

Information Storage and Retrieval

CSCE 670

Texas A&M University

Department of Computer Science & Engineering

Instructor: Prof. James Caverlee

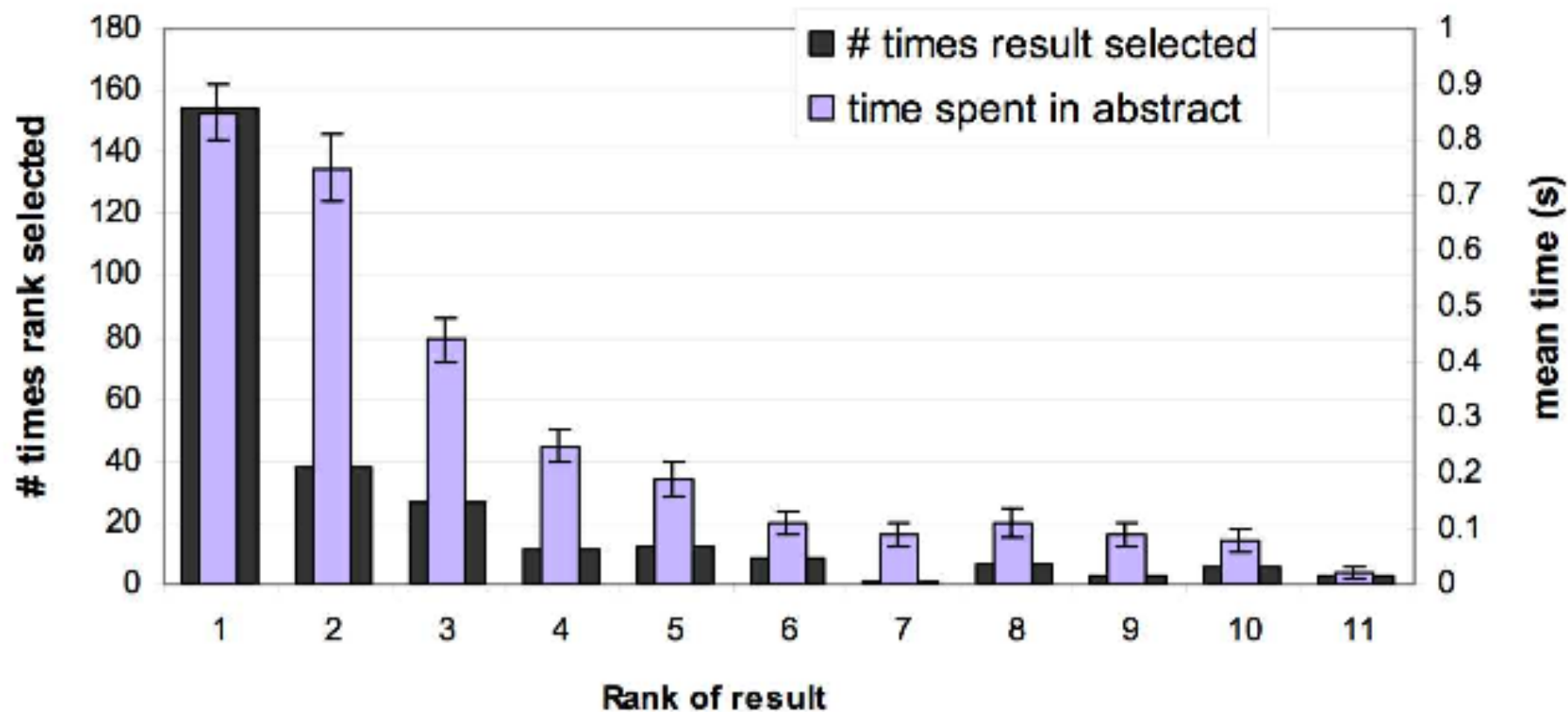
Vector Space Retrieval

6 February 2018

Today

- Ranking is important (evidence from Google)
- Vector space retrieval
 - TF
 - IDF
 - Cosine

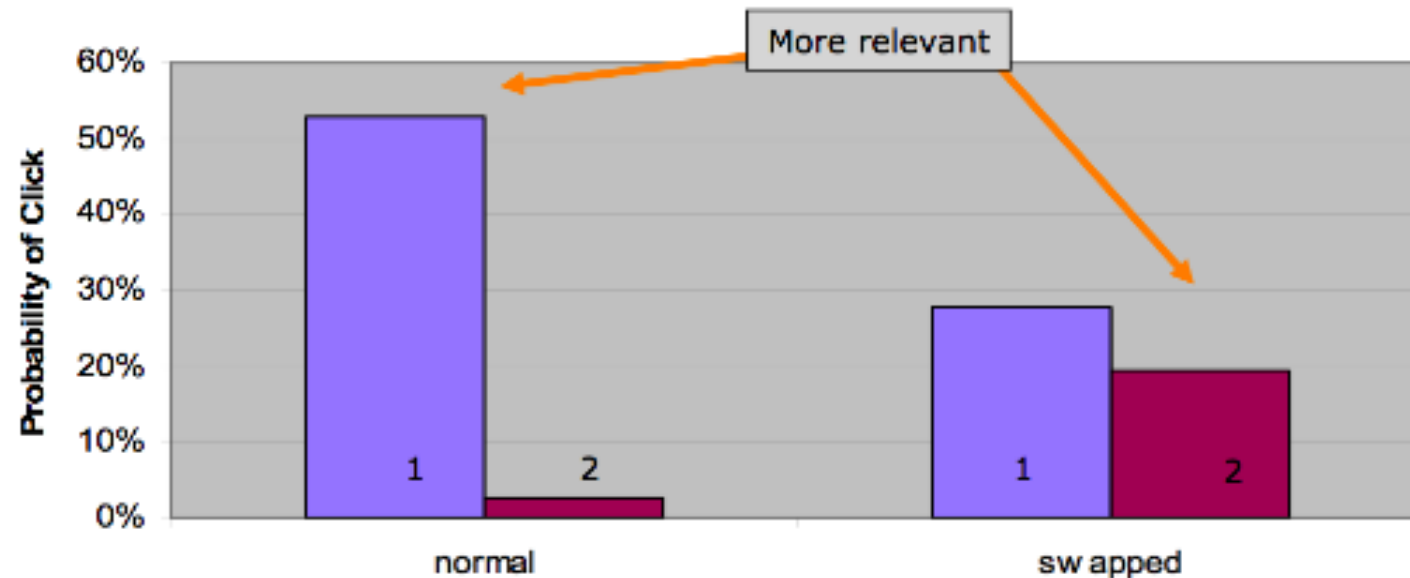
Looking vs. Clicking



- Users view results one and two more often / thoroughly
- Users click most frequently on result one

Presentation bias – reversed results

- Order of presentation influences where users look **AND** where they click



2009 / 2010 Introducing LETOR 4.0 Datasets

Tao Qin and Tie-Yan Liu

{taoqin,tyliu}@microsoft.com

Microsoft Research Asia

| Column in Output | Description |
|------------------|---|
| 1 | TF(Term frequency) of body |
| 2 | TF of anchor |
| 3 | TF of title |
| 4 | TF of URL |
| 5 | TF of whole document |
| 6 | IDF(Inverse document frequency) of body |
| 7 | IDF of anchor |
| 8 | IDF of title |
| 9 | IDF of URL |
| 10 | IDF of whole document |
| 11 | TF*IDF of body |
| 12 | TF*IDF of anchor |
| 13 | TF*IDF of title |
| 14 | TF*IDF of URL |
| 15 | TF*IDF of whole document |
| 16 | DL(Document length) of body |
| 17 | DL of anchor |

| | |
|----|----------------|
| 21 | BM25 of body |
| 22 | BM25 of anchor |
| 23 | BM25 of title |
| 24 | BM25 of URL |

Term Frequency

Binary incidence matrix

| | Anthony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth ... |
|-----------|-----------------------------|------------------|----------------|--------|---------|----------------|
| ANTHONY | 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 as a binary vector $\in \{0, 1\}^{|M|}$.

Count matrix

| | Anthony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth ... |
|-----------|-----------------------------|------------------|----------------|--------|---------|----------------|
| ANTHONY | 157 | 73 | 0 | 0 | 0 | 1 |
| BRUTUS | 4 | 157 | 0 | 2 | 0 | 0 |
| CAESAR | 232 | 227 | 0 | 2 | 1 | 0 |
| CALPURNIA | 0 | 10 | 0 | 0 | 0 | 0 |
| CLEOPATRA | 57 | 0 | 0 | 0 | 0 | 0 |
| MERCY | 2 | 0 | 3 | 8 | 5 | 8 |
| WORSER | 2 | 0 | 1 | 1 | 1 | 5 |
| ... | | | | | | |

Each document is now represented as a count vector
 $\in \mathbb{N}^M$.

Bag of words model

- We do not consider the **order** of words in a document.
- *John is quicker than Mary and Mary is quicker than John* are represented the same way.
- This is called a **bag of words model**.

Hans Peter Luhn



TF (1957)

The weight of a term that occurs in a document is simply proportional to the term frequency.

IBM Researcher

Foundational work in the 1950s

(invented an algorithm to checksum your credit card numbers!)

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 because:
- A document with **$tf = 10$** occurrences of the term is more relevant than a document with **$tf = 1$** occurrence of the term.
- But not 10 times more relevant.
- Relevance does not increase proportionally with term frequency.

Instead of raw frequency: Log frequency weighting

- The log frequency weight of term t in d is defined as follows

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

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

Frequency in document vs. frequency in collection

- In addition, to term frequency (the frequency of the term in the document) ...
- ...we also want to use the frequency of the term in the collection for weighting and ranking.

Desired weight for rare terms

- Rare terms are more informative than frequent terms.
- 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.
- → We want **high weights for rare terms** like ARACHNOCENTRIC.

Desired weight for frequent terms

- Frequent terms are less informative than rare terms.
- Consider a term in the query that is **frequent** in the collection (e.g., GOOD, INCREASE, LINE).
- A document containing this term is more likely to be relevant than a document that doesn't ...
- ... but words like GOOD, INCREASE and LINE are not sure indicators of relevance.
- → **For frequent terms** like GOOD, INCREASE and LINE, we want positive weights ...
- ... but **lower weights** than for rare terms.

Document frequency

- We want **high weights for rare terms** like ARACHNOCENTRIC.
- We want **low (positive) weights for frequent words** like GOOD, INCREASE and LINE.
- We will use **document frequency** to factor this into computing the matching score.
- The document frequency is **the number of documents in the collection that the term occurs in.**

Karen Spärck Jones



IDF (1972)

A STATISTICAL INTERPRETATION OF TERM SPECIFICITY AND ITS APPLICATION IN RETRIEVAL

KAREN SPARCK JONES

University of Cambridge Computer Laboratory

The exhaustivity of document descriptions and the specificity of index terms are usually regarded as independent. It is suggested that specificity should be interpreted statistically, as a function of term use rather than of term meaning. The effects on retrieval of variations in term specificity are examined, experiments with three test collections showing in particular that frequently-occurring terms are required for good overall performance. It is argued that terms should be weighted according to collection frequency, so that matches on less frequent, more specific, terms are of greater value than matches on frequent terms. Results for the test collections show that considerable improvements in performance are obtained with this very simple procedure.

INDEX TERM WEIGHTING

KAREN SPARCK JONES

Computer Laboratory, University of Cambridge, Cambridge, England

(Received 17 August 1973)

Summary—Various approaches to index term weighting have been investigated. In particular, claims have been made for the value of statistically-based indexing in automatic retrieval systems. The paper discusses the logic of different types of weighting, and describes experiments testing weighting schemes of these types. The results show that one type of weighting leads to material performance improvements in quite different collection environments.

- (1) The occurrence(s) of a term in a document are significant:
if they are treated as equally significant, one term will have the same weight in different documents, and different terms will have the same weight in one document;
if they are treated as variably significant, one term may have different weights in different documents, and different terms may have different weights in one document.
- (2) The occurrence(s) of a term in a short document are more significant than its occurrence(s) in a long document:
assuming equal significance under 1, one term may have different weights in different documents, but different terms will have the same weight in one document;
otherwise different terms may have different weights in one document.
- (3) The occurrence(s) of a rare term in a document are more significant than the occurrence(s) of a frequent one:
assuming equal significance under 1, one term will have the same weight in different documents, but different terms will have different weights in the same document;
otherwise the same term may have different weights in different documents.

idf weight

- df_t is the document frequency, the number of documents that t occurs in.
- df_t is an inverse measure of the **informativeness** of term t .
- We define the **idf weight** of term t as follows:

$$idf_t = \log_{10} \frac{N}{df_t}$$

(N is the number of documents in the collection.)

- idf_t is a measure of the **informativeness** of the term.
- $[\log N/df_t]$ instead of $[N/df_t]$ to “dampen” the effect of idf
- Note that we use the log transformation for both term frequency and document frequency.

Examples for idf

- Compute idf_t using the formula: $\text{idf}_t = \log_{10} \frac{1,000,000}{\text{df}_t}$

| term | df_t | idf_t |
|-----------|---------------|----------------|
| calpurnia | 1 | 6 |
| animal | 100 | 4 |
| sunday | 1000 | 3 |
| fly | 10,000 | 2 |
| under | 100,000 | 1 |
| the | 1,000,000 | 0 |

TF-IDF Weighting

tf-idf weighting

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

$$w_{t,d} = (1 + \log \text{tf}_{t,d}) \cdot \log \frac{N}{\text{df}_t}$$

- tf-weight
- idf-weight
- 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

Summary: tf-idf

- Assign a tf-idf weight for each term t in each document d :

$$w_{t,d} = (1 + \log \text{tf}_{t,d}) \cdot \log \frac{N}{\text{df}_t}$$

- The tf-idf weight ...
 - ...increases with the number of occurrences within a document. (term frequency)
 - ...increases with the rarity of the term in the collection. (inverse document frequency)

The Vector Space Model

Gerard Salton

Father of Information Retrieval



TF-IDF (1975)
Vector Space Model

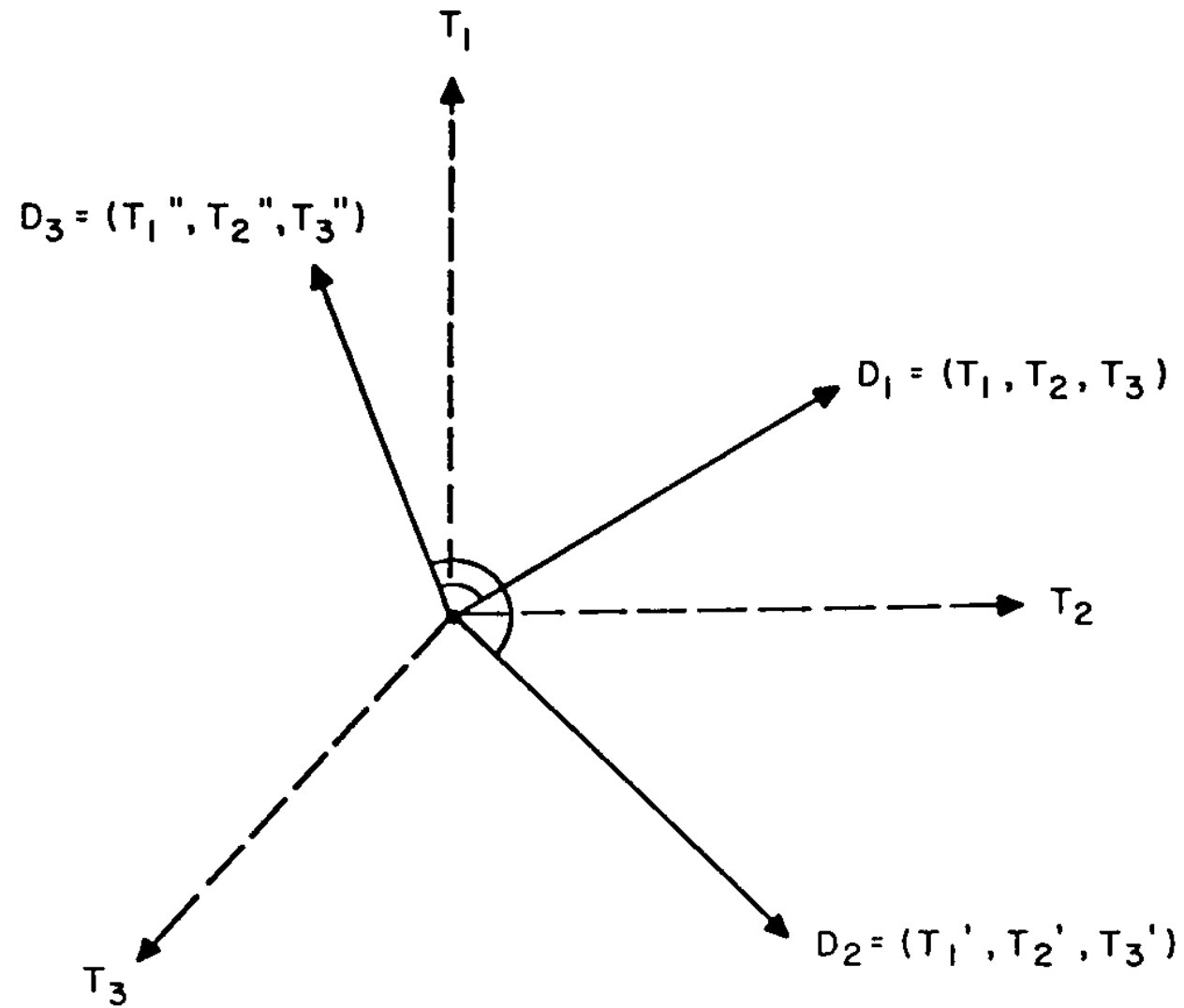
SMART system (at Cornell)

A Vector Space Model for Automatic Indexing

G. Salton, A. Wong
and C. S. Yang
Cornell University

In a document retrieval, or other pattern matching environment where stored entities (documents) are compared with each other or with incoming patterns (search requests), it appears that the best indexing (property) space is one where each entity lies as far away from the others as possible; in these circumstances the value of an indexing system may be expressible as a function of the density of the object space; in particular, retrieval performance may correlate inversely with space density. An approach based on space density computations is used to choose an optimum indexing vocabulary for a collection of documents. Typical evaluation results are shown, demonstrating the usefulness of the model.

Fig. 1. Vector representation of document space.



Binary incidence matrix

| | Anthony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth ... |
|-----------|-----------------------------|------------------|----------------|--------|---------|----------------|
| ANTHONY | 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 as a binary vector $\in \{0, 1\}^{|M|}$.

Count matrix

| | Anthony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth ... |
|-----------|-----------------------------|------------------|----------------|--------|---------|----------------|
| ANTHONY | 157 | 73 | 0 | 0 | 0 | 1 |
| BRUTUS | 4 | 157 | 0 | 2 | 0 | 0 |
| CAESAR | 232 | 227 | 0 | 2 | 1 | 0 |
| CALPURNIA | 0 | 10 | 0 | 0 | 0 | 0 |
| CLEOPATRA | 57 | 0 | 0 | 0 | 0 | 0 |
| MERCY | 2 | 0 | 3 | 8 | 5 | 8 |
| WORSER | 2 | 0 | 1 | 1 | 1 | 5 |
| ... | | | | | | |

Each document is now represented as a count vector
 $\in \mathbb{N}^M$.

Binary \rightarrow count \rightarrow weight matrix

| | Anthony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth ... |
|-----------|-----------------------------|------------------|----------------|--------|---------|----------------|
| ANTHONY | 5.25 | 3.18 | 0.0 | 0.0 | 0.0 | 0.35 |
| BRUTUS | 1.21 | 6.10 | 0.0 | 1.0 | 0.0 | 0.0 |
| CAESAR | 8.59 | 2.54 | 0.0 | 1.51 | 0.25 | 0.0 |
| CALPURNIA | 0.0 | 1.54 | 0.0 | 0.0 | 0.0 | 0.0 |
| CLEOPATRA | 2.85 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| MERCY | 1.51 | 0.0 | 1.90 | 0.12 | 5.25 | 0.88 |
| WORSER | 1.37 | 0.0 | 0.11 | 4.15 | 0.25 | 1.95 |
| ... | | | | | | |

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

Documents as vectors

- Each document is now represented as a real-valued vector of tf-idf weights $\in \mathbb{R}^{|V|}$.
- So we have a $|V|$ -dimensional real-valued vector space.
- Terms are **axes** of the space.
- Documents are **points** or **vectors** in this space.
- Very high-dimensional: tens of millions of dimensions when you apply this to web search engines
- Each vector is very sparse - most entries are zero.

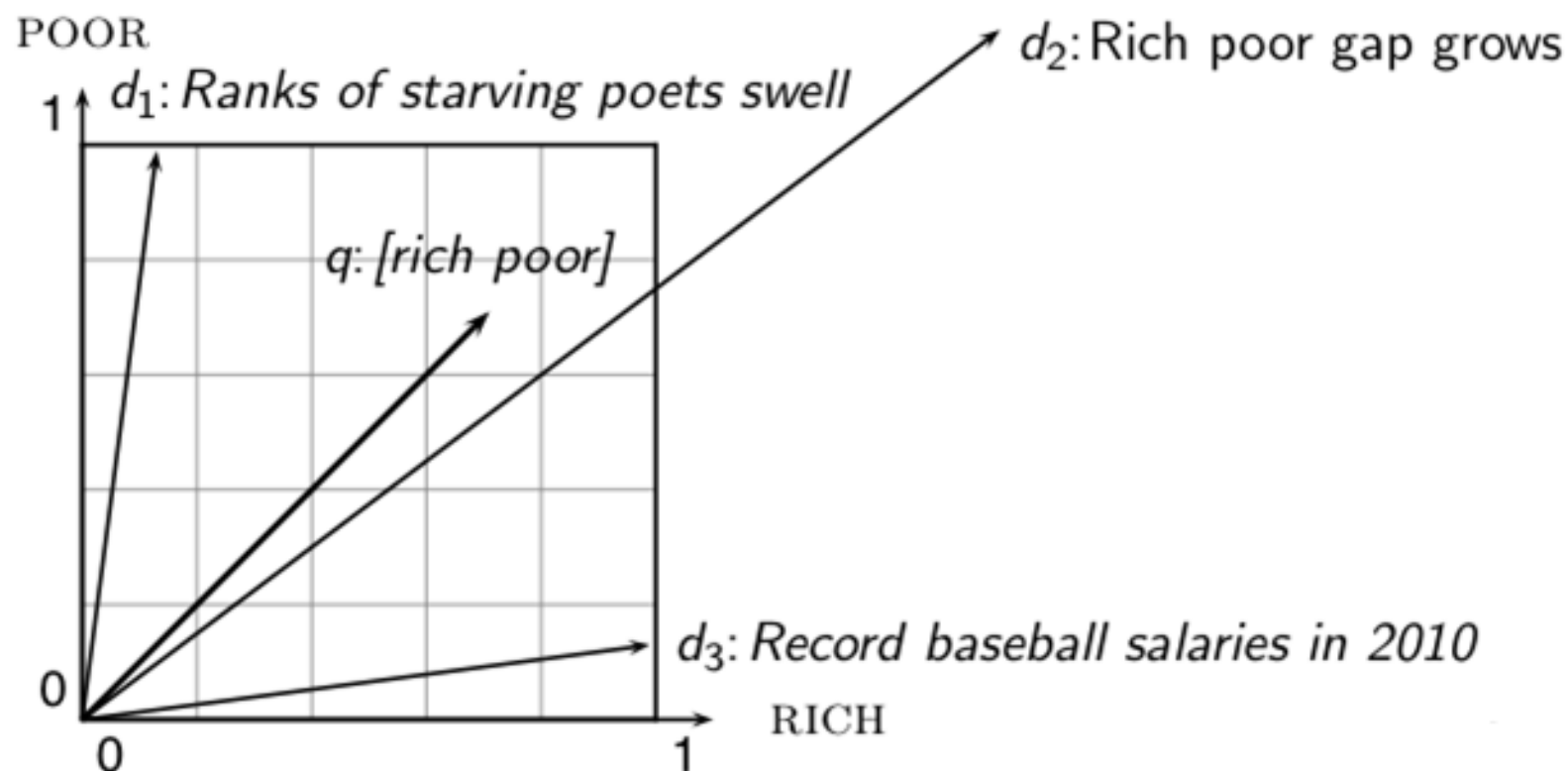
Queries as vectors

- Key idea 1: do the same for queries: represent them as vectors in the high-dimensional space
- Key idea 2: Rank documents according to their proximity to the query
- proximity = similarity
- proximity \approx negative distance
- Recall: We're doing this because we want to get away from the you're-either-in-or-out, feast-or-famine Boolean model.
- Instead: rank relevant documents higher than nonrelevant documents

How do we formalize vector space similarity?

- First cut: (negative) distance between two points
- (= distance between the end points of the two vectors)
- Euclidean distance?
- Euclidean distance is a bad idea ...
- ... because Euclidean distance is large for vectors of different lengths.

Why distance is a bad idea



The Euclidean distance of \vec{q} and \vec{d}_2 is large although the distribution of terms in the query q and the distribution of terms in the document d_2 are very similar.

Questions about basic vector space setup?

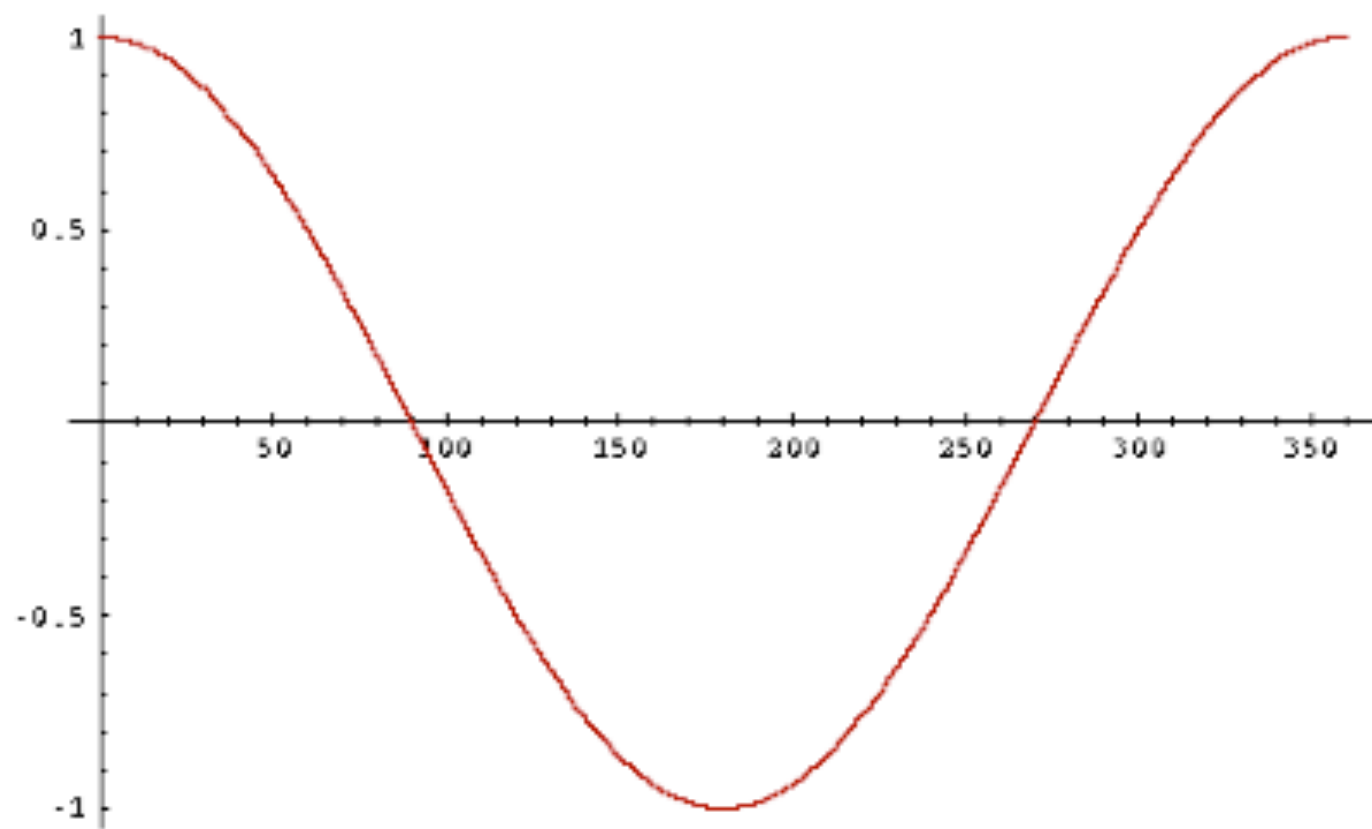
Use angle instead of distance

- Rank documents according to angle with query
- Thought experiment: take a document d and append it to itself. Call this document d' . d' is twice as long as d .
- “Semantically” d and d' have the same content.
- The angle between the two documents is 0, corresponding to maximal similarity ...
- ... even though the Euclidean distance between the two documents can be quite large.

From angles to cosines

- The following two notions are equivalent.
 - Rank documents according to the **angle** between query and document in decreasing order
 - Rank documents according to **cosine**(query,document) in increasing order
- Cosine is a monotonically decreasing function of the angle for the interval $[0^\circ, 180^\circ]$

Cosine



Length normalization

- How do we compute the cosine?
- A vector can be (length-) normalized by dividing each of its components by its length – here we use the L_2 norm:

$$||x||_2 = \sqrt{\sum_i x_i^2}$$

- This maps vectors onto the unit sphere ...
- ... since after normalization: $||x||_2 = \sqrt{\sum_i x_i^2} = 1.0$
- As a result, longer documents and shorter documents have weights of the same order of magnitude.
- Effect on the two documents d and d' (d appended to itself) from earlier slide: they have **identical vectors** after length-normalization.

Cosine similarity between query and document

$$\cos(\vec{q}, \vec{d}) = \text{SIM}(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}| |\vec{d}|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

- q_i is the tf-idf weight of term i in the query.
- d_i is the tf-idf weight of term i in the document.
- $|\vec{q}|$ and $|\vec{d}|$ are the lengths of \vec{q} and \vec{d} .
- This is the **cosine similarity** of \vec{q} and \vec{d} or, equivalently, the cosine of the angle between \vec{q} and \vec{d} .

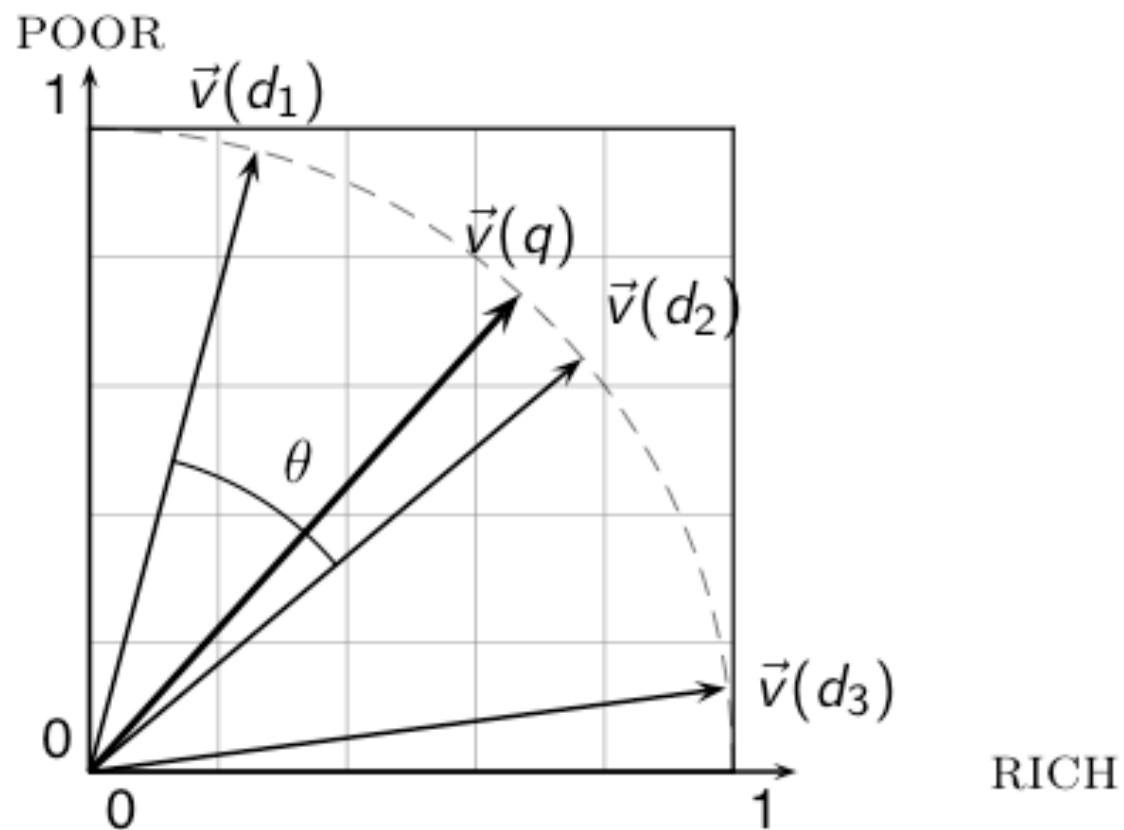
Cosine for normalized vectors

- For normalized vectors, the cosine is equivalent to the dot product or scalar product.

$$\cos(\vec{q}, \vec{d}) = \vec{q} \cdot \vec{d} = \sum_i q_i \cdot d_i$$

- (if \vec{q} and \vec{d} are length-normalized).

Cosine similarity illustrated



Cosine: Example

| | (counts) term | term frequencies | | |
|--|------------------|------------------|-----|----|
| | | SaS | PaP | WH |
| How similar are these novels? SaS: Sense and Sensibility PaP: Pride and Prejudice WH: Wuthering Heights | AFFECTION | 115 | 58 | 20 |
| | JEALOUS | 10 | 7 | 11 |
| | GOSSIP | 2 | 0 | 6 |
| | WUTHERING | 0 | 0 | 38 |

Cosine: Example

term frequencies (counts)

log frequency weighting

| term | SaS | PaP | WH | term | SaS | PaP | WH |
|-----------|-----|-----|----|-----------|------|------|------|
| AFFECTION | 115 | 58 | 20 | AFFECTION | 3.06 | 2.76 | 2.30 |
| JEALOUS | 10 | 7 | 11 | JEALOUS | 2.0 | 1.85 | 2.04 |
| GOSSIP | 2 | 0 | 6 | GOSSIP | 1.30 | 0 | 1.78 |
| WUTHERING | 0 | 0 | 38 | WUTHERING | 0 | 0 | 2.58 |

(To simplify this example, we don't do idf weighting.)

Cosine: Example

log frequency weighting

| term | SaS | PaP | WH |
|-----------|------|------|------|
| AFFECTION | 3.06 | 2.76 | 2.30 |
| JEALOUS | 2.0 | 1.85 | 2.04 |
| GOSSIP | 1.30 | 0 | 1.78 |
| WUTHERING | 0 | 0 | 2.58 |

log frequency weighting &
cosine normalization

| term | SaS | PaP | WH |
|-----------|-------|-------|-------|
| AFFECTION | 0.789 | 0.832 | 0.524 |
| JEALOUS | 0.515 | 0.555 | 0.465 |
| GOSSIP | 0.335 | 0.0 | 0.405 |
| WUTHERING | 0.0 | 0.0 | 0.588 |

- $\cos(\text{SaS}, \text{PaP}) \approx 0.789 * 0.832 + 0.515 * 0.555 + 0.335 * 0.0 + 0.0 * 0.0 \approx 0.94$.
- $\cos(\text{SaS}, \text{WH}) \approx 0.79$
- $\cos(\text{PaP}, \text{WH}) \approx 0.69$
- Why do we have $\cos(\text{SaS}, \text{PaP}) > \cos(\text{SAS}, \text{WH})$?

Computing the cosine score

COSINESCORE(q)

```
1  float Scores[ $N$ ] = 0
2  float Length[ $N$ ]
3  for each query term  $t$ 
4  do calculate  $w_{t,q}$  and fetch postings list for  $t$ 
5      for each pair( $d, tf_{t,d}$ ) in postings list
6      do Scores[ $d$ ] + =  $w_{t,d} \times w_{t,q}$ 
7  Read the array Length
8  for each  $d$ 
9  do Scores[ $d$ ] = Scores[ $d$ ] / Length[ $d$ ]
10 return Top  $K$  components of Scores[]
```

Components of tf-idf weighting

| Term frequency | | Document frequency | | Normalization | |
|----------------|---|--------------------|---------------------------------------|--------------------|--|
| n (natural) | $tf_{t,d}$ | n (no) | 1 | n (none) | 1 |
| l (logarithm) | $1 + \log(tf_{t,d})$ | t (idf) | $\log \frac{N}{df_t}$ | c (cosine) | $\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_M^2}}$ |
| a (augmented) | $0.5 + \frac{0.5 \times tf_{t,d}}{\max_t(tf_{t,d})}$ | p (prob idf) | $\max\{0, \log \frac{N-df_t}{df_t}\}$ | u (pivoted unique) | $1/u$ |
| b (boolean) | $\begin{cases} 1 & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$ | | | b (byte size) | $1/CharLength^\alpha$, $\alpha < 1$ |
| L (log ave) | $\frac{1 + \log(tf_{t,d})}{1 + \log(\text{ave}_{t \in d}(tf_{t,d}))}$ | | | | |

tf-idf example

- We often use **different weightings** for queries and documents.
- Notation: ddd.qqq
- Example: Inc.ltn
- document: logarithmic tf, no df weighting, cosine normalization
- query: logarithmic tf, idf, no normalization
- **Isn't it bad to not idf-weight the document?**
- Example query: “best car insurance”
- Example document: “car insurance auto insurance”

tf-idf example: Inc.ltn

Query: “best car insurance”. Document: “car insurance auto insurance”.

| word | query | | | | | document | | | | product |
|-----------|--------|---------|-------|-----|--------|----------|---------|--------|---------|---------|
| | tf-raw | tf-wght | df | idf | weight | tf-raw | tf-wght | weight | n'lized | |
| auto | 0 | 0 | 5000 | 2.3 | 0 | 1 | 1 | 1 | 0.52 | 0 |
| best | 1 | 1 | 50000 | 1.3 | 1.3 | 0 | 0 | 0 | 0 | 0 |
| car | 1 | 1 | 10000 | 2.0 | 2.0 | 1 | 1 | 1 | 0.52 | 1.04 |
| insurance | 1 | 1 | 1000 | 3.0 | 3.0 | 2 | 1.3 | 1.3 | 0.68 | 2.04 |

Key to columns: tf-raw: raw (unweighted) term frequency, tf-wght: logarithmically weighted term frequency, df: document frequency, idf: inverse document frequency, weight: the final weight of the term in the query or document, n'lized: document weights after cosine normalization, product: the product of final query weight and final document weight

$$\sqrt{1^2 + 0^2 + 1^2 + 1.3^2} \approx 1.92$$

$$1/1.92 \approx 0.52$$

$$1.3/1.92 \approx 0.68$$

Final similarity score between query and document: $\sum_i w_{qi} \cdot w_{di} = 0 + 0 + 1.04 + 2.04 = 3.08$ Questions?

Summary: Ranked retrieval in the vector space model

- Represent the query as a weighted tf-idf vector
- Represent each document as a weighted tf-idf vector
- Compute the cosine similarity between the query vector and each document vector
- Rank documents with respect to the query
- Return the top K (e.g., $K = 10$) to the user

Table 2. Typical term-weighting formulas

| Weighting System | Document term weight | Query Term weight |
|--|---|--|
| Best fully weighted system $tfc \cdot nfx$ | $\frac{tf \cdot \log \frac{N}{n}}{\sqrt{\sum_{vector} \left(tf_i \cdot \log \frac{N}{n_i} \right)^2}}$ | $\left(0.5 + \frac{0.5 \text{ tf}}{\max \text{ tf}} \right) \cdot \log \frac{N}{n}$ |
| Best weighted probabilistic weight $nxx \cdot bpx$ | $0.5 + \frac{0.5 \text{ tf}}{\max \text{ tf}}$ | $\log \frac{N - n}{n}$ |
| Classical idf weight $bfx \cdot bfx$ | $\log \frac{N}{n}$ | $\log \frac{N}{n}$ |
| Binary term independence $bxx \cdot bpx$ | 1 | $\log \frac{N - n}{n}$ |
| Standard tf weight: $txc \cdot txx$ | $\frac{tf}{\sqrt{\sum_{vector} (tf_i)^2}}$ | tf |
| Coordination level $bxx \cdot bxx$ | 1 | 1 |