

# 3-1: Scoring, Term Weighting and the Vector Space Model

Fatemeh Azimzadeh

# This lecture; IIR Sections 6.2

- ▶ Ranked retrieval
- ▶ Scoring documents
- ▶ Term frequency
- ▶ Collection statistics
- ▶ Weighting schemes

# Ranked retrieval

- ▶ Thus far, our queries have all been Boolean.
- ▶ Documents either match or don't.
- ▶ Good for expert users with precise understanding of their needs and the collection.
- ▶ Also good for applications: Applications can easily consume 1000s of results.
- ▶ Not good for the majority of users.
- ▶ Most users incapable of writing Boolean queries (or they are, but they think it's too much work).
- ▶ Most users don't want to wade through 1000s of results.
- ▶ This is particularly true of web search.

# Problem with Boolean search: feast or famine

- ▶ Boolean queries often result in either too few (=0) or too many (1000s) results.
- ▶ Query 1: “standard user dlink 650” → 200,000 hits
- ▶ Query 2: “standard user dlink 650 no card found”: 0 hits
- ▶ It takes a lot of skill to come up with a query that produces a manageable number of hits.
- ▶ AND gives too few; OR gives too many

# Ranked retrieval models

- ▶ Rather than a set of documents satisfying a query expression, in ranked retrieval, the system returns an ordering over the (top) documents in the collection for a query
- ▶ Free text queries: Rather than a query language of operators and expressions, the user's query is just one or more words in a human language
- ▶ In principle, there are two separate choices here, but in practice, ranked retrieval has normally been associated with free text queries and vice versa

# Feast or famine: not a problem in ranked retrieval

- ▶ When a system produces a ranked result set, large result sets are not an issue
  - ▶ Indeed, the size of the result set is not an issue
  - ▶ We just show the top  $k$  ( $\approx 10$ ) results
  - ▶ We don't overwhelm the user
- ▶ Premise: the ranking algorithm works

# Feast or famine: not a problem in ranked retrieval

- ▶ When a system produces a ranked result set, large result sets are not an issue
  - ▶ Indeed, the size of the result set is not an issue
  - ▶ We just show the top  $k$  ( $\approx 10$ ) results
  - ▶ We don't overwhelm the user
- ▶ Premise: the ranking algorithm works





# Scoring as the basis of ranked retrieval

- ▶ We wish to return in order the documents most likely to be useful to the searcher
- ▶ How can we rank-order the documents in the collection with respect to a query?
- ▶ Assign a score - say in  $[0, 1]$  - to each document
- ▶ This score measures how well document and query “match”.

# Query-document matching scores

- ▶ We need a way of assigning a score to a query/document pair
- ▶ **Let's start with a one-term query**
- ▶ If the query term does not occur in the document: score should be 0
- ▶ **The more frequent the query term in the document, the higher the score (should be)**
- ▶ We will look at a number of alternatives for this.

# Take 1: Jaccard coefficient

- ▶ A commonly used measure of overlap of two sets  $A$  and  $B$
- ▶  $\text{jaccard}(A,B) = |A \cap B| / |A \cup B|$
- ▶  $\text{jaccard}(A,A) = 1$
- ▶  $\text{jaccard}(A,B) = 0$  if  $A \cap B = 0$
- ▶  $A$  and  $B$  don't have to be the same size.
- ▶ Always assigns a number between 0 and 1.

# Jaccard coefficient: Scoring example

- ▶ What is the query-document match score that the Jaccard coefficient computes for each of the two documents below?
- ▶ Query: *ides of march*
- ▶ Document 1: *caesar died in march*
- ▶ Document 2: *the long march*

# Issues with Jaccard for scoring

- ▶ It doesn't consider *term frequency* (how many times a term occurs in a document)
- ▶ Rare terms in a collection are more informative than frequent terms. Jaccard doesn't consider this information
- ▶ We need a more sophisticated way of normalizing for length
- ▶ Later in this lecture, we'll use
- ▶ . . . instead of  $|A \cap B| / |A \cup B|$  (Jaccard) for length normalization.

$$|A \cap B| / \sqrt{|A \cup B|}$$

# Recall (Lecture 1): Binary term-document incidence matrix

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0

Each document is represented by a binary vector  $\in \{0,1\}^{|V|}$

# Term-document count matrices

- Consider the number of occurrences of a term in a document:
  - Each document is a count vector in  $\mathbb{N}^V$ : a column below

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	157	73	0	0	0	0
Brutus	4	157	0	1	0	0
Caesar	232	227	0	2	1	1
Calpurnia	0	10	0	0	0	0
Cleopatra	57	0	0	0	0	0
mercy	2	0	3	5	5	1
worser	2	0	1	1	1	0

# Bag of words model

- ▶ Vector representation doesn't consider the ordering of words in a document
- ▶ *John is quicker than Mary* and *Mary is quicker than John* have the same vectors
- ▶ This is called the bag of words model.
- ▶ In a sense, this is a step back: The positional index was able to distinguish these two documents.
- ▶ We will look at “recovering” positional information later in this course.
- ▶ For now: bag of words model



# Coordination level matching

- ▶ Coordination level matching (CLM) is a straightforward approach to bag-of-words queries
  - Idea: Documents whose index records have  $n$  different terms in common with the query are more relevant than documents with  $n - 1$  different terms held in common
- The coordination level (also called “size of overlap”) between a query  $Q$  and a document  $D$  is the number of terms they have in common
- How to answer a query?
  1. Sort the document collection by coordination level
  2. Return the head of this sorted list to the user (say, the best 20 documents)

# Example

- ▶ Document1 = {step, man, mankind}
- ▶ Document2 = {step, man, China}
- ▶ Document3 = {step, mankind}
  
- Query1 = {man, mankind}
  - ▶ Result: 1. Document1 (2)      2. Document2 , Document3 (1)
- Query2 = {China, man, mankind}
  - ▶ Result: 1. Document1 , Document2 (2)      2. Document3(1)

# Term frequency tf

- ▶ The term frequency  $tf_{t,d}$  of term  $t$  in document  $d$  is defined as the number of times that  $t$  occurs in  $d$ .
- ▶ We want to use tf when computing query-document match scores. But how?
- ▶ Raw term frequency is not what we want:
  - ▶ A document with 10 occurrences of the term is more relevant than a document with 1 occurrence of the term.
  - ▶ But not 10 times more relevant.
- ▶ Relevance does not increase proportionally with term frequency.

NB: frequency = count in IR

# TF Weighting

- ▶ Idea: A term is more important if it occurs more frequently in a document
- ▶ Formulas: Let  $f(t,d)$  be the frequency count of term  $t$  in doc  $d$ 
  - ▶ Raw TF:  $TF(t,d) = f(t,d)$
  - ▶ Log TF:  $TF(t,d) = \log f(t,d)$
  - ▶ Maximum frequency normalization:  
$$TF(t,d) = 0.5 + 0.5 * f(t,d) / \text{MaxFreq}(d)$$
- ▶ Normalization of TF is very important!

# Log-frequency weighting

- ▶ The log frequency weight of term  $t$  in  $d$  is  $w_{t,d} = \begin{cases} 1 + \log_{10} \text{tf}_{t,d}, & \text{if } \text{tf}_{t,d} > 0 \\ 0, & \text{otherwise} \end{cases}$

- ▶  $0 \rightarrow 0, 1 \rightarrow 1, 2 \rightarrow 1.3, 10 \rightarrow 2, 1000 \rightarrow 4$ , etc.
- ▶ Score for a document-query pair: sum over terms  $t$  in both  $q$  and  $d$ :
- ▶ score

$$= \sum_{t \in q \cap d} (1 + \log \text{tf}_{t,d})$$

- ▶ The score is 0 if none of the query terms is present in the document.

# Document frequency

- ▶ Rare terms are more informative than frequent terms
  - ▶ Recall stop words
- ▶ Consider a term in the query that is rare in the collection (e.g., *arachnocentric*)
- ▶ A document containing this term is very likely to be relevant to the query *arachnocentric*
- ▶ → We want a high weight for rare terms like *arachnocentric*.

# Document frequency, continued

- ▶ Frequent terms are less informative than rare terms
- ▶ Consider a query term that is frequent in the collection (e.g., *high*, *increase*, *line*)
- ▶ A document containing such a term is more likely to be relevant than a document that doesn't
- ▶ But it's not a sure indicator of relevance.
- ▶ → For frequent terms, we want high positive weights for words like *high*, *increase*, and *line*
- ▶ But lower weights than for rare terms.
- ▶ We will use document frequency (df) to capture this.

# TF Normalization

- ▶ Why?
  - ▶ Document length variation
  - ▶ “Repeated occurrences” are less informative than the “first occurrence”
- ▶ Two views of document length
  - ▶ A doc is long because it uses more words
  - ▶ A doc is long because it has more contents
- ▶ Generally penalize long doc, but avoid over-penalizing



# idf weight

- ▶  $df_t$  is the document frequency of  $t$ : the number of documents that contain  $t$ 
  - ▶  $df_t$  is an inverse measure of the informativeness of  $t$
  - ▶  $df_t \leq N$
- ▶ We define the idf (inverse document frequency) of  $t$  by

$$idf_t = \log_{10} (N/df_t)$$

- ▶  $N$  : total number of docs
- ▶ We use  $\log (N/df_t)$  instead of  $N/df_t$  to “dampen” the effect of idf.

Idea: A term is more discriminative if it occurs only in fewer documents

idf example, suppose  $N = 1$  million

term	$df_t$	$idf_t$
calpurnia	1	
animal	100	
sunday	1,000	
fly	10,000	
under	100,000	
the	1,000,000	

$$idf_t = \log_{10} (N/df_t)$$

There is one idf value for each term  $t$  in a collection.

# Collection vs. Document frequency

- ▶ The collection frequency of  $t$  is the number of occurrences of  $t$  in the collection, counting multiple occurrences.
- ▶ Example:

Word	Collection frequency	Document frequency
<i>insurance</i>	10440	3997
<i>try</i>	10422	8760

# tf-idf weighting

- ▶ The tf-idf weight of a term is the product of its tf weight and its idf weight.

- ▶  $w_{t,d} = \log(1 + \text{tf}_{t,d}) \times \log_{10}(N / \text{df}_t)$

- ▶ Best known weighting scheme in information retrieval
  - ▶ Note: the “-” in tf-idf is a hyphen, not a minus sign!
  - ▶ Alternative names: tf.idf, tf x idf
- ▶ Increases with the number of occurrences within a document
- ▶ Increases with the rarity of the term in the collection

## Score for a document given a query

$$\text{Score}(q, d) = \sum_{t \in q \cap d} \text{tf.idf}_{t,d}$$

- ▶ There are many variants
  - ▶ How “tf” is computed (with/without logs)
  - ▶ Whether the terms in the query are also weighted
  - ▶ ...

# Binary $\rightarrow$ count $\rightarrow$ weight matrix

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	5.25	3.18	0	0	0	0.35
Brutus	1.21	6.1	0	1	0	0
Caesar	8.59	2.54	0	1.51	0.25	0
Calpurnia	0	1.54	0	0	0	0
Cleopatra	2.85	0	0	0	0	0
mercy	1.51	0	1.9	0.12	5.25	0.88
worser	1.37	0	0.11	4.15	0.25	1.95

Each document is now represented by a real-valued vector of tf-idf weights  $\in \mathbb{R}^{|V|}$

# The Notion of Relevance

