

# Introduction to Information Retrieval

<http://informationretrieval.org>

## IIR 6&7: Vector Space Model

Hinrich Schütze

Institute for Natural Language Processing, University of Stuttgart

2011-08-29

# Models and Methods

- 1 Boolean model and its limitations (30)
- 2 **Vector space model (30)**
- 3 Probabilistic models (30)
- 4 Language model-based retrieval (30)
- 5 Latent semantic indexing (30)
- 6 Learning to rank (30)

# Take-away

# Take-away

- **tf-idf weighting**: Quick review of tf-idf weighting

# Take-away

- **tf-idf weighting**: Quick review of tf-idf weighting
- **Vector space model** – represents queries and documents in a high-dimensional space.

# Take-away

- **tf-idf weighting**: Quick review of tf-idf weighting
- **Vector space model** – represents queries and documents in a high-dimensional space.
- **Pivot normalization** (or “pivoted document length normalization”): alternative to cosine normalization that removes a bias inherent in standard length normalization

# Outline

- 1 tf-idf weighting
- 2 Vector space model
- 3 Pivot length normalization

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





# 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**  $\in \{0, 1\}^{|V|}$ .



# 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  $\in \mathbb{N}^{|V|}$ .



# 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**  $\in \mathbb{N}^{|V|}$ .



# 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  $d$** .

# 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 rank documents according to query-document matching scores and use tf as a component in these matching scores.

# 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 rank documents according to query-document matching scores and use tf as a component in these matching scores.
- But how?

# 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 rank documents according to query-document matching scores and use tf as a component in these matching scores.
- But how?
- Raw term frequency is not what we want because:

# 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 rank documents according to query-document matching scores and use tf as a component in these matching 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.



# 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 rank documents according to query-document matching scores and use tf as a component in these matching 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.

# 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 rank documents according to query-document matching scores and use tf as a component in these matching 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

# Instead of raw frequency: Log frequency weighting

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

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

# Instead of raw frequency: Log frequency weighting

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

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

- $\text{tf}_{t,d} \rightarrow w_{t,d}$ :  
 $0 \rightarrow 0, 1 \rightarrow 1, 2 \rightarrow 1.3, 10 \rightarrow 2, 1000 \rightarrow 4, \text{ etc.}$

# 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.
- Matching score for a document-query pair: sum over terms  $t$  in both  $q$  and  $d$ :  
 $\text{tf-matching-score}(q, d) = \sum_{t \in q \cap d} (1 + \log \text{tf}_{t,d})$  □

# Frequency in document vs. frequency in collection

- In addition, to term frequency (the frequency of the term 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. □



# idf weight

# idf weight

- $df_t$  is the document frequency, the number of documents that  $t$  occurs in.

# 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  $t$ .

# 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  $t$ .
- Inverse document frequency,  $idf_t$ , is a direct measure of the **informativeness** of the term.

# 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  $t$ .
- Inverse document frequency,  $idf_t$ , is a direct measure of the **informativeness** of the term.
- The **idf weight** of term  $t$  is defined as follows:

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

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

# 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  $t$ .
- Inverse document frequency,  $idf_t$ , is a direct measure of the **informativeness** of the term.
- The **idf weight** of term  $t$  is defined as follows:

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

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

- $[\log N/df_t]$  instead of  $[N/df_t]$  to “dampen” the effect of idf



# Examples for idf

# Examples for idf

$$\text{idf}_t = \log_{10} \frac{1,000,000}{\text{df}_t}$$

term	$\text{df}_t$	$\text{idf}_t$
calpurnia	1	6
animal	100	4
sunday	1000	3
fly	10,000	2
under	100,000	1
the	1,000,000	0



# Effect of idf on ranking

# Effect of idf on ranking

- idf gives high weights to rare terms like ARACHNOCENTRIC.

# Effect of idf on ranking

- idf gives **high weights to rare terms** like ARACHNOCENTRIC.
- idf gives **low weights to frequent words** like GOOD, INCREASE, and LINE.

# Effect of idf on ranking

- idf gives **high weights to rare terms** like ARACHNOCENTRIC.
- idf gives **low weights to frequent words** like GOOD, INCREASE, and LINE.
- idf affects the ranking of documents for **queries with at least two terms**.

# Effect of idf on ranking

- idf gives **high weights to rare terms** like ARACHNOCENTRIC.
- idf gives **low weights to frequent words** like GOOD, INCREASE, and LINE.
- idf affects the ranking of documents for **queries with at least two terms**.
- For example, in the query “arachnocentric line”, idf weighting **increases** the relative weight of ARACHNOCENTRIC and **decreases** the relative weight of LINE.

# Effect of idf on ranking

- idf gives **high weights to rare terms** like ARACHNOCENTRIC.
- idf gives **low weights to frequent words** like GOOD, INCREASE, and LINE.
- idf affects the ranking of documents for **queries with at least two terms**.
- For example, in the query “arachnocentric line”, idf weighting **increases** the relative weight of ARACHNOCENTRIC and **decreases** the relative weight of LINE.
- idf has **little effect** on ranking for **one-term queries**. □

# Summary: tf-idf weighting

# Summary: tf-idf weighting

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

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



# Summary: tf-idf weighting

- 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 ...

# Summary: tf-idf weighting

- 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 component)

# Summary: tf-idf weighting

- 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 component)
  - ... increases with the rarity of the term in the collection. (inverse document frequency component)



# Outline

- 1 tf-idf weighting
- 2 **Vector space model**
- 3 Pivot length normalization

# 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**  $\in \{0, 1\}^{|V|}$ .



# 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**  $\in \mathbb{N}^{|V|}$ .



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  $\in \mathbb{R}^{|V|}$ .



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  $\in \mathbb{R}^{|V|}$ .

# 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.

# 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 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.

# 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

# 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

# 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



# 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

# 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

# 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.

# 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?

# How do we formalize vector space similarity?

- First cut: (negative) distance between two points

# How do we formalize vector space similarity?

- First cut: (negative) distance between two points
- ( = distance between the end points of the two vectors)

# 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?



# How do we formalize vector space similarity?

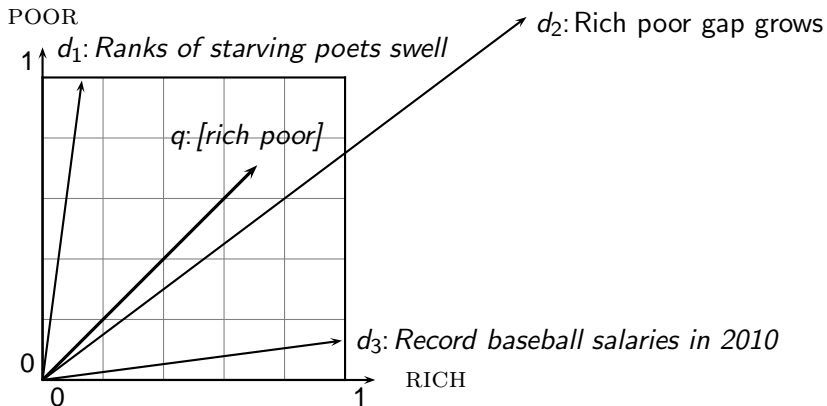
- First cut: (negative) distance between two points
- ( = distance between the end points of the two vectors)
- Euclidean distance?
- Euclidean distance is a bad idea . . .

# 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

# 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. □

# Use angle instead of distance

# Use angle instead of distance

- Rank documents according to angle with query

# Use angle instead of distance

- Rank documents according to angle with query
- The following two notions are equivalent.

# Use angle instead of distance

- Rank documents according to angle with query
- The following two notions are equivalent.
  - Rank documents according to the **angle** between query and document in decreasing order



# Use angle instead of distance

- Rank documents according to angle with query
- 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

# Use angle instead of distance

- Rank documents according to angle with query
- 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  $[0^\circ, 180^\circ]$

# Use angle instead of distance

- Rank documents according to angle with query
- 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  $[0^\circ, 180^\circ]$
- → do ranking according to cosine □

# Cosine similarity between query and document

# Cosine similarity between query and document



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

# Cosine similarity between query and document



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

- $q_i$  is the tf-idf weight of term  $i$  in the query.

# Cosine similarity between query and document



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

- $q_i$  is the tf-idf weight of term  $i$  in the query.
- $d_i$  is the tf-idf weight of term  $i$  in the document.

# Cosine similarity between query and document



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

- $q_i$  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  $\vec{d}$ .



# Cosine similarity between query and document



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

- $q_i$  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 similarity between query and document



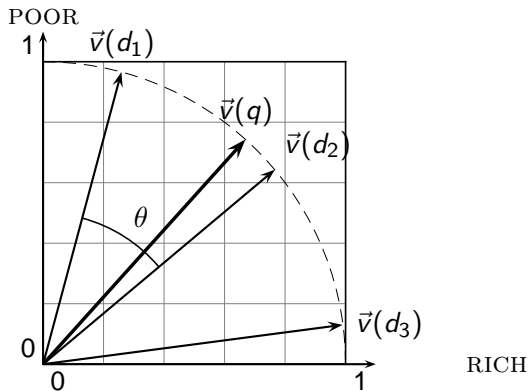
$$\cos(\vec{q}, \vec{d}) = \text{SIM}(\vec{q}, \vec{d}) = \frac{\vec{q}}{|\vec{q}|} \cdot \frac{\vec{d}}{|\vec{d}|} = \sum_{i=1}^{|\mathcal{V}|} \frac{q_i}{|\vec{q}|} \cdot \frac{d_i}{|\vec{d}|}$$

- $q_i$  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 similarity = dot product of length-normalized vectors



# Cosine similarity illustrated

# Cosine similarity illustrated



# 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)	$\frac{1 + \log(tf_{t,d})}{1 + \log(\text{ave}_{t \in d}(tf_{t,d}))}$				

# 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)	$\frac{1 + \log(tf_{t,d})}{1 + \log(\text{ave}_{t \in d}(tf_{t,d}))}$				

Best known combination of weighting options

# tf-idf example

# tf-idf example

- We often use **different weightings** for queries and documents.



# tf-idf example

- We often use **different weightings** for queries and documents.
- Notation: ddd.qqq

# tf-idf example

- We often use **different weightings** for queries and documents.
- Notation: ddd.qqq
- Example: Inc.ltn

# 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

# 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

# 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?**

# 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”

# 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										
best										
car										
insurance										

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



# 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									
best	1									
car	1									
insurance	1									

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

# 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					1				
best	1					0				
car	1					1				
insurance	1					2				

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

# 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				1				
best	1	1				0				
car	1	1				1				
insurance	1	1				2				

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

# 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				1	1			
best	1	1				0	0			
car	1	1				1	1			
insurance	1	1				2	1.3			

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

# 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			1	1			
best	1	1	50000			0	0			
car	1	1	10000			1	1			
insurance	1	1	1000			2	1.3			

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

# 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	
auto	0	0	5000	2.3		1	1		
best	1	1	50000	1.3		0	0		
car	1	1	10000	2.0		1	1		
insurance	1	1	1000	3.0		2	1.3		

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

# 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			
best	1	1	50000	1.3	1.3	0	0			
car	1	1	10000	2.0	2.0	1	1			
insurance	1	1	1000	3.0	3.0	2	1.3			

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

# 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			
best	1	1	50000	1.3	1.3	0	0			
car	1	1	10000	2.0	2.0	1	1			
insurance	1	1	1000	3.0	3.0	2	1.3			

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



# 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		
best	1	1	50000	1.3	1.3	0	0	0		
car	1	1	10000	2.0	2.0	1	1	1		
insurance	1	1	1000	3.0	3.0	2	1.3	1.3		

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

# 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	
best	1	1	50000	1.3	1.3	0	0	0	0	
car	1	1	10000	2.0	2.0	1	1	1	0.52	
insurance	1	1	1000	3.0	3.0	2	1.3	1.3	0.68	

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

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

$$1/1.92 \approx 0.52$$

$$1.3/1.92 \approx 0.68$$

# 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**

# 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

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

# Outline

- 1 tf-idf weighting
- 2 Vector space model
- 3 Pivot length normalization

# A problem for cosine normalization

# A problem for cosine normalization

- Query  $q$ : “anti-doping rules Beijing 2008 olympics”

# A problem for cosine normalization

- Query  $q$ : “anti-doping rules Beijing 2008 olympics”
- Compare three documents



# A problem for cosine normalization

- Query  $q$ : “anti-doping rules Beijing 2008 olympics”
- Compare three documents
  - $d_1$ : a short document on anti-doping rules at 2008 Olympics

# A problem for cosine normalization

- Query  $q$ : “anti-doping rules Beijing 2008 olympics”
- Compare three documents
  - $d_1$ : a short document on anti-doping rules at 2008 Olympics
  - $d_2$ : a long document that consists of a copy of  $d_1$  and 5 other news stories, all on topics different from Olympics/anti-doping

# A problem for cosine normalization

- Query  $q$ : “anti-doping rules Beijing 2008 olympics”
- Compare three documents
  - $d_1$ : a short document on anti-doping rules at 2008 Olympics
  - $d_2$ : a long document that consists of a copy of  $d_1$  and 5 other news stories, all on topics different from Olympics/anti-doping
  - $d_3$ : a short document on anti-doping rules at the 2004 Athens Olympics

# A problem for cosine normalization

- Query  $q$ : “anti-doping rules Beijing 2008 olympics”
- Compare three documents
  - $d_1$ : a short document on anti-doping rules at 2008 Olympics
  - $d_2$ : a long document that consists of a copy of  $d_1$  and 5 other news stories, all on topics different from Olympics/anti-doping
  - $d_3$ : a short document on anti-doping rules at the 2004 Athens Olympics
- What ranking do we expect in the vector space model?

# Pivot normalization

- Cosine normalization produces weights that are **too large for short documents** and **too small for long documents** (on average).

# Pivot normalization

- Cosine normalization produces weights that are **too large for short documents** and **too small for long documents** (on average).
- Adjust cosine normalization by linear adjustment: “turning” the average normalization on the **pivot**

# Pivot normalization

- Cosine normalization produces weights that are **too large for short documents** and **too small for long documents** (on average).
- Adjust cosine normalization by linear adjustment: “turning” the average normalization on the **pivot**
- Effect: Similarities of short documents with query **decrease**; similarities of long documents with query **increase**.

# Pivot normalization

- Cosine normalization produces weights that are **too large for short documents** and **too small for long documents** (on average).
- Adjust cosine normalization by linear adjustment: “turning” the average normalization on the **pivot**
- Effect: Similarities of short documents with query **decrease**; similarities of long documents with query **increase**.
- This removes the unfair advantage that short documents have.

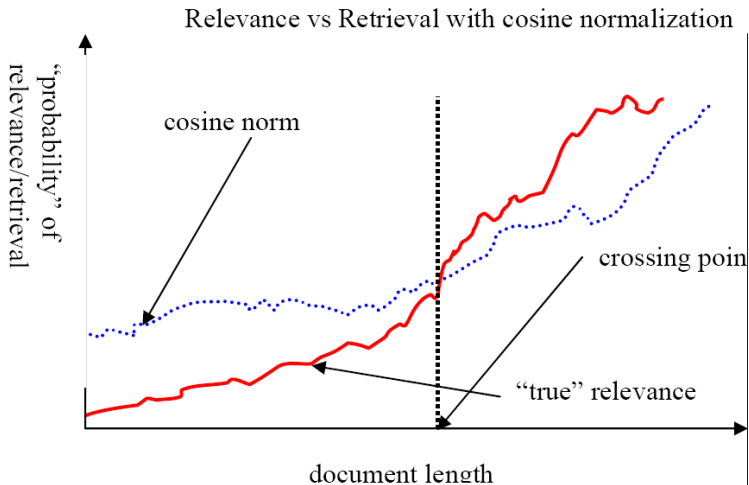


# Pivot normalization

- Cosine normalization produces weights that are **too large for short documents** and **too small for long documents** (on average).
- Adjust cosine normalization by linear adjustment: “turning” the average normalization on the **pivot**
- Effect: Similarities of short documents with query **decrease**; similarities of long documents with query **increase**.
- This removes the unfair advantage that short documents have.
- Singhal’s study is also interesting from the point of view of methodology. □

# Predicted and true probability of relevance

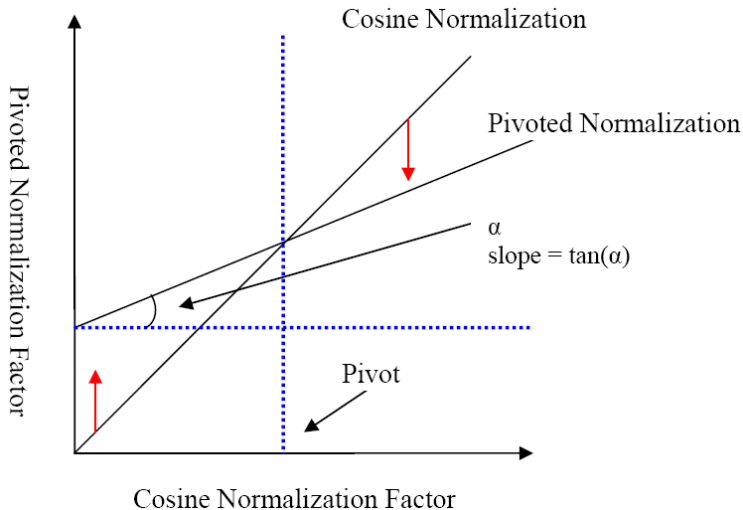
# Predicted and true probability of relevance



source:  
Lillian Lee

# Pivot normalization

# Pivot normalization



source:  
Lillian Lee

# Pivoted normalization: Amit Singhal's experiments

# Pivoted normalization: Amit Singhal's experiments

Cosine	Pivoted Cosine Normalization				
	Slope				
	0.60	0.65	0.70	<b>0.75</b>	0.80
6,526	6,342	6,458	6,574	<b>6,629</b>	6,671
0.2840	0.3024	0.3097	0.3144	<b>0.3171</b>	0.3162
Improvement	+ 6.5%	+ 9.0%	+10.7%	<b>+11.7%</b>	+11.3%

(relevant documents retrieved and (change in) average precision)



# Summary: Ranked retrieval in the vector space model



# Summary: Ranked retrieval in the vector space model

- Represent each document as a weighted tf-idf vector

# Summary: Ranked retrieval in the vector space model

- Represent each document as a weighted tf-idf vector
- Represent the query as a weighted tf-idf vector

## Summary: Ranked retrieval in the vector space model

- Represent each document as a weighted tf-idf vector
- Represent the query as a weighted tf-idf vector
- Compute the cosine similarity between the query vector and each document vector

# Summary: Ranked retrieval in the vector space model

- Represent each document as a weighted tf-idf vector
- Represent the query as a weighted tf-idf vector
- Compute the cosine similarity between the query vector and each document vector
  - Alternatively, use pivot normalization

# Summary: Ranked retrieval in the vector space model

- Represent each document as a weighted tf-idf vector
- Represent the query as a weighted tf-idf vector
- Compute the cosine similarity between the query vector and each document vector
  - Alternatively, use pivot normalization
- Rank documents with respect to the query

# Summary: Ranked retrieval in the vector space model

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



# Take-away

- **tf-idf weighting**: Quick review of tf-idf weighting
- **Vector space model** – represents queries and documents in a high-dimensional space.
- **Pivot normalization** (or “pivoted document length normalization”): alternative to cosine normalization that removes a bias inherent in standard length normalization

# Resources

- Chapters 6 and 7 of Introduction to Information Retrieval
- Resources at <http://informationretrieval.org/essir2011>
  - Gerard Salton (main proponent of vector space model in 70s, 80s, 90s)
  - Exploring the similarity space (Moffat and Zobel, 2005)
  - Pivot normalization (original paper)