevant information by following a path enabled by the structures on the documents
 - Works well when the user wants to explore information, doesn't know what keywords to use, or can't conveniently enter a query

2.3 Recommended reading

- N. J. Belkin and W. B. Croft. 1992. «Information filtering and information retrieval: two sides of the same coin?» Commun. ACM 35, 12 (Dec. 1992), 29-38.

3 Text Retrieval Problem

3.1 What Is Text Retrieval?

TR = Text Retrieval¹

- Collection of text documents exists
- User gives a query to express the information need
- Search engine system returns relevant documents to users
- Often called “information retrieval” (IR), but IR is actually much broader
- Known as «search technology» in industry

TR is an empirically defined problem:

- Can't mathematically prove one method is better than another
- Must rely on empirical evaluation involving users!

3.2 Formal Formulation of TR

- **Vocabulary:** $V = \{w_1, w_2, \dots, w_N\}$ of language
- **Query:** $q = q_1, \dots, q_m$, where $q_i \in V$
- **Document:** $d_i = d_{i1}, \dots, d_{im_i}$, where $d_{ij} \in V$
- **Collection:** $C = \{d_1, \dots, d_M\}$
- **Set of relevant documents:** $R(q) \subseteq C$
 - Generally unknown and user-dependent
 - Query is a «hint» on which doc is in $R(q)$
- **Task:** compute $R'(q)$, an approximation of $R(q)$

3.3 How to Compute $R'(q)$

- Strategy 1: Document selection
 - $R'(q) = \{d \in C \mid f(d, q) = 1\}$, where $f(d, q) \in \{0, 1\}$ is an indicator function or binary classifier
 - System must decide if a doc is relevant or not (absolute relevance)
- Strategy 2 (generally preferred): Document ranking
 - $R'(q) = \{d \in C \mid f(d, q) > \theta\}$, where $f(d, q) \in \mathfrak{R}$ is a relevance measure function; θ is a cutoff determined by the user
 - System only needs to decide if one doc is more likely relevant than another (relative relevance)

¹Retrieval - поиск

3.4 Theoretical Justification for Ranking

Probability Ranking Principle [Robertson 77]: Returning a ranked list of documents in descending order of probability that a document is relevant to the query is the optimal strategy under the following two assumptions:

- The utility of a document (to a user) is independent of the utility of any other document
- A user would browse the results sequentially

3.5 Recommended reading

- S.E. Robertson, «The probability ranking principle in IR». Journal of Documentation 33, 294-304, 1977
- **C. J. van Rijsbergen, «Information Retrieval», 2nd Edition**, Butterworth-Heinemann, Newton, MA, USA, 1979

4 Overview of Text Retrieval Methods

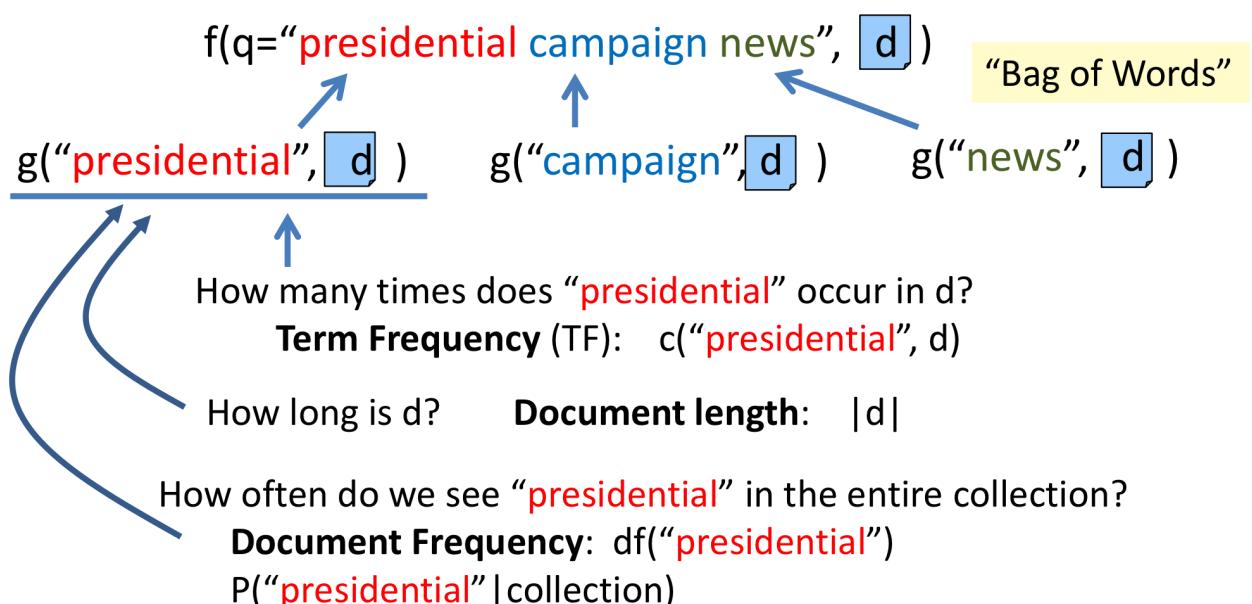
4.1 How to Design a Ranking Function

- **Query:** $q = q_1, \dots, q_m$, where $q_i \in V$
- **Document:** $d = d_1, \dots, d_n$, where $d_i \in V$
- **Ranking function:** $f(q, d) \in \mathfrak{R}$
- **Key challenge:** how to measure the likelihood that document d is relevant to query q
- **Retrieval model:** formalization of relevance (give a computational definition of relevance)

4.2 Retrieval Models

- **Similarity-based models:** $f(q, d) = \text{similarity}(q, d)$
 - Vector space model
- **Probabilistic models:** $f(d, q) = p(R = 1 \mid d, q)$, where $R \in \{0, 1\}$
 - Classic probabilistic model
 - Language model
 - Divergence-from-randomness model
- **Probabilistic inference model:** $f(q, d) = p(d \rightarrow q)$
- **Axiomatic model:** $f(q, d)$ must satisfy a set of constraints

4.3 Common Ideas in State of the Art Retrieval Models



State of the art ranking functions tend to rely on:

- Bag of words representation
- Term Frequency (TF) and Document Frequency (DF) of words
- Document length

4.4 Which Model Works the Best?

When optimized, the following models tend to perform equally well [Fang et al. 11]:

- **Pivoted length normalization – BM25**
- Query likelihood
- PL2

4.5 Recommended reading

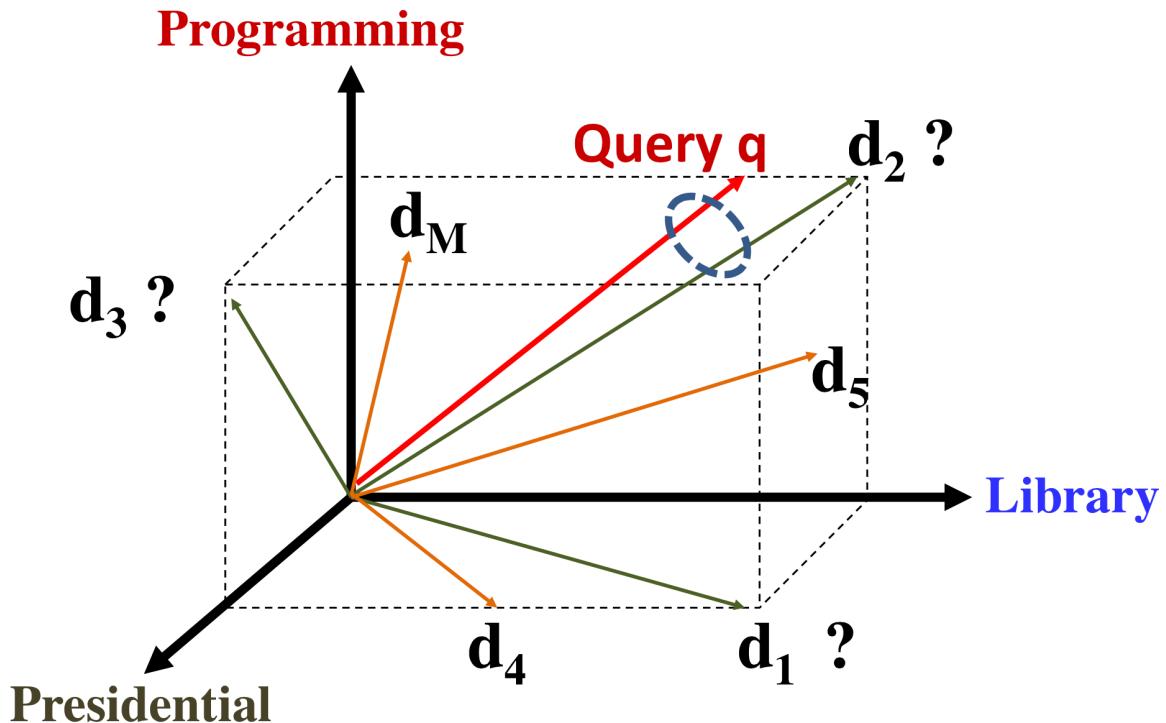
- Hui Fang, Tao Tao, and Chengxiang Zhai. 2011. «Diagnostic Evaluation of Information Retrieval Models». ACM Trans. Inf. Syst. 29, 2, Article 7 (April 2011)
- ChengXiang Zhai, «Statistical Language Models for Information Retrieval», Morgan & Claypool Publishers, 2008. (Chapter 2)

5 Vector Space Retrieval Model

VSM - Vector Space Model

5.1 Basic Idea

5.1.1 Vector Space Model (VSM): Illustration

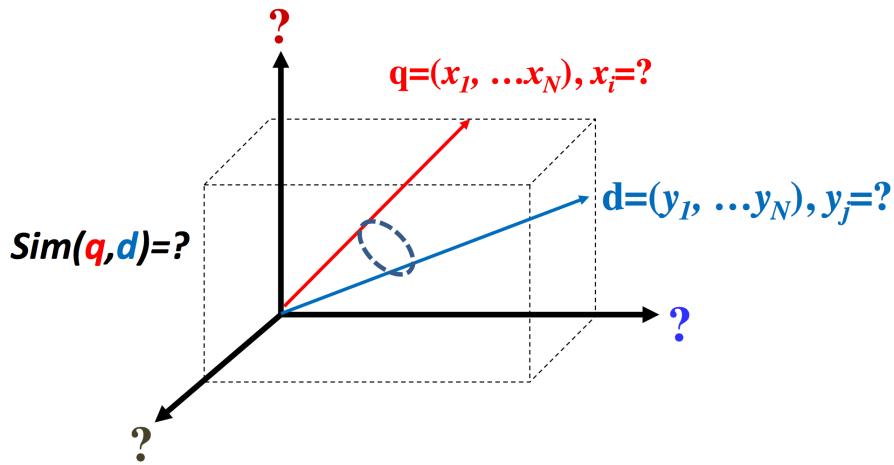


5.1.2 VSM Is a Framework

- Represent a doc/query by a term vector
 - **Term:** basic concept, e.g., word or phrase
 - Each term defines one dimension
 - N terms define an **N-dimensional space**
 - **Query vector:** $q = (x_1, \dots, x_N)$, $x_i \in \mathfrak{R}$ is query term weight
 - **Doc vector:** $d = (y_1, \dots, y_N)$, $y_j \in \mathfrak{R}$ is doc term weight
- $relevance(q, d) \propto similarity(q, d) = f(q, d)$

5.1.3 What VSM Doesn't Say

- How to define/select the “basic concept” – Concepts are assumed to be orthogonal
- How to place docs and query in the space (= how to assign term weights)
 - Term weight in query indicates importance of term
 - Term weight in doc indicates how well the term characterizes the doc
- How to define the similarity measure



5.2 Simplest VSM = Bit-Vector + Dot-Product + BOW

$q = (x_1, \dots, x_N)$ $x_i, y_i \in \{0, 1\}$
 $d = (y_1, \dots, y_N)$ 1: word W_i is present
 0: word W_i is absent

$$Sim(q, d) = q \cdot d = x_1 y_1 + \dots + x_N y_N = \sum_{i=1}^N x_i y_i$$

Simplest VSM:

- Dimension = word
- Vector = 0-1 bit vector (word presence/absence)
- Similarity = dot product
- $f(q, d)$ = number of distinct query words matched in d

5.3 Improved Instantiation

Improved VSM:

- Dimension = word
- Vector = TF-IDF weight vector
- Similarity = dot product

5.3.1 Improved VSM with Term Frequency (TF) Weighting

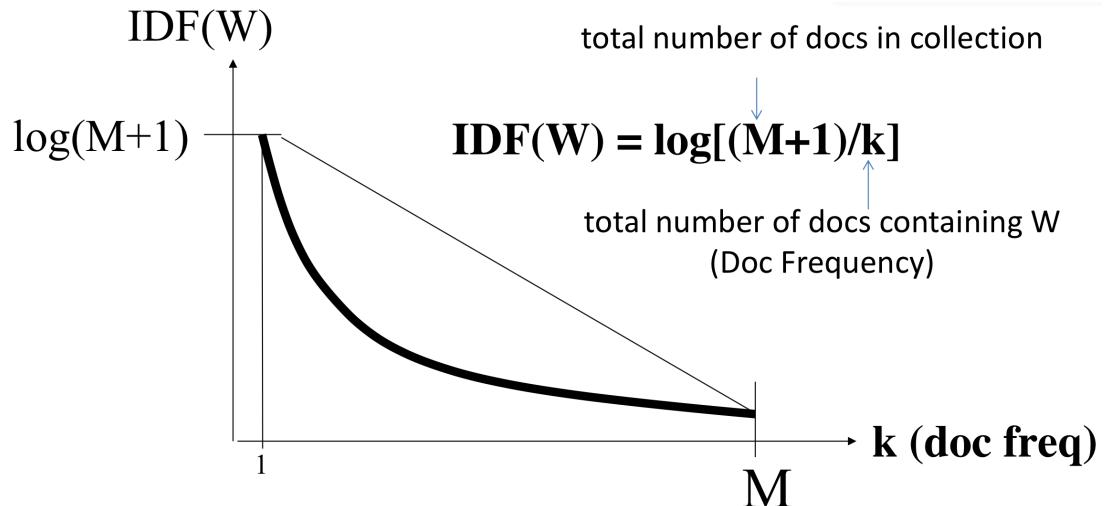
$q = (x_1, \dots, x_N)$ $x_i = \text{count of word } W_i \text{ in query}$

$d = (y_1, \dots, y_N)$ $y_i = \text{count of word } W_i \text{ in doc}$

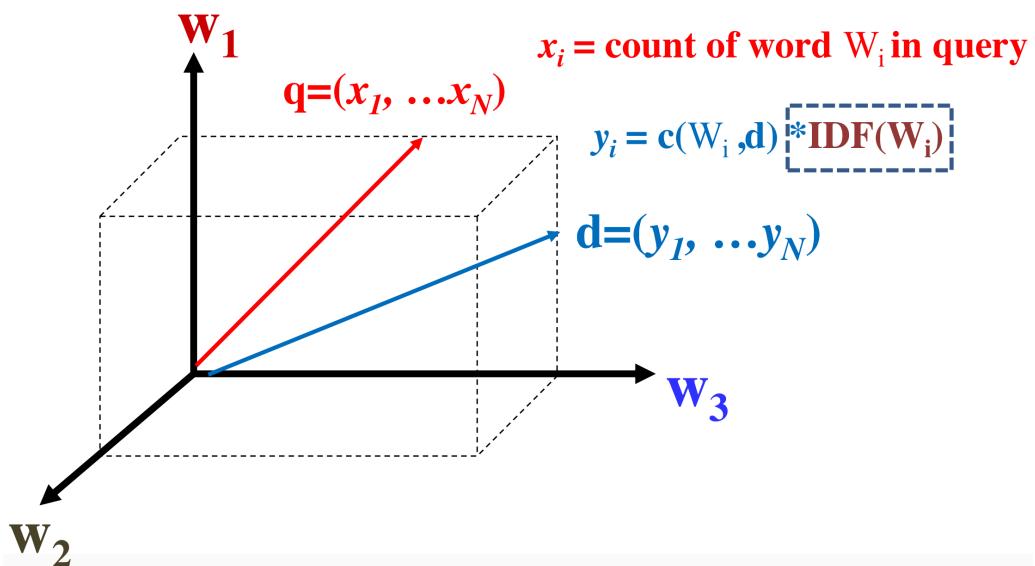
$$Sim(q, d) = q \cdot d = x_1 y_1 + \dots + x_N y_N = \sum_{i=1}^N x_i y_i$$

5.3.2 IDF Weighting: Penalizing Popular Terms

IDF — inverse document frequency



5.3.3 Adding Inverse Document Frequency (IDF)



5.4 TF Transformation

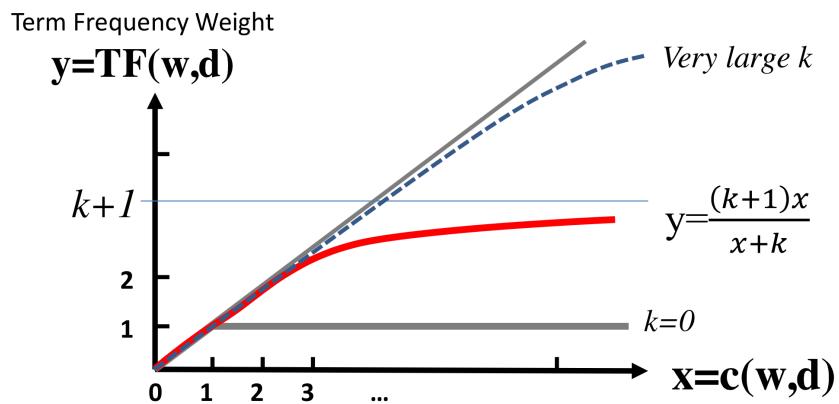
5.4.1 Ranking Function with TF-IDF Weighting

$$f(q, d) = \sum_{i=1}^N x_i y_i = \sum_{w \in q \cap d} c(w, q) c(w, d) \log \frac{M+1}{df(w)}$$

- $w \in q \cap d$ - all matched query (q) words in document (d)
- $c(w, q)$ - count of word w in document d
- M - total number of documents in collection
- $df(w)$ - Doc Frequency (total number of documents containing word w)

5.4.2 TF Transformation: BM25 Transformation

BM = Best Matching



5.4.3 Summary

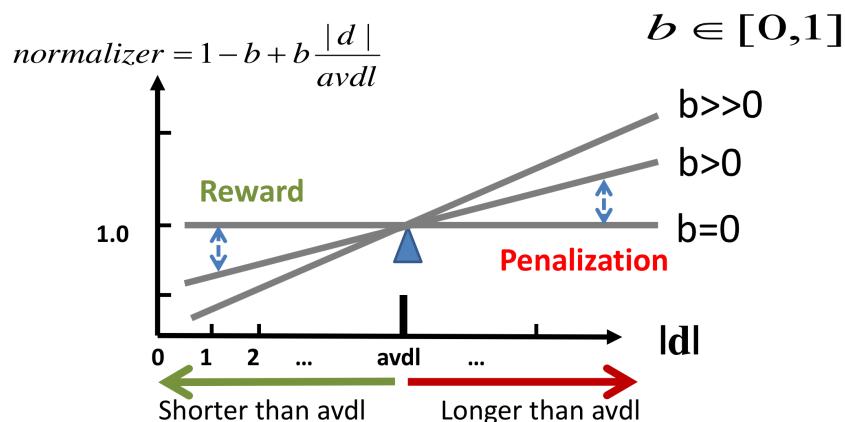
- Sublinear TF Transformation is needed to
 - capture the intuition of «diminishing return» from higher TF
 - avoid dominance by one single term over all others
- BM25 Transformation
 - has an upper bound
 - is robust and effective
- Ranking function with BM25 TF ($k \geq 0$):

$$f(q, d) = \sum_{i=1}^N x_i y_i = \sum_{w \in q \cap d} c(w, q) \cdot \frac{(k+1)c(w, d)}{c(w, d) + k} \log \frac{M+1}{df(w)}$$

5.5 Doc Length Normalization

5.5.1 Pivoted Length Normalization

Pivoted length normalizer: use average doc length as «pivot»². Normalizer = 1 if $|d|$ = average doc length (avdl).



²Pivot - стержень; точка опоры, вращения

5.5.2 State of the Art VSM Ranking Functions

Pivoted Length Normalization VSM [Singhal et al 96]:

$$f(q, d) = \sum_{w \in q \cap d} c(w, q) \frac{\ln[1 + \ln(1 + c(w, d))]}{1 - b + b \frac{|d|}{avdl}} \log \frac{M + 1}{df(w)}$$

BM25/Okapi [Robertson & Walker 94]:

$$f(q, d) = \sum_{w \in q \cap d} c(w, q) \frac{(k + 1) c(w, d)}{c(w, d) + k \left(1 - b + b \frac{|d|}{avdl}\right)} \log \frac{M + 1}{df(w)}$$

5.6 Further Improvement of VSM?

- Improved instantiation of dimension?
 - stemmed words, stop word removal, phrases, latent semantic indexing (word clusters), character n-grams, ...
 - bag-of-words with phrases is often sufficient in practice
 - Language-specific and domain-specific tokenization is important to ensure “normalization of terms”
- Improved instantiation of similarity function?
 - cosine of angle between two vectors?
 - Euclidean?
 - dot product seems still the best (sufficiently general especially with appropriate term weighting)

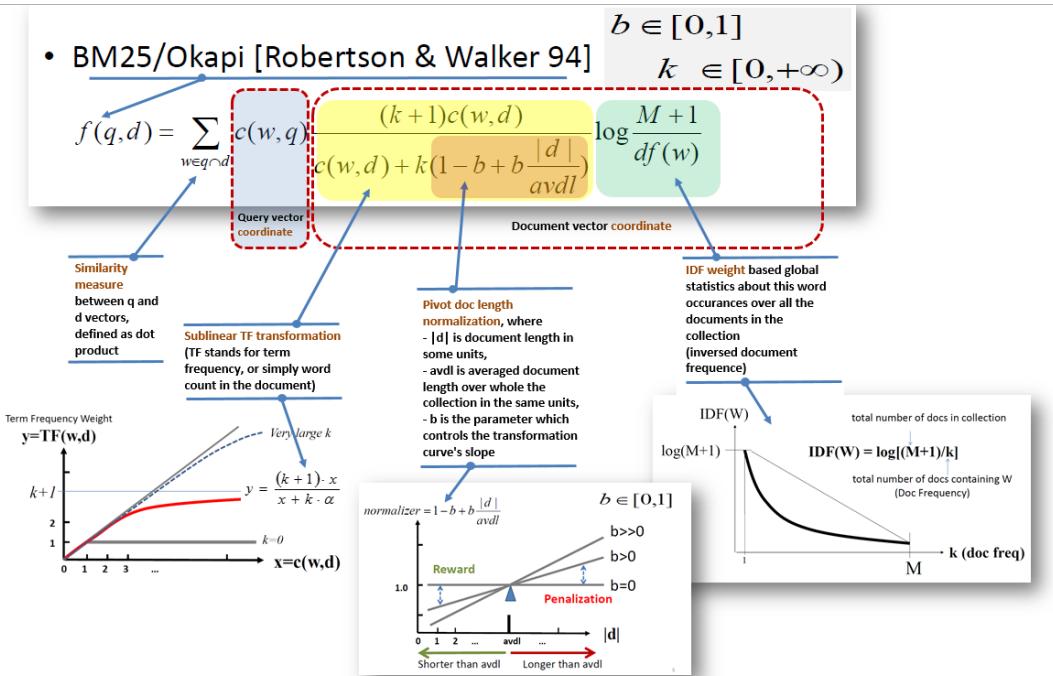
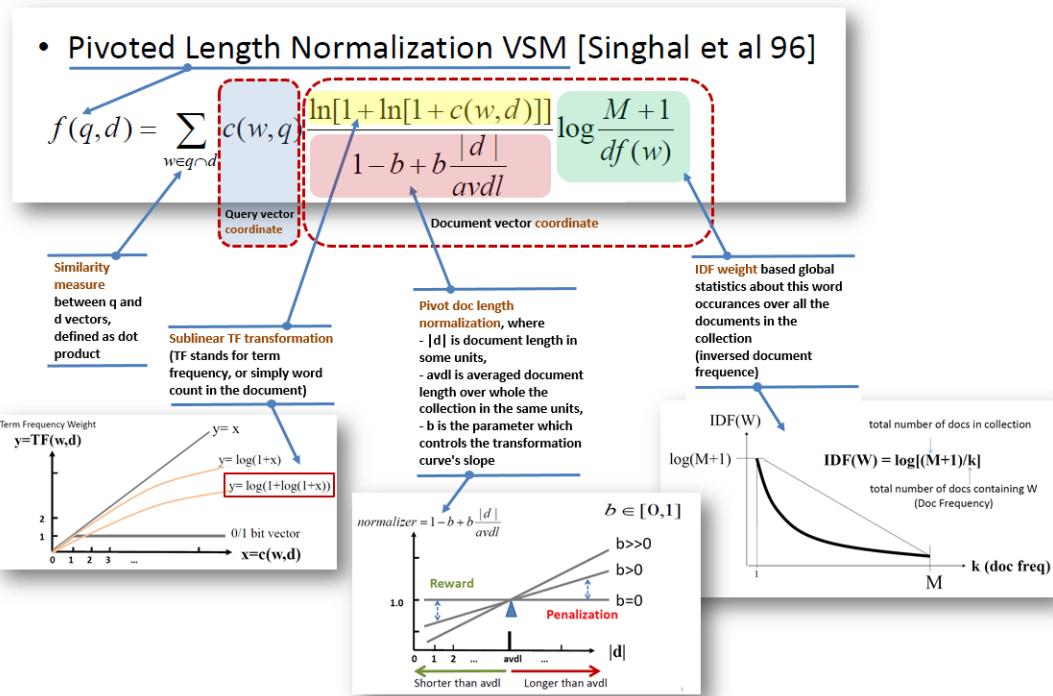
5.7 Further Improvement of BM25

- BM25F [Robertson & Zaragoza 09]
 - Use BM25 for documents with structures («F»=fields)
 - Key idea: combine the frequency counts of terms in all fields and then apply BM25 (instead of the other way)
- BM25+ [Lv & Zhai 11]
 - Address the problem of over penalization of long documents by BM25 by adding a small constant to TF
 - Empirically and analytically shown to be better than BM25

5.8 Summary of Vector Space Model

- Relevance(q,d) = similarity(q,d)
- Query and documents are represented as vectors

- Heuristic³ design of ranking function
- Major term weighting heuristics
 - TF weighting and transformation
 - IDF weighting
 - Document length normalization
- BM25 and Pivoted normalization seem to be most effective



³Heuristic - эвристический

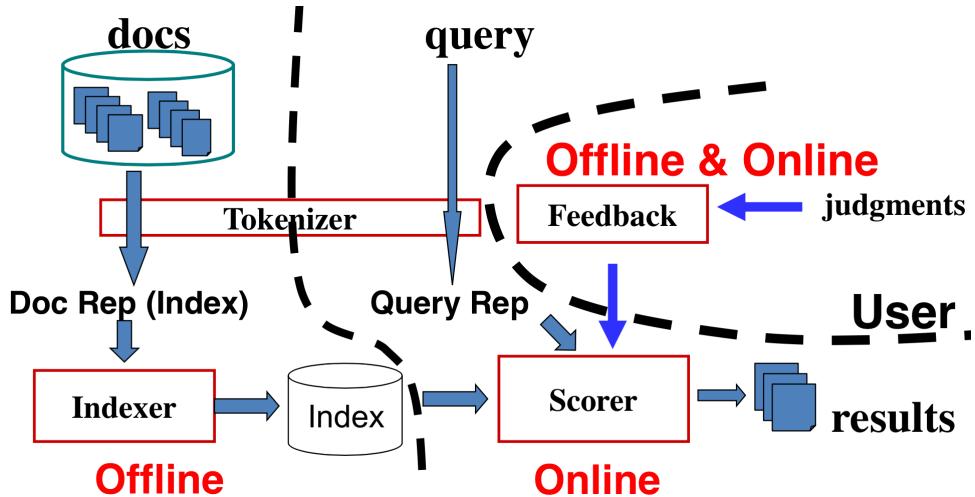
5.9 Recommended reading

- A.Singhal, C.Buckley, and M.Mitra. «Pivoted document length normalization». In Proceedings of ACM SIGIR 1996.
- S. E. Robertson and S. Walker. «Some simple effective approximations to the 2-Poisson model for probabilistic weighted retrieval», Proceedings of ACM SIGIR 1994.
- S. Robertson and H. Zaragoza. «The Probabilistic Relevance Framework: BM25 and Beyond», Found. Trends Inf. Retr. 3, 4 (April 2009).
- Y. Lv, C. Zhai, «Lower-bounding term frequency normalization». In Proceedings of ACM CIKM 2011.

6 System Implementation

6.1 Implementation of TR Systems

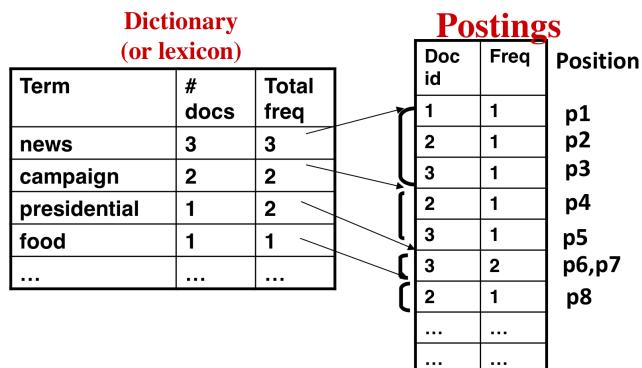
6.1.1 Typical TR System Architecture



6.1.2 Tokenization

- Normalize lexical units: words with similar meanings should be mapped to the same indexing term
- Stemming: mapping all inflectional forms of words to the same root form
- Some languages (e.g., Chinese) pose challenges in word segmentation

6.1.3 Inverted Index

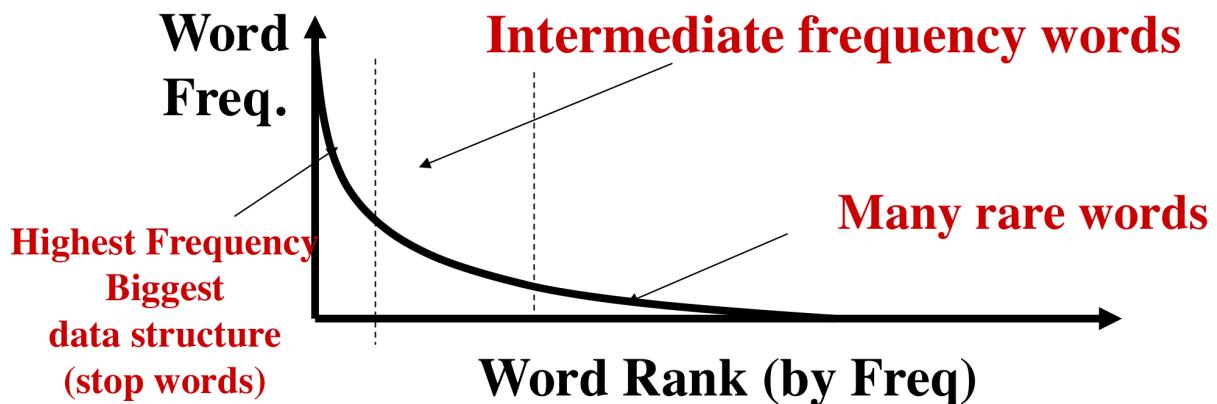


6.1.4 Empirical Distribution of Words

There are stable language-independent patterns in how people use natural languages:

- A few words occur very frequently; most occur rarely. E.g., in news articles:
 - Top 4 words: 10–15% word occurrences
 - Top 50 words: 35–40% word occurrences
- The most frequent word in one corpus may be rare in another

6.1.5 Zipf's Law



$$F(w) = \frac{C}{r(w)^\alpha}, \alpha \approx 1, C \approx 0.1$$

rank \times frequency \approx constant:

- $F(w)$ - word frequency
- $r(w)$ - word rank

6.1.6 Data Structures for Inverted Index

- Dictionary: modest size
 - Needs fast random access
 - Preferred to be in memory
 - Hash table, B-tree, trie, ...
- Postings: huge
 - Sequential access is expected
 - Can stay on disk
 - May contain docID, term freq., term pos, etc
 - Compression is desirable

6.2 Inverted Index Construction

Sort-based method:

- Step 1: Collect local (termID, docID, freq) tuples from documents
- Step 2: Sort local tuples by termID (to make «runs») and save to files
- Step 3: Pair-wise merge runs
- Step 4: Output inverted file

6.2.1 Inverted Index Compression

In general, leverage skewed distribution of values and use variable-length encoding:

- TF compression:
 - Small numbers tend to occur far more frequently than large numbers (Zipf's law)
 - Fewer bits for small (high frequency) integers at the cost of more bits for large integers
- Doc ID compression:
 - «d-gap» (store difference): $d_1, d_2 - d_1, d_3 - d_2, \dots$
 - Feasible due to sequential access

6.2.2 Integer Compression Methods

- **Binary**: equal-length coding
- **Unary**: $x \geq 1$ is coded as $x - 1$ one bits followed by 0, e.g., 3=>110; 5=>11110
- **γ -code**: $x \Rightarrow$ unary code for $1 + \lfloor \log x \rfloor$ followed by uniform code for $x - 2^{\lfloor \log x \rfloor}$ in $\lfloor \log x \rfloor$ bits, e.g., 3=>101, 5=>11001
- **δ -code**: same as γ -code, but replace the unary prefix with γ -code. E.g., 3=>1001, 5=>10101

6.3 Fast Search

6.3.1 General Form of Scoring Function

$$f(q, d) = f_a \left(h \left(g(t_1, d, q), \dots, g(t_k, d, q) \right), f_d(d), f_q(q) \right)$$

- $f_d(d), f_q(q)$ - adjustment factors of document and query
- $g(t_i, d, q)$ - weight of a **matched** query term t_i in d
- $h()$ - weights aggregation function
- $f_a()$ - final score adjustment function

6.3.2 A General Algorithm for Ranking Documents

- $f_d(d)$ - can be precomputed at index time, $f_q(q)$ - at query time
- Maintain a score accumulator for each d to compute h
- For each query term t_i
 - Fetch the inverted list $\{(d_1, f_1), \dots, (d_n, f_n)\}$
 - For each entry (d_j, f_j) , compute $g(t_i, d_j, q)$, and update score accumulator for doc d_i to incrementally compute h
- Adjust the score to compute f_a , and sort

6.3.3 Further Improving Efficiency

- Caching (e.g., query results, list of inverted index)
- Keep only the most promising accumulators
- Scaling up to the Web-scale? (need parallel processing)

6.4 Some Text Retrieval Toolkits

- [Lucene](#)
- [Lemur/Indri](#)
- [Terrier](#)
- [MeTA](#)
- More can be found [here](#)

6.5 Summary of System Implementation

- Inverted index and its construction
 - Preprocess data as much as we can
 - Compression when appropriate
- Fast search using inverted index
 - Exploit inverted index to accumulate scores for documents matching a query term
 - Exploit Zipf's law to avoid touching many documents not matching any query term
 - Can support a wide range of ranking algorithms
- Further scaling up using distributed file system, parallel processing, and caching

6.6 Recommended reading

- Ian H. Witten, Alistair Moffat, Timothy C. Bell: «Managing Gigabytes: Compressing and Indexing Documents and Images», Second Edition. Morgan Kaufmann, 1999.
- Stefan Büttcher, Charles L. A. Clarke, Gordon V. Cormack: «Information Retrieval - Implementing and Evaluating Search Engines». MIT Press, 2010.

7 Evaluation of Text Retrieval Systems

7.1 The Cranfield Evaluation Methodology

A methodology for laboratory testing of system components developed in 1960s. General idea is to build reusable test collections and define measures. A test collection can then be reused many times to compare different systems.

- A sample collection of documents (simulate real document collection)
- A sample set of queries/topics (simulate user queries)
- Relevance judgments (ideally made by users who formulated the queries) => ideal ranked list
- Measures to quantify how well a system's result matches the ideal ranked list

7.2 Basic Measures

7.2.1 Evaluating a Set of Retrieved Docs

	Retrieved	Not Retrieved
Relevant	a	b
Not Relevant	c	d

- Precision: are the retrieved results all relevant?

$$Precision = \frac{a}{a + c}$$

- Recall: have all the relevant documents been retrieved?

$$Recall = \frac{a}{a + b}$$

- In reality, high recall tends to be associated with low precision

7.2.2 Combine Precision and Recall: F-Measure

$$F_{\beta} = \frac{1}{\frac{\beta^2}{\beta^2 + 1} \frac{1}{R} + \frac{1}{\beta^2 + 1} \frac{1}{P}} = \frac{(\beta^2 + 1) \cdot P \cdot R}{\beta^2 \cdot P + R}$$

- P - precision
- R - recall
- β - parameter, often set to 1:

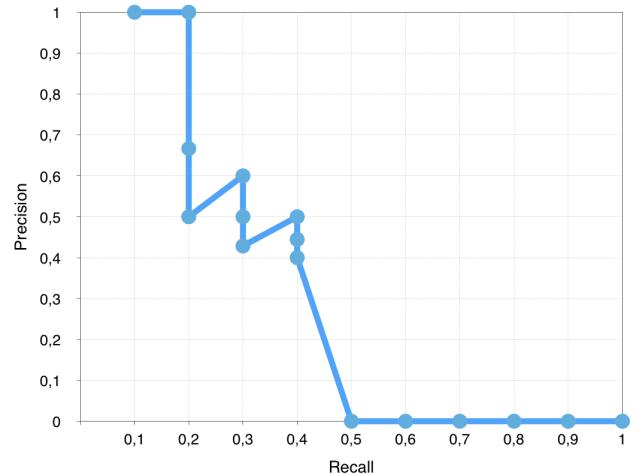
$$F_1 = \frac{2 \cdot P \cdot R}{P + R}$$

7.3 Evaluating a Ranked List

7.3.1 Evaluating Ranking: Precision-Recall (PR) Curve

- Total number of relevant documents in collection: $a + b = 10$
- Number of retrieved documents: $a + c = 10$

Relevance	Precision	Recall
D1 +	1/1	1/10
D2 +	2/2	2/10
D3 -	2/3	2/10
D4 -	2/4	2/10
D5 +	3/5	3/10
D6 -	3/6	3/10
D7 -	3/7	3/10
D8 +	4/8	4/10
D9 -	4/9	4/10
D10 -	4/10	4/10



7.3.2 How to Summarize a Ranking

Average Precision is sensitive to the rank of each relevant document:

$$AveP = \frac{\sum_{k=1}^{a+c} P(k) \cdot rel(k)}{a+b} = \sum_{k=1}^{a+c} P(k) \cdot \Delta r(k)$$

- $a + c$ - number of retrieved documents
- $a + b$ - total number of relevant documents in collection
- $P(k)$ - the precision at cut-off k in the list
- $rel(k)$ - indicator function equaling 1 if the item at rank k is a relevant document, zero otherwise
- $\Delta r(k)$ - the change in recall that happened between cut-off $k - 1$ and cut-off k

In special case, when there's only one relevant document in the collection (e.g., known item search):

- Average Precision = **Reciprocal Rank** = $1/r$, where r is the rank position of the single relevant doc

7.3.3 Mean Average Precision (MAP)

In case of multiple queries:

- MAP = arithmetic mean of average precision over a set of queries

$$MAP = \frac{\sum_{q=1}^N AveP(q)}{N}, \text{ where } N \text{ is the number of queries}$$

- **gMAP** = geometric mean of average precision over a set of queries

$$gMAP = \sqrt[N]{\prod_{q=1}^N AveP(q)} = \exp \frac{\sum_{q=1}^N \log(AveP(q))}{N}$$

7.3.4 Summary

- Precision-Recall curve characterizes the overall accuracy of a ranked list
- The **actual** utility of a ranked list depends on how many top-ranked results a user would examine
- Average Precision is the standard measure for comparing two ranking methods
 - Combines precision and recall
 - Sensitive to the rank of **every** relevant document

7.4 Multi-level Relevance Judgments

Discounted cumulative gain (DCG) is a measure of ranking quality. Two assumptions are made in using DCG and its related measures:

- Highly relevant documents are more useful when appearing earlier in a search engine result list (have higher ranks)
- Highly relevant documents are more useful than marginally relevant documents, which are in turn more useful than irrelevant documents.

For a rank position p:

- **Cumulative Gain:** $CG_p = \sum_{i=1}^p rel_i$, where rel_i is the graded relevance of the result at position i
- **Discounted Cumulative Gain:** $DCG_p = rel_1 + \sum_{i=2}^p \frac{rel_i}{\log_2 i}$
- Alternative version of **Discounted Cumulative Gain:** $DCG_p = \sum_{i=1}^p \frac{2^{rel_i} - 1}{\log_2(i + 1)}$
- **Normalized DCG:** $nDCG_p = \frac{DCG_p}{IDCG_p}$, where $IDCG_p$ is an Ideal DCG (the maximum possible DCG till position p)

For example, each document is to be judged on a scale of 0-3 with 0 meaning irrelevant, 3 meaning completely relevant, and 1 and 2 meaning «somewhere in between»

Document	rel_i	$\frac{rel_i}{\log_2 i}$
D1	3	-
D2	2	2
D3	3	1.892
D4	0	0
D5	1	0.431
D6	2	0.774

- $DCG_6 = rel_1 + \sum_{i=2}^6 \frac{rel_i}{\log_2 i} = 3 + (2 + 1.892 + 0 + 0.431 + 0.774) = 8.10$
- $IDCG_6 = 8.69$ ($rel_i = 3, 3, 2, 2, 1, 0$)
- $nDCG_6 = \frac{DCG_6}{IDCG_6} = \frac{8.10}{8.69} = 0.932$

7.5 Statistical Significance Tests

<u>Query</u>	<u>System A</u>	<u>System B</u>	<u>Sign Test</u>	<u>Wilcoxon</u>
1	0.02	0.76	+	+0.74
2	0.39	0.07	-	- 0.32
3	0.16	0.37	+	+0.21
4	0.58	0.21	-	- 0.37
5	0.04	0.02	-	- 0.02
6	0.09	0.91	+	+0.82
7	0.12	0.46	+	+0.34
Average	0.20	0.40	$p=1.0$	$p=0.9375$

Нулевая гипотеза - гипотеза об отсутствии взаимосвязи или корреляции между исследуемыми переменными, об отсутствии различий (однородности) в распределениях (параметрах распределений) двух и/или более выборках.

- H_0 : median difference between the pairs is zero
- H_1 : median difference is not zero.

7.5.1 Sign test

Критерий знаков используется при проверке нулевой гипотезы о равенстве медиан двух непрерывно распределенных случайных величин

Рассмотрим две непрерывно распределенные случайные величины X и Y , и пусть нулевая гипотеза выполняется, то есть их медианы равны. Тогда $p = \mathbb{P}(X > Y) = 0.5$. Иными словами, каждая из случайных величин равновероятно больше другой.

Рассмотрим пару связных выборок $\{(x_1, y_1), \dots, (x_n, y_n)\}$. Будем считать, что в выборке нет элементов, для которых $x_i = y_i$ (иначе уберем эти элементы из выборки). Построим статистику w , равную числу элементов в выборке, при которых $x_i > y_i$. При выполнении нулевой гипотезы, эта величина имеет биномиальное распределение: $w \sim B(n, 0.5)$ с функцией вероятности

$$p_Y(k) \equiv \mathbb{P}(Y = k) = \binom{n}{k} p^k q^{n-k}, \quad k = 0, \dots, n,$$

где $\binom{n}{k} = C_n^k = \frac{n!}{(n-k)! k!}$ – биномиальный коэффициент.

Для применения критерия необходимо вычислить «левый хвост» биномиального распределения до w :

$$b = 2^{-n} \sum_{i=0}^w \binom{n}{i}$$

Согласно критерию, при уровне значимости α : если $b \notin [\alpha/2, 1 - \alpha/2]$, то нулевая гипотеза $p \neq 0.5$ отвергается.

7.5.2 Wilcoxon signed-rank test

The Wilcoxon signed-rank test is a non-parametric statistical hypothesis test used when comparing two related samples or repeated measurements on a single sample to assess whether their population mean ranks differ.

Let N be the sample size, the number of pairs. Thus, there are a total of $2N$ data points. For $i = 1, \dots, N$, let $x_{1,i}$ and $x_{2,i}$ denote the measurements.

- For $i = 1, \dots, N$, calculate $|x_{2,i} - x_{1,i}|$ and $\text{sign}(x_{2,i} - x_{1,i})$.
- Exclude pairs with $|x_{2,i} - x_{1,i}| = 0$. Let N_r be the reduced sample size.
- Order the remaining N_r pairs from smallest absolute difference to largest absolute difference, $|x_{2,i} - x_{1,i}|$.
- Rank the pairs, starting with the smallest as 1. Ties receive a rank equal to the average of the ranks they span. Let R_i denote the rank.
- Calculate the test statistic W , the absolute value of the sum of the signed ranks:

$$W = \left| \sum_{i=1}^{N_r} [\text{sign}(x_{2,i} - x_{1,i}) \cdot R_i] \right|$$

- As N_r increases, the sampling distribution of W converges to a normal distribution. Thus,
 - For $N_r \geq 10$, a z-score can be calculated as $z = \frac{W - 0.5}{\sigma_W}$, $\sigma_W = \sqrt{\frac{N_r(N_r+1)(2N_r+1)}{6}}$. If $z > z_{critical}$ then reject H_0
 - For $N_r < 10$, W is compared to a critical value from a reference table. If $W \geq W_{critical, N_r}$ then reject H_0

Example:

			$x_{2,i} - x_{1,i}$	
i	$x_{2,i}$	$x_{1,i}$	sgn	abs
1	125	110	1	15
2	115	122	-1	7
3	130	125	1	5
4	140	120	1	20
5	140	140		0
6	115	124	-1	9
7	140	123	1	17
8	125	137	-1	12
9	140	135	1	5
10	135	145	-1	10

order by absolute difference

			$x_{2,i} - x_{1,i}$			
i	$x_{2,i}$	$x_{1,i}$	sgn	abs	R_i	$\text{sgn} \cdot R_i$
5	140	140			0	
3	130	125	1	5	1.5	1.5
9	140	135	1	5	1.5	1.5
2	115	122	-1	7	3	-3
6	115	124	-1	9	4	-4
10	135	145	-1	10	5	-5
8	125	137	-1	12	6	-6
1	125	110	1	15	7	7
7	140	123	1	17	8	8
4	140	120	1	20	9	9

Notice that pairs 3 and 9 are tied in absolute value. They would be ranked 1 and 2, so each gets the average of those ranks, 1.5.

$$N_r = 10 - 1 = 9, W = |1.5 + 1.5 - 3 - 4 - 5 - 6 + 7 + 8 + 9| = 9 \\ W < W_{\alpha=0.05,9} = 39 \therefore \text{fail to reject } H_0$$

7.6 Pooling: Avoid Judging all Documents

Pooling strategy:

- Choose a diverse set of ranking methods (TR systems)
- Have each to return top-K documents
- Combine all the top-K sets to form a pool for human assessors to judge
- Other (unjudged) documents are usually assumed to be non-relevant (though they don't have to)

Pooling strategy is okay for comparing systems that contributed to the pool, but problematic for evaluating new systems.

7.7 Summary of TR Evaluation

Evaluation is extremely important:

- TR is an empirically defined problem
- Inappropriate experiment design misguides research and applications
- Make sure to get it right for your research or application

Cranfield evaluation methodology is the main paradigm:

- MAP and nDCG: appropriate for comparing ranking algorithms
- Precision@10docs is easier to interpret from a user's perspective

Not covered:

- A-B Test [Sanderson 10]
- User studies [Kelly 09]

7.8 Recommended reading

- Donna Harman, «Information Retrieval Evaluation. Synthesis Lectures on Information Concepts, Retrieval, and Services», Morgan & Claypool Publishers 2011
- Mark Sanderson, «Test Collection Based Evaluation of Information Retrieval Systems». Foundations and Trends in Information Retrieval 4(4): 247-375 (2010)
- Diane Kelly, «Methods for Evaluating Interactive Information Retrieval Systems with Users». Foundations and Trends in Information Retrieval 3(1-2): 1-224 (2009)

8 Probabilistic Model

8.1 Basic Idea of Probabilistic Model

- Probabilistic models ranking function:

$$f(d, q) = p(R = 1 \mid d, q), R \in \{0, 1\}$$

- **Query Likelihood:** if a user likes document d , how likely would the user enter query q (in order to retrieve d)?
- Assumption: a user formulates a query based on an «imaginary relevant document»:

$$p(R = 1 \mid d, q) \approx p(q \mid d, R = 1)$$

- Basic idea based on user clicks ($R = 1$):

$$f(q, d) = p(R = 1 \mid d, q) = \frac{\text{count}(q, d, R = 1)}{\text{count}(q, d)}$$

- How to compute $p(q \mid d)$? How to compute probability of text in general? → **Language Model**

8.2 Language Model

The term language model (LM) refers to a probabilistic model of text (i.e., it defines a probability distribution over sequences of words).

Uses of a Language Model:

- Representing topics
- Discovering word associations

8.2.1 The Simplest Language Model: Unigram LM

Unigram Language Model = word distribution

- Generate text by generating each word **independently**
- Thus, $p(w_1, w_2 \dots w_n) = p(w_1)p(w_2) \dots p(w_n)$
- Parameters: $\{p(w_i)\} : p(w_1) + \dots + p(w_N) = 1$ (N is vocabulary size)
- Text = sample drawn according to this **word distribution**

Maximum Likelihood (ML) Estimator:

$$p(w \mid \theta) = p(w \mid d) = \frac{c(w, d)}{|d|}$$

- θ - document language model
- $c(w, d)$ - count on word w in document d
- $|d|$ - length of document d

8.2.2 Recommended reading

- Chris Manning and Hinrich Schütze, «Foundations of Statistical Natural Language Processing», MIT Press. Cambridge, MA: May 1999.
- Rosenfeld, R., «Two decades of statistical language modeling: where do we go from here?», Proceedings of the IEEE , vol.88, no.8, pp.1270,1278, Aug. 2000

8.3 Ranking based on Query Likelihood

- Query: $q = w_1 w_2 \dots w_n$
- Vocabulary of the language of the documents: $V = \{w_1, \dots, w_{|V|}\}$
- $c(w, q)$ - count of word w in query q

How likely would we observe this query from this document model?

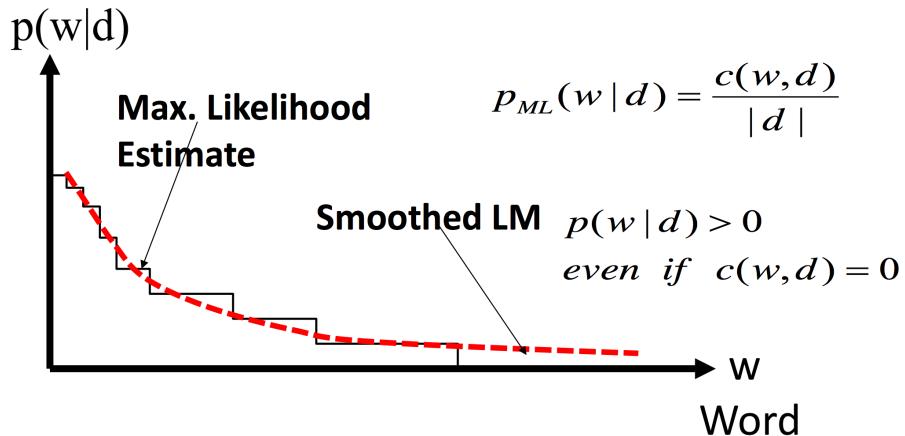
$$p(q | d) = p(w_1 | d) \times \dots \times p(w_n | d)$$

Retrieval problem → estimation of $p(w_i | d)$

$$f(q, d) = \log p(q | d) = \sum_{i=1}^n \log p(w_i | d) = \sum_{w \in V} c(w, q) \cdot \log p(w | d)$$

8.3.1 How to Estimate $p(w | d)$

Smoothing of $p(w | d)$ is necessary for query likelihood:



Key question: what probability should be assigned to an unseen word? Let the probability of an unseen word be proportional to its probability given by a reference LM. One possibility: Reference LM = Collection LM:

$$p(w | d) = \begin{cases} p_{seen}(w | d) & \text{if } w \text{ is seen in } d \\ \alpha_d p(w | C) & \text{otherwise} \end{cases}$$

Discounted ML estimate
if w is seen in d
otherwise

Collection language model

8.3.2 Rewriting the Ranking Function with Smoothing

$$\begin{aligned}
 f(q, d) &= \sum_{w \in V} c(w, q) \log(p(w|d)) \\
 &= \sum_{w \in V, c(w,d) > 0} c(w, q) \log(p_{seen}(w|d)) + \sum_{w \in V, c(w,d) = 0} c(w, q) \log(\alpha_d p(w|\mathcal{C})) \\
 &\quad \text{Word found in document} \quad \text{Modified } p(w|d) \quad \text{Weighting from document collection} \\
 &\quad \text{Word NOT found in document} \quad \text{Weighted probability of word in the collection} \\
 &\quad \text{All words in the collection} \quad \text{Words found in the document} \\
 &\quad \text{This is equal to the count of all the words in the query times log } (\alpha) \\
 f(q, d) &= \sum_{w \in V, c(w,d) > 0} c(w, q) \log \left(\frac{P_{seen}(w|d)}{\alpha_d p(w|\mathcal{C})} \right) + |q| \log(\alpha_d) + \sum_{w \in V} c(w, q) \log(p(w|\mathcal{C}))
 \end{aligned}$$

Smoothing with $p(w|d)$ leads to a general ranking formula for query likelihood with TF-IDF weighting and document length normalization:

$$\log p(q | d) = \sum_{\substack{w_i \in d \\ w_i \in q}} [\log \frac{p_{seen}(w_i | d)}{\alpha_d p(w_i | C)}] + n \log \alpha_d + \sum_{i=1}^n \log p(w_i | C)$$

TF weighting **Doc length normalization**
matched query terms **IDF weighting** **Ignore for ranking**

8.3.3 Smoothing Methods

$$f(q, d) = \sum_{w_i \in d, w_i \in q} c(w_i, q) \log \frac{p_{seen}(w_i | d)}{\alpha_d p(w_i | \mathcal{C})} + n \log \alpha_d$$

Jelinek-Mercer Smoothing

Jelinek-Mercer: Fixed coefficient linear interpolation

$$p_{seen}(w \mid d) = (1 - \lambda) \frac{c(w, d)}{|d|} + \lambda p(w \mid \mathcal{C}), \quad \alpha_d = \lambda, \quad \lambda \in [0, 1]$$

Ranking Function:

$$f_{JM}(q, d) = \sum_{w \in d, w \in q} c(w, q) \log \left(1 + \frac{1 - \lambda}{\lambda} \frac{c(w, d)}{|d| p(w \mid \mathcal{C})} \right)$$

Dirichlet Prior (Bayesian) Smoothing

Dirichlet Prior: Adding pseudo counts; adaptive interpolation

$$\begin{aligned} p_{seen}(w \mid d) &= \frac{c(w, d) + \mu p(w \mid \mathcal{C})}{|d| + \mu} = \frac{|d|}{|d| + \mu} \cdot \frac{c(w, d)}{|d|} + \frac{\mu}{|d| + \mu} \cdot p(w \mid \mathcal{C}) \\ \alpha_d &= \frac{\mu}{|d| + \mu}, \quad \mu \in [0, +\infty) \end{aligned}$$

Ranking Function:

$$f_{DIR}(q, d) = \sum_{w \in d, w \in q} c(w, q) \log \left(1 + \frac{c(w, d)}{\mu p(w \mid \mathcal{C})} \right) + n \log \frac{\mu}{\mu + |d|}$$

8.3.4 Summary of Query Likelihood Probabilistic Model

- Effective ranking functions obtained using pure probabilistic modeling
 - Assumption 1: $Relevance(q, d) = p(R = 1 \mid q, d) \approx p(q \mid d, R = 1) \approx p(q \mid d)$
 - Assumption 2: Query words are generated independently
 - Assumption 3: Smoothing with $p(w \mid \mathcal{C})$
 - Assumption 4: JM or Dirichlet prior smoothing
- Less heuristic compared with VSM
- Many extensions have been made [Zhai 08]

8.3.5 Recommended reading

- ChengXiang Zhai, «Statistical Language Models for Information Retrieval» (Synthesis Lectures Series on Human Language Technologies), Morgan & Claypool Publishers, 2008.

9 Feedback in Text Retrieval

Feedback = learn from examples

9.1 Types of Feedback

- Relevance Feedback

Users make explicit relevance judgments on the initial results. Judgments are reliable, but users don't want to make extra effort.

- Pseudo/Blind/Automatic Feedback.

Top-k initial results are simply assumed to be relevant. Judgments aren't reliable, but no user activity is required.

- Implicit Feedback.

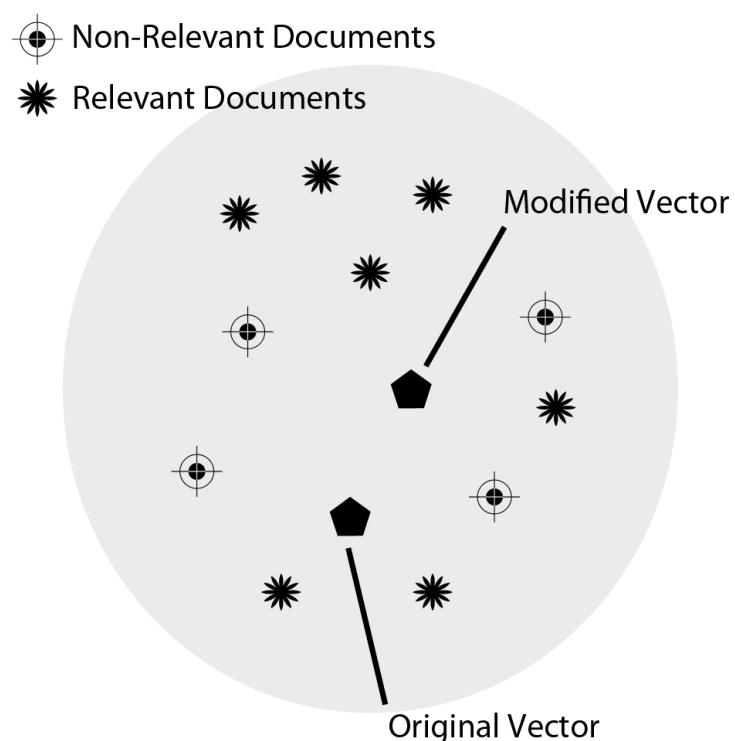
User-clicked docs are assumed to be relevant; skipped ones non-relevant. Judgments aren't completely reliable, but no extra effort from users.

9.2 Feedback in Vector Space Model

How can a TR system learn from examples to improve retrieval accuracy? General method
- query modification:

- Adding new (weighted) terms (query expansion)
- Adjusting weights of old terms

9.2.1 Rocchio Feedback



$$\overrightarrow{Q_m} = (\alpha \cdot \overrightarrow{Q_o}) + \left(\frac{\beta}{|D_r|} \cdot \sum_{\overrightarrow{D_j} \in D_r} \overrightarrow{D_j} \right) - \left(\frac{\gamma}{|D_{nr}|} \cdot \sum_{\overrightarrow{D_k} \in D_{nr}} \overrightarrow{D_k} \right)$$

- $\overrightarrow{Q_m}$ - Modified Query Vector
- $\overrightarrow{Q_o}$ - Original Query Vector
- $\overrightarrow{D_j}$ - Related Document Vector
- $\overrightarrow{D_k}$ - Non-Related Document Vector
- α - Original Query Weight
- β - Related Documents Weight
- γ - Non-Related Documents Weight
- D_r - Set of Related Documents
- D_{nr} - Set of Non-Related Documents

9.2.2 Rocchio in Practice

- Negative (non-relevant) examples are not very important
- Often truncate the vector (i.e., consider only a small number of words that have highest weights in the centroid vector) (efficiency concern)
- Avoid «over-fitting» (keep relatively high weight on the original query weights)
- Can be used for relevance feedback and pseudo feedback (β should be set to a larger value for relevance feedback than for pseudo feedback)
- Usually robust and effective

9.3 Feedback in Language Model

9.3.1 The KL-divergence measure

Given two probability mass functions $p(x)$ and $q(x)$, the Kullback-Leibler divergence (or relative entropy) between p and q is defined as:

$$D(p \parallel q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

9.3.2 Kullback-Leibler (KL) Divergence Retrieval Model

KL-divergence (cross entropy):

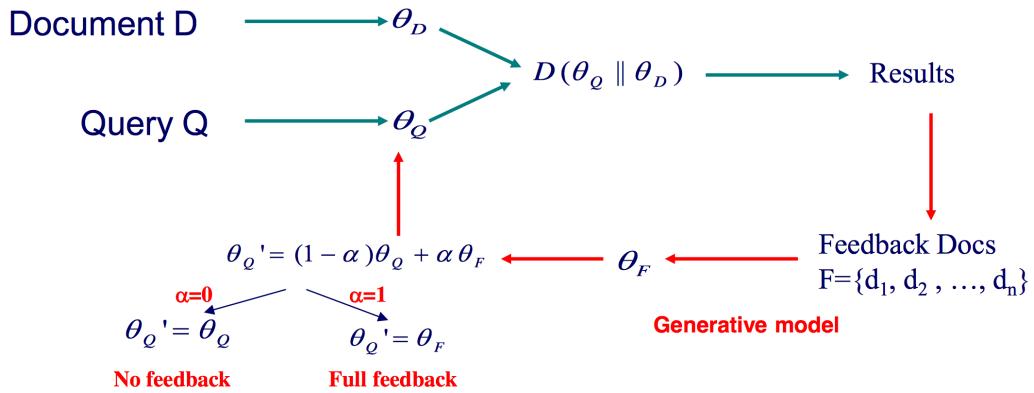
$$f(q, d) = \sum_{\substack{w \in d, \\ p(w|\hat{\theta}_Q) > 0}} p(w \mid \hat{\theta}_Q) \log \frac{p_{seen}(w_i \mid d)}{\alpha_d p(w_i \mid \mathcal{C})} + \log \alpha_d$$

- $\hat{\theta}_Q$ - estimated query unigram language model

If query language model is maximum likelihood $p_{ml}(w | \hat{\theta}_Q) = \frac{c(w, q)}{|q|}$, then KL-divergence becomes a query likelihood retrieval formula:

$$f(q, d) = \sum_{\substack{w_i \in d, \\ w_i \in q}} c(w_i, q) \log \frac{p_{seen}(w_i | d)}{\alpha_d p(w_i | \mathcal{C})} + n \log \alpha_d$$

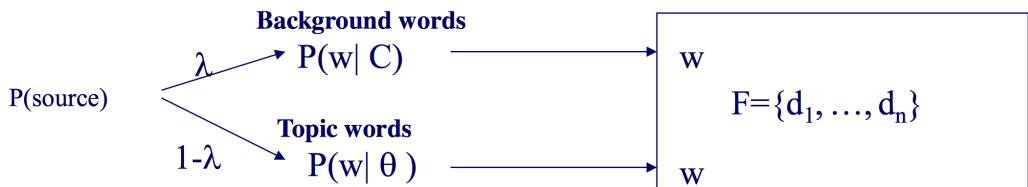
9.3.3 Feedback as Model Interpolation



$$p(w | \hat{\theta}_Q) = (1 - \alpha) p_{ml}(w | \hat{\theta}_Q) + \alpha p(w | \theta_F)$$

- α is a parameter that needs to be set empirically
- $p(w | \theta_F)$ is a feedback model

9.3.4 Generative Mixture Model



$$\log p(\mathcal{F} | \theta_F) = \sum_{i=1}^n \sum_w c(w, d_i) \log ((1 - \lambda)p(w | \theta_F) + \lambda p(w | \mathcal{C}))$$

- $F = \{d_1, \dots, d_n\}$ - the set of feedback documents
- λ - parameter that indicates the amount of «background noise» in the feedback documents

Maximum likelihood estimate of $\hat{\theta}_F$:

$$\hat{\theta}_F = \operatorname{argmax}_{\theta_F} \log p(\mathcal{F} | \theta_F)$$

9.4 Recommended reading

- ChengXiang Zhai, «Statistical Language Models for Information Retrieval» (Synthesis Lectures Series on Human Language Technologies), Morgan & Claypool Publishers, 2008.
- Victor Lavrenko and W. Bruce Croft. 2001. «Relevance based language models». In Proceedings of ACM SIGIR 2011.

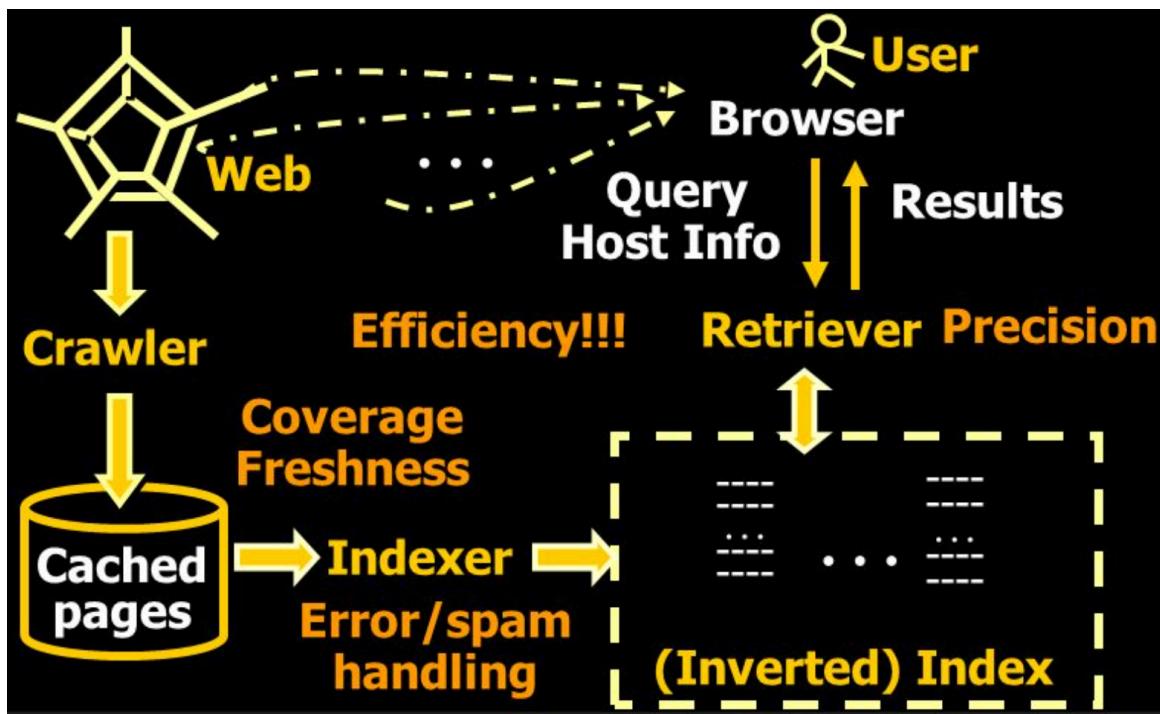
10 Web Search

10.1 Web Search: Challenges & Opportunities

Web search is one of the most important applications of text retrieval.

- Challenges
 - Scalability (the size of the Web, completeness of coverage, many user queries) *to* Parallel indexing & searching (MapReduce)
 - Low quality information and spams *to* Spam detection & Robust ranking
 - Dynamics of the Web (new pages are constantly created, some pages may be updated)
- Opportunities
 - many additional heuristics (e.g., link information, layout) can be leveraged to improve search accuracy *to* Link analysis & multi-feature ranking

10.2 Basic Search Engine Technologies



10.3 Crawler/Spider/Robot

- Building a «toy crawler» is easy
 - Start with a set of “seed pages” in a priority queue
 - Fetch pages from the web
 - Parse fetched pages for hyperlinks; add them to the queue
 - Follow the hyperlinks in the queue

- A real crawler is much more complicated...
 - Robustness (server failure, trap, etc.)
 - Crawling courtesy (server load balance, robot exclusion, etc.)
 - Handling file types (images, PDF files, etc.)
 - URL extensions (cgi script, internal references, etc.)
 - Recognize redundant pages (identical and duplicates)
 - Discover «hidden» URLs (e.g., truncating a long URL)

10.3.1 Major Crawling Strategies

- Breadth-First⁴ is common (balance server load)
- Parallel crawling is natural
- Variation: focused crawling
 - Targeting at a subset of pages (e.g., all pages about «automobiles»)
 - Typically given a query
- How to find new pages (they may not linked to an old page!)
- Incremental/repeated crawling
 - Need to minimize resource overhead
 - Can learn from the past experience (updated daily vs. monthly)
 - Target at : 1) frequently updated pages; 2) frequently accessed pages

10.4 Web Index

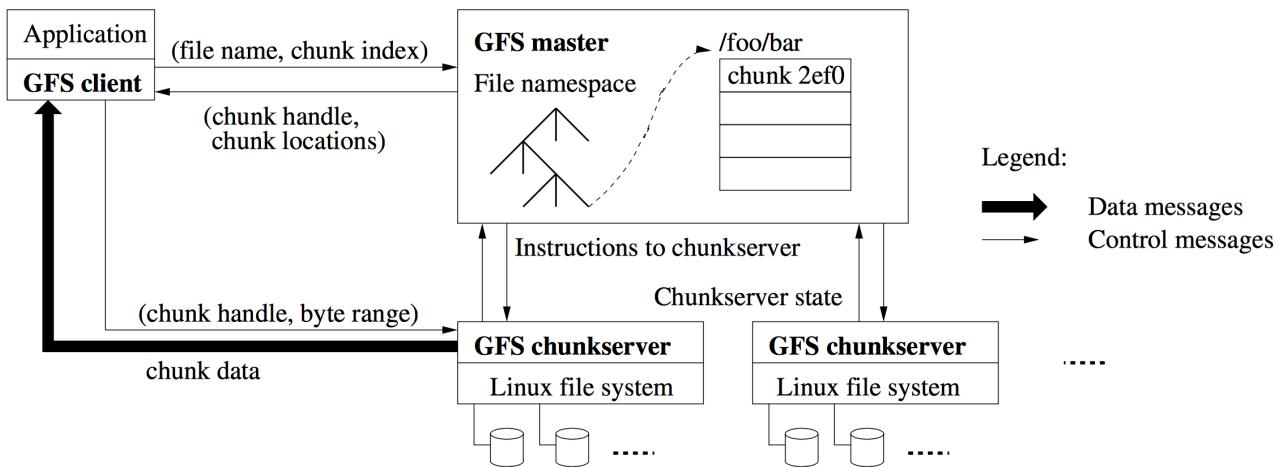
Standard IR techniques are the basis, but insufficient – they lack scalability and efficiency.

Google's contributions:

- Google File System (GFS): distributed file system
- MapReduce: Software framework for parallel computation
- Hadoop: Open source implementation of MapReduce

⁴**Breadth-first search (BFS)** is an algorithm for traversing or searching tree or graph data structures. It starts at the tree root (or some arbitrary node of a graph, sometimes referred to as a «search key») and explores the neighbor nodes first, before moving to the next level neighbors.

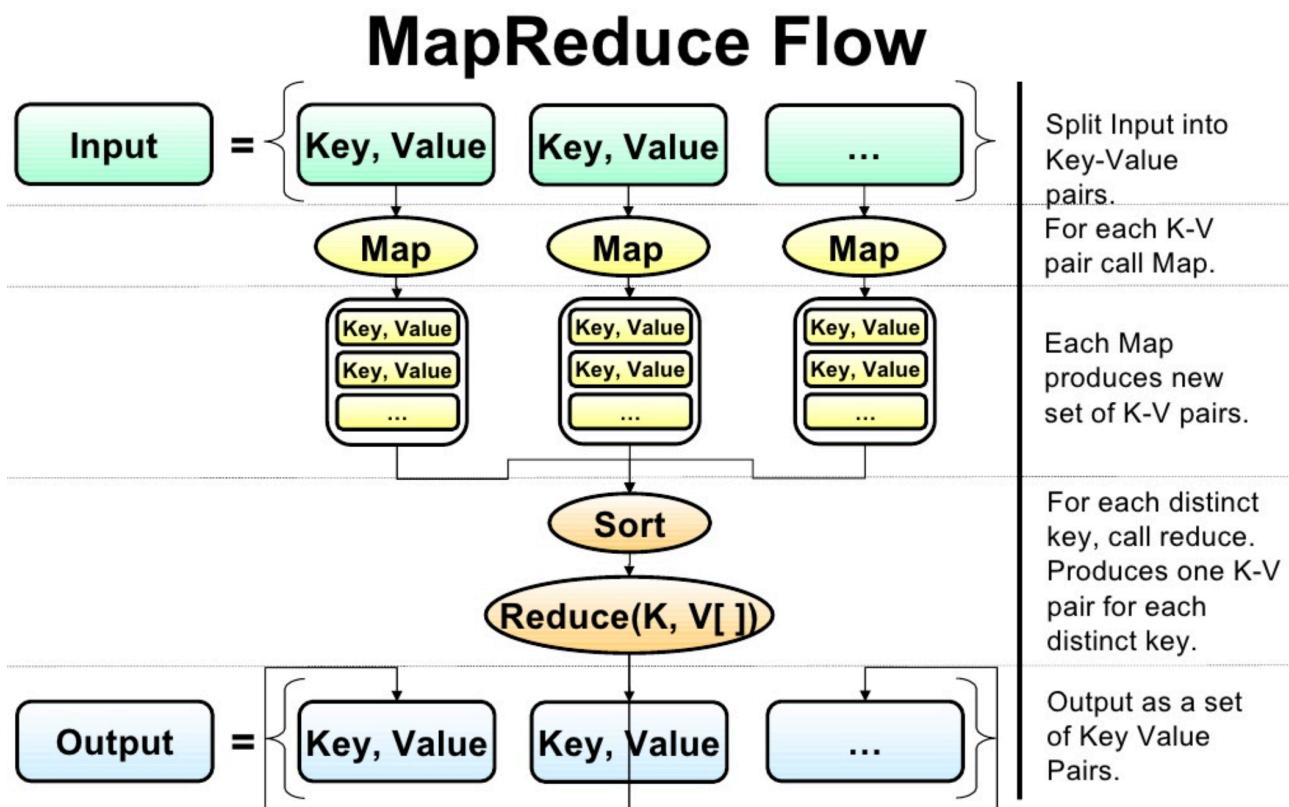
10.4.1 GFS Architecture



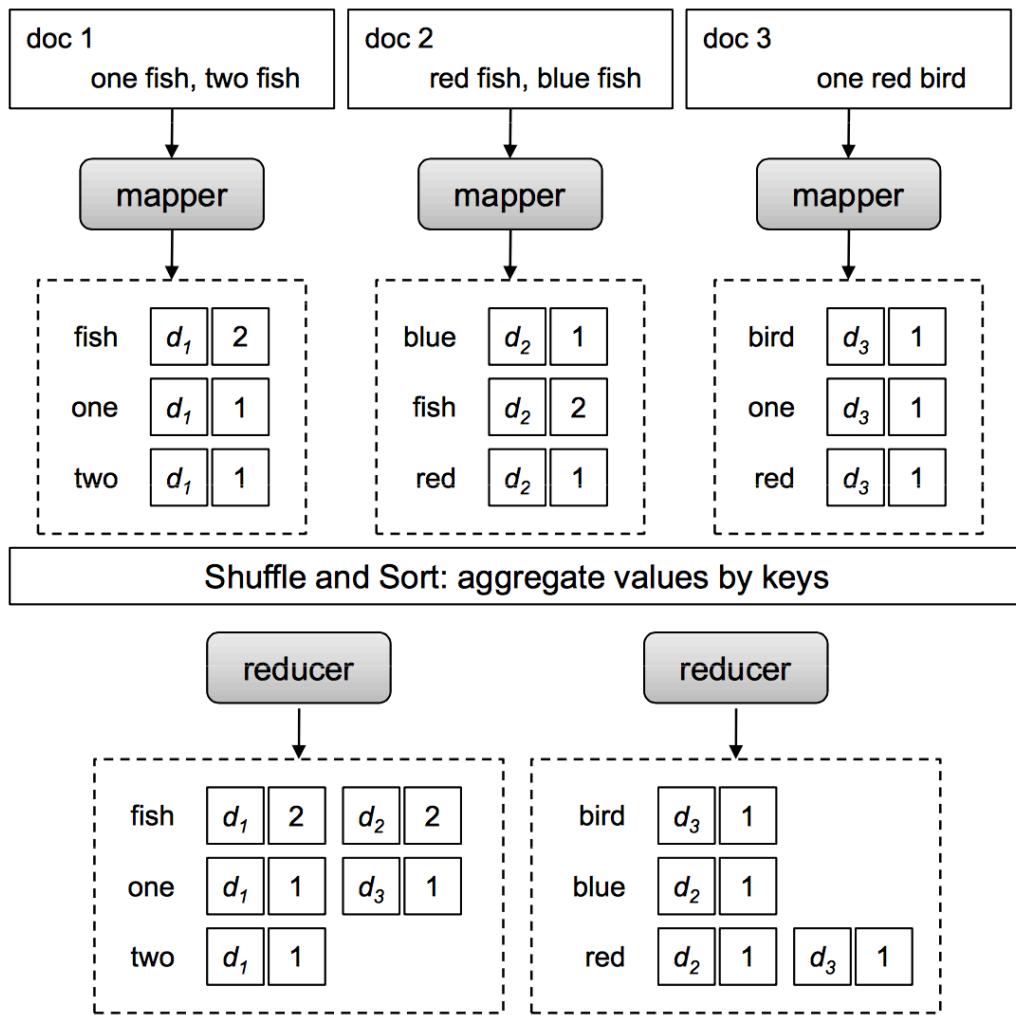
10.4.2 MapReduce: A Framework for Parallel Programming

- Minimize effort of programmer for simple parallel processing tasks
 - Features
 - Hide many low-level details (network, storage)
 - Built-in fault tolerance
 - Automatic load balancing

10.4.3 MapReduce: Computation Pipeline



10.4.4 Inverted Indexing with MapReduce



Algorithm 1 Pseudo-code of the baseline inverted indexing algorithm in MapReduce

```

class Mapper
procedure Map(docid n, doc d)
    H  $\leftarrow$  new AssociativeArray
    for all term t  $\in$  doc d do
        H{t}  $\leftarrow$  H{t} + 1
    for all term t  $\in$  H do
        Emit(term t, posting <n, H{t}>)

class Reducer
procedure Reduce(term t, postings [<n1, f1>, <n2, f2> ...])
    P  $\leftarrow$  new List
    for all posting <a, f>  $\in$  postings [<n1, f1>, <n2, f2> ...] do
        Append(P, <a, f>)
    Sort(P)
    Emit(term t, postings P)

```

10.5 Link Analysis

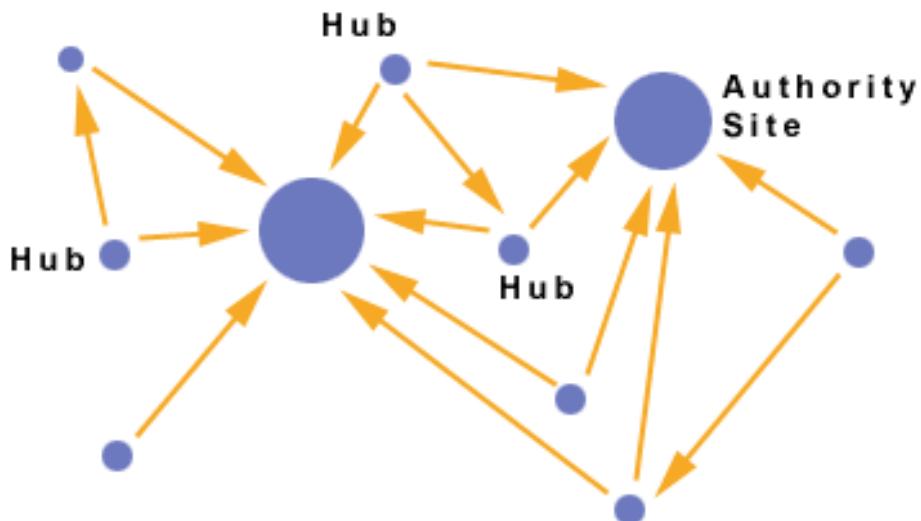
10.5.1 Ranking Algorithms for Web Search

- Standard IR models apply but aren't sufficient
 - Different information needs (search for a particular web page instead of text information)
 - Documents have additional information (links, layout)
 - Information quality varies a lot
- Major extensions
 - Exploiting links to improve scoring
 - Exploiting clickthroughs for massive implicit feedback
 - In general, rely on machine learning to combine all kinds of features

10.5.2 Exploiting Inter-Document Links

Description of a link (**«anchor text»**) is a summary and a query example for a target document.

An **authority** is a web page containing valuable information with respect to a specific subject. A **hub** is a web page not actually authoritative in the information it holds, but contains useful links toward an authoritative page – basically, advertising the authoritative web page.

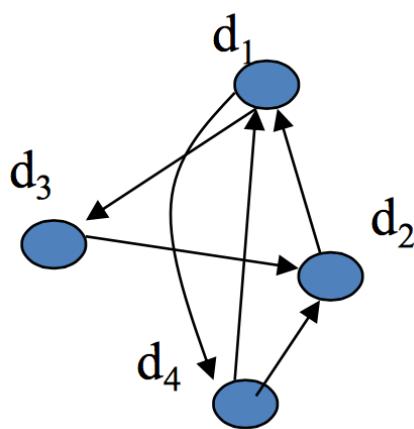


10.5.3 PageRank: Capturing Page «Popularity»

- Intuitions
 - Links are like citations in literature
 - A page that is cited often can be expected to be more useful in general

- PageRank is essentially «citation counting», but improves over simple counting
 - Consider «indirect citations» (being cited by a highly cited paper counts a lot...)
 - Smoothing of citations (every page is assumed to have a non-zero pseudo citation count)
- PageRank can also be interpreted as random surfing (thus capturing popularity)

10.5.4 The PageRank Algorithm



Transition matrix:

$$M = \begin{pmatrix} 0 & 0 & 1/2 & 1/2 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 \end{pmatrix}$$

M_{ij} - probability of going from d_i to d_j :

$$\forall i : \sum_{j=1}^N M_{ij} = 1$$

«Equilibrium Equation»:

$$p_{t+1}(d_j) = (1 - \alpha) \sum_{i=1}^N M_{ij} p_t(d_i) + \alpha \sum_{i=1}^N \frac{1}{N} p_t(d_i),$$

where

- N - number of pages
- $p_t(d_i)$ - probability of visiting page d_i at time t
- α - probability of jumping to a random page

For a converged⁵ state we can drop the time index:

$$p(d_j) = \sum_{i=1}^N \left[\frac{1}{N} \alpha + (1 - \alpha) M_{ij} \right] p(d_i)$$

and come to this equation:

$$\vec{p} = \left(\alpha I + (1 - \alpha) M \right)^T \vec{p}, \text{ where } I_{ij} = \frac{1}{N},$$

which can be solved with an iterative algorithm starting with initial value $p(d) = 1/N$.

⁵Converge - сходиться

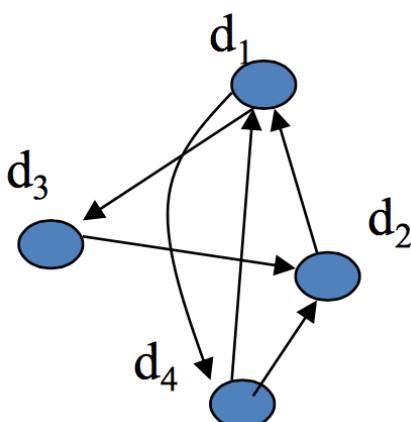
10.5.5 PageRank in Practice

- Computation can be quite efficient since M is usually sparse
- Normalization doesn't affect ranking, leading to some variants of the formula
- The zero-outlink problem: $p(d_i)$'s don't sum to 1
 - One possible solution = page-specific damping factor⁶ ($\alpha = 1.0$ for a page with no outlink)
- Many extensions (e.g., topic-specific PageRank)
- Many other applications (e.g., social network analysis)

10.5.6 HITS: Capturing Authorities & Hubs

- Intuitions
 - Pages that are widely cited are good authorities
 - Pages that cite many other pages are good hubs
- The key idea of HITS (Hypertext-Induced Topic Search)
 - Good authorities are cited by good hubs
 - Good hubs point to good authorities
 - Iterative reinforcement...
- Many applications in graph/network analysis

10.5.7 The HITS Algorithm



Algorithm:

- Initial values: $a(d_i) = h(d_i) = 1$

- Iteration:

$$\begin{cases} h(d_i) = \sum_{d_j \in OUT(d_i)} a(d_j) \\ a(d_i) = \sum_{d_j \in IN(d_i)} h(d_j) \end{cases}$$

- Iteration in matrix form:

«Adjacency matrix»⁷:

$$A = \begin{pmatrix} 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{pmatrix}$$

$$\begin{cases} \vec{h} = A \vec{a} = AA^T \vec{h} \\ \vec{a} = A^T \vec{h} = A^T A \vec{a} \end{cases}$$

- Normalize: $\sum_i a(d_i)^2 = \sum_i h(d_i)^2 = 1$

⁶Damping factor - коэффициент затухания

⁷«Adjacency matrix» ≈ матрица связей («соседства»)

10.6 Learning to Rank

10.6.1 How Can We Combine Many Features?

- Given a query-doc pair (Q, D) , define various kinds of features $X_i(Q, D)$
- Examples of feature:
 - the number of overlapping terms,
 - BM25 score of Q and D ,
 - $p(Q | D)$,
 - PageRank of D ,
 - $p(Q | D_i)$, where D_i may be anchor text or big font text,
 - «does the URL contain ‘?’»
 - etc.
- Hypothesize $p(R = 1 | Q, D) = s(X_1(Q, D), \dots, X_n(Q, D), \lambda)$, where λ is a set of parameters
- Learn λ by fitting function s with training data, i.e., 3-tuples like $(D, Q, 1)$ (D is relevant to Q) or $(D, Q, 0)$ (D is non-relevant to Q)

10.6.2 Regression-Based Approaches

Logistic Regression:

- $X_i(Q, D)$ is feature
- β 's are parameters

$$\log \frac{p(R = 1 | Q, D)}{1 - p(R = 1 | Q, D)} = \beta_0 + \sum_{i=1}^N \beta_i X_i$$
$$p(R = 1 | Q, D) = \frac{1}{1 + \exp(-\beta_0 - \sum_{i=1}^N \beta_i X_i)}$$

Estimate β 's by maximizing the likelihood of training data:

$$\vec{\beta}^* = \underset{\vec{\beta}}{\operatorname{argmax}} p\left(\{(Q_1, D_{11}, R_{11}), (Q_1, D_{12}, R_{12}), \dots, (Q_n, D_{n1}, R_{n1}), \dots\}\right) \quad (1)$$

10.6.3 More Advanced Learning Algorithms

- Attempt to directly optimize a retrieval measure (e.g. MAP, nDCG)
 - More difficult as an optimization problem
 - Many solutions were proposed [Liu 09]
- Can be applied to many other ranking problems beyond search
 - Recommender systems
 - Computational advertising
 - Summarization

10.6.4 Recommended reading

- Tie-Yan Liu. «Learning to Rank for Information Retrieval». Foundations and Trends in Information Retrieval 3, 3 (2009): 225-331.
- Hang Li. «A Short Introduction to Learning to Rank», IEICE Trans. Inf. & Syst. E94-D, 10 (Oct. 2011): n.p.

10.7 Future of Web Search

10.7.1 Next Generation Search Engines

- More specialized/customized (vertical search engines):
 - Special group of users (community engines, e.g., Citeseer)
 - Personalized (better understanding of users)
 - Special genre/domain (better understanding of documents)
- Learning over time (evolving)
- Integration of search, navigation, and recommendation/filtering (full-fledged information management)
- Beyond search to support tasks (e.g., shopping)
- Many opportunities for innovations!

10.7.2 Future Intelligent Information Systems

- Search ⇒ Access ⇒ Mining ⇒ Task support
- Keyword queries ⇒ Search history ⇒ Complete user model
- Bag of words ⇒ Entities-Relations ⇒ Knowledge presentation

11 Recommender Systems

11.1 Recommender \approx Filtering System

- Stable & long term interest, dynamic info source
- System must make a delivery decision immediately as a document «arrives»

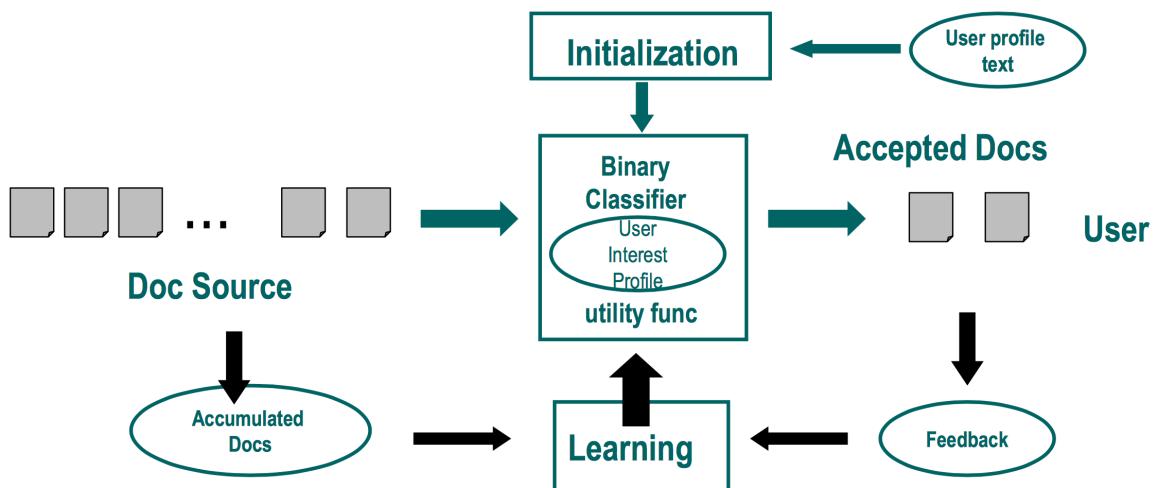
11.1.1 Basic Filtering Question

Will User U Like Item X?

- Look at what items U likes, and then check if X is similar
 - Item similarity \Rightarrow **content-based filtering**
- item Look at who likes X, and then check if U is similar
 - User similarity \Rightarrow **collaborative filtering**

11.2 Content-Based Filtering

11.2.1 A Typical Content-Based Filtering System



Example of linear utility:

$$\text{Linear Utility} = 3 \times \#\text{good} - 2 \times \#\text{bad}$$

11.2.2 Three Basic Problems in Content-Based Filtering

- Making filtering decision (Binary classifier)
 - Doc text, profile text \rightarrow yes/no
- Initialization
 - Initialize the filter based on only the profile text or very few examples
- Learning from

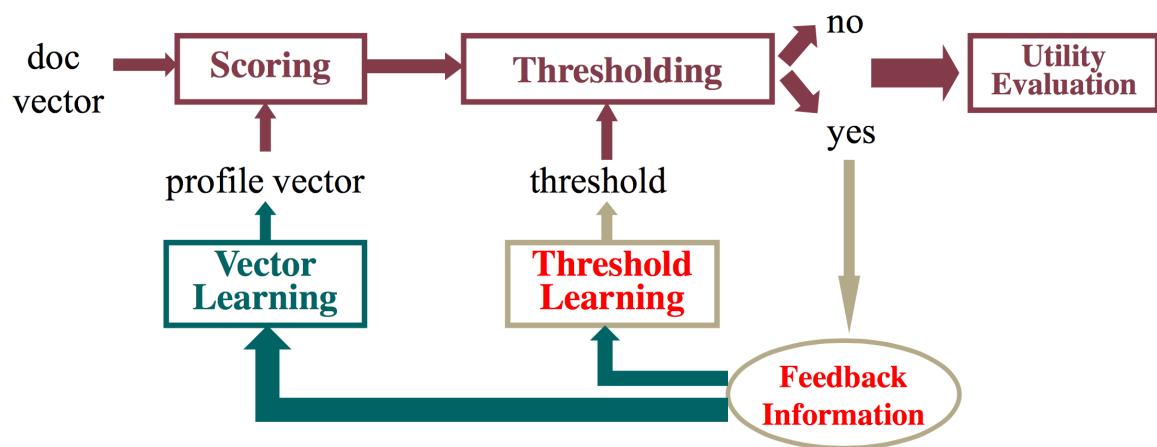
- Limited relevance judgments (only on «yes» docs) – Accumulated documents
- All trying to maximize the utility

11.2.3 Extend a Retrieval System for Information Filtering

Content-based recommender system can be built based on a search engine system:

- «Reuse» retrieval techniques to score documents
- Use a score threshold for filtering decision
- Learn to improve scoring with traditional feedback
- New approaches to threshold setting and learning

11.2.4 A General Vector-Space Approach



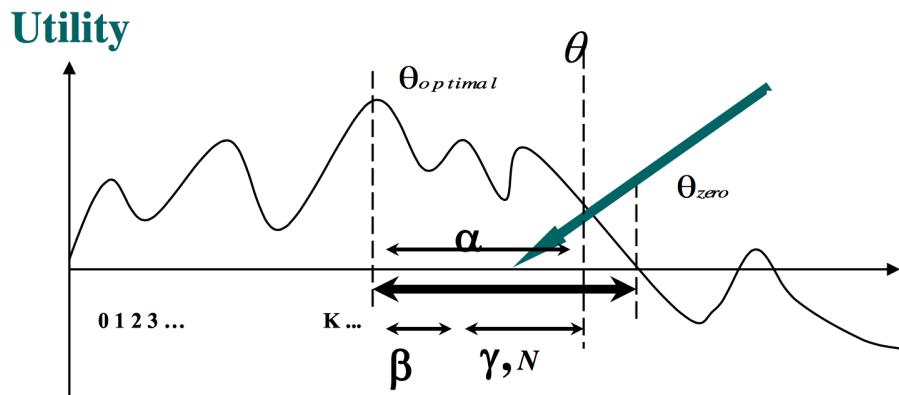
11.2.5 Difficulties in Threshold Learning

- Censored data (judgments only available on delivered documents)
- Little or none labeled data
- Exploration vs. Exploitation

11.2.6 Empirical Utility Optimization

- Basic idea
 - Compute the utility on the training data for each candidate score threshold
 - Choose the threshold that gives the maximum utility on the training data set
- Difficulty: Biased training sample!
 - We can only get an upper bound for the true optimal threshold
 - Could a discarded item be possibly interesting to the user?
- Solution:
 - Heuristic adjustment (lowering) of threshold

11.2.7 Beta-Gamma Threshold Learning



Encourage exploration of threshold up to θ_{zero} :

$$\theta = \alpha \times \theta_{zero} + (1 - \alpha) \times \theta_{optimal}$$

The more examples, the less exploration (closer to $\theta_{optimal}$):

$$\alpha = \beta + (1 - \beta) \times e^{-N\gamma}, \text{ where } \beta, \gamma \in [0, 1]$$

- Pros:
 - Explicitly addresses exploration-exploitation tradeoff («Safe» exploration)
 - Arbitrary utility (with appropriate lower bound)
 - Empirically effective
- Cons:
 - Purely heuristic
 - Zero utility lower bound often too conservative

11.3 Collaborative Filtering

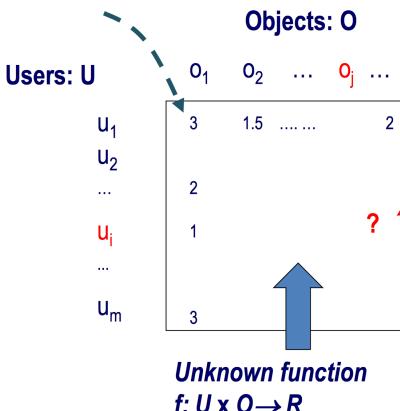
11.3.1 What is Collaborative Filtering (CF)?

- Making filtering decisions for an individual user based on the judgments of other users
- Inferring individual's interest/preferences from that of other similar users
- User similarity can be judged based on their similarity in preferences on a common set of items

CF assumptions

- Users with the same interest will have similar preferences
- Users with similar preferences probably share the same interest
- Sufficiently large number of user preferences are available (if not, there will be a «cold start» problem)

11.3.2 Memory-based Approaches to Collaboration Filtering Problem



- Assume we have m users (u_i) and n objects o_j
- x_{ij} : rating of object o_j by user u_i
- n_i : average rating of all objects by user u_i
- Normalized ratings: $v_{ij} = x_{ij} - n_i$

Prediction of rating of object o_j by user u_a :

$$\hat{v}_{aj} = \frac{\sum_{i=1}^m v_{a,i} v_{ij}}{\sum_{i=1}^m v_{a,i}}, \quad \hat{x}_{aj} = \hat{v}_{aj} + n_a$$

11.3.3 User Similarity Measures

- Pearson correlation

$$w_p(a, j) = \frac{\sum_j (x_{aj} - n_a)(x_{ij} - n_i)}{\sqrt{\sum_j (x_{aj} - n_a)^2 \sum_j (x_{ij} - n_i)^2}}$$

- Cosine measure

$$w_c(a, j) = \frac{\sum_j x_{aj} x_{ij}}{\sqrt{\sum_j x_{aj}^2 \sum_j x_{ij}^2}}$$

11.3.4 Improving User Similarity Measures

- Dealing with missing values: set to default ratings
 - average ratings
 - iterate unknown ratings calculations followed by similarity measure calculation
- Inverse User Frequency (IUF): similar to IDF

11.3.5 Summary

- Filtering/Recommendation is «easy»
 - The user's expectation is low
 - Any recommendation is better than none
- Filtering is «hard»
 - Must make a binary decision, though ranking is also possible
 - Data sparseness (limited feedback information)
 - «Cold start» (little information about users at the beginning)

11.3.6 Recommended reading

- Francesco Ricci, Lior Rokach, Bracha Shapira, Paul B. Kantor. «**Recommender Systems Handbook**». Springer 2011.

12 Course Summary

12.1 Recommended reading

- [Synthesis Digital Library](#) has many excellent short books/long tutorials on relevant topics
 - [Information Concepts, Retrieval and Services](#)
 - [Human Language Technology](#)
 - [Artificial Intelligence & Machine Learning](#)
- Journals: ACM TOIS, IRJ, IPM, ...
- Conferences: SIGIR, CIKM, ECIR, WSDM, WWW, KDD, ACL, ...
- [More info](#)
- [Search User Interface](#), by Marti Hearst, Cambridge University Press, 2009

12.2 Main Techniques for Harnessing Big Text Data: Text Retrieval + Text Mining

