An Application of Vector Space: Ranking Web Search Results



Query Engine's operations

- Process query
- Look-up the index
- Retrieve list of documents
- Order documents
- Prepare results page

MOTIVATING QUESTION:

Given a large list of documents that match a query, how to **order** them according to their relevance?

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 $\in \{0,1\}^{|V|}$

Why Ranking Search Results?

- Thus far, our queries have all been Boolean.
 Documents either match or don't. No "top-10" results.
- Boolean queries exhibit the "feast-or-famine" problem:
 Adding a term to a query goes from '000's of results to no results
- Boolean queries are good for expert users with precise understanding of their needs and the collection.
 - Also good for applications: they can easily handle 1000s of results.
- Boolean are 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.

Solution: Score Documents to Rank

Given document d_i ; Given query q

- Calculate score(q, d_i), usually a number [0, 1]
 score measures how well document and query "match"
- Rank documents in decreasing order of *score(q, d_i)*

Assign 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: what the score of the document should be? _____
- If the query term occurs 5 times in d1 and 50 times in d2, which document should have a higher score? _____
- We will look at a number of alternatives for this.

But first: The Bag of Words Model

- "Bag of words" Model:
 - Document = bag of [unordered] words
 - d1 = "John is quicker than Mary" and d2 = "Mary is quicker than John" have the same words.
 - Is this a good idea?
- In a sense, this is a step back: The positional index was able to distinguish these two. But that's ok for now.

Term-document count matrices

6

8

- Consider the number of occurrences of a term in a document:
 - Each document is a count vector in ℕ^v: a column

	Antony and cleopatra	Julius Caesai	The rempest	namiet	Otheno	Machell
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

Idea: Query terms that appear more often in a document, match it better

5



Log-frequency weighting

• The log frequency weight of term t in d is

 $w_{t,d} = \begin{cases} 1 + \log_{10} tf_{t,d}, & \text{if } tf_{t,d} > 0\\ 0, & \text{otherwise} \end{cases}$

- $0 \rightarrow 0, 1 \rightarrow 1, 2 \rightarrow 1.3, 10 \rightarrow 2, 1000 \rightarrow 4$, etc.
- Score for a document-query pair: sum over terms t in both q and d:
- score = $\sum_{t \in q \cap d} (1 + \log t f_{t,d})$
- The score is 0 if none of the query terms is present in the document.

Why using just $tf_{t,d}$ is not good?

- Bigger documents have more terms, thus their score is larger.
 Will need to fix this: "length normalization"
- A document having 10 occurrences of a term is not 10 times better than another that has it once.
 - Will need to fix this
- All terms are equally important.
 How frequency can indicate importance?

Postulate:

A word that appears in (almost) every document, is not very important i.e., it has no discriminatory power.

10

12

Method 2: Weights according to rarity

- 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*.
- df_t document frequency for term t (= num of docs contain t)
- idf_t inverse document frequency for term t

 $idf_t = \log \frac{N}{dt}$

11

idf example

suppose N = 1,000,000 documents

term	df _t	idf _t
arachnocentric	1	
capricious	100	
sunday	1,000	
person	10,000	
under	100,000	
the	1,000,000	

$$\mathrm{idf}_t = \log_{10} \left(N/\mathrm{df}_t \right)$$

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

Method 3: Better tf.idf weighting

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

 $\mathbf{W}_{t,d} = (1 + \log t \mathbf{f}_{t,d}) \times \log_{10}(N/d\mathbf{f}_t)$

- Best known weighting scheme in information retrieval
 - Alternative names: tf x idf, tf-idf (with hyphen, not minus)
- What happens to tf.idf weight when...
 - the **number of occurrences** within a document increases?
 - the **rarity of the term** in the collection increases?

15

Binary \rightarrow count \rightarrow tf.idf weight matrix

Term-Document 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 $\in \mathbb{R}^{|V|}$ of |V| terms

How to rank results for query: "In the mercy of Caesar"

Ranking of documents for a query **DOCUMENTS** are vectors • So we have a |V|-dimensional vector space • Terms are axes of the space • Documents are points or vectors in this space Score(q,d) = $\sum_{t \in q \cap d} \text{tf.idf}_{t,d}$ • Very high-dimensional: tens of millions of dimensions d_3 when you apply this to a web search engine • These are very sparse vectors: • Rank by summing tf.idf weights. most entries are zero. But how do we normalize lengths? We need to think in terms of vector arithmetic... 17 18 **QUERIES** are also vectors Euclidean distance is not a good idea • Key idea 1: Represent gueries as GOSSIP d_2 vectors in the space d_1 The Euclidean distance • Key idea 2: Rank documents between q and d_2 is large according to their proximity to query in this space even though the - proximity = similarity of vectors $d_3 \times$ distribution of terms in the proximity ≈ inverse of "distant" query *q* and the θ • How to determine vector distribution of terms proximity? in d_2 are very similar. JEALOUS Let's try shortest distance b/w vectors - Euclidean distance? d_4 19 20



- Effect on the two documents d and d²: they have identical vectors after length-normalization.
 - Long and short documents now have comparable weights

 d_i is the tf-idf weight of term *i* in the document

 $\cos(\vec{q}, \vec{d})$ is the cosine similarity of \vec{q} and \vec{d} ... or,

equivalently, the cosine of the angle between \vec{q} and \vec{d} .

Cosine similarity illustrated

• The effect of length normalization



3 documents example contd.

Log frequency weighting

After length normalization

term	SaS	PaP	WH	term	SaS	PaP
affection	3.06	2.76	2.30	affection	0.789	0.832
jealous	2.00	1.85	2.04	jealous	0.515	0.555
gossip	1.30	0	1.78	gossip	0.335	C
wuthering	0	0	2.58	wuthering	0	C

$\cos(SaS, PaP) \approx$

 $0.789 \times 0.832 + 0.515 \times 0.555 + 0.335 \times 0.0 + 0.0 \times 0.0$ ≈ 0.94 $\cos(SaS,WH) \approx 0.79$ $\cos(PaP,WH) \approx 0.69$

Cosine similarity amongst 3 documents

How similar are the novels SaS: Sense and Sensibility PaP: Pride and Prejudice, and WH: Wuthering Heights?

term	SaS	PaP	WH
affection	115	58	20
jealous	10	7	11
gossip	2	0	6
wuthering	0	0	38

Term frequencies (counts)

Note: To simplify this example, we don't do idf weighting. 26

TERM-DOCUM	MENT MATRIX	C			
term vectors	SaS	PaP	Wh Hei	query = jeal	ous gossip
affection	115	58	20	0	
jealous	10	7	11	1	
gossip	2	0	6	1	
EUCLIDIAN N	NORMS				
euclidian nor	115.451288	58.4208867	23.6008474	1.41421356	
NORMALIZED	D TERM-DOCU	MENT MATRI	x		
norm term veo	SaS	PaP	Wh Hei		
affection	0.99609109	0.99279561	0.8474272	0	
jealous	0.08661662	0.11982016	0.46608496	0.70710678	
gossip	0.01732332	0	0.25422816	0.70710678	
SIMILARITIE	S CALCULAT	ION	SIMILARITIE	ES TO QUERY	
sim(SAS, PaP)	0.99929328		sim(q,SaS)	0.07349664	
sim(SaS, WH)	0.88888946		sim(q,PaP)	0.08472565	
sim(PaP,WH)	0.89716838		sim(q,WH)	0.50933829	best match
SaS is much n	nore similar to	PaP			
than it is to W	н.				

Calculation of similarities

WH

0.524

0.465

0.405 0.588

Alternative 1: Jaccard coefficient

- A commonly used measure of overlap of two sets A and B
 - $\operatorname{jaccard}(A,B) = |A \cap B| / |A \cup B|$
 - jaccard(A,A) = ____
 - jaccard(A,B) = ____ if A \cap B = 0
- Jaccard() takes values between ____ and ____.
- Example: Query q: ides of march
 - What is the score that the Jaccard coefficient computes for each of the two documents below?
 - $\underline{d}1$: caesar died in march
 - $\underline{d}2$: the long march
- Score(q,d1) = ____ Score(q,d2) = ____ More relevant: ____

29

Issues with Jaccard for scoring

- + Easy to implement
- It doesn't consider how many times a term occurs in a document
 - Ignores "term frequency"
- - It considers every word to be equally important. But rare terms in a collection are more informative than frequent terms.

- Ignores "term rarity"

- Long documents have an advantage.
 We need a sophisticated way of normalizing for length
 - Ignores "length normalization"

30