

Information Retrieval Overview

PRI 23/24 · Information Processing and Retrieval
M.EIC · Master in Informatics Engineering and Computation

Sérgio Nunes
Dept. Informatics Engineering
FEUP · U.Porto

Based on Chapters 1, 2, 6, 11, 12 from Introduction to Information Retrieval, Manning, C. et al. (2008)
Based on Chapter 3 from Modern Information Retrieval, Baeza-Yates et al. (2011)

PRI Status

→ Lecture: IR techniques and models.

→ Lab:

→ 1. Indexing and retrieval of project collection in Solr (last week);

→ 2. Evaluation of results (**this week**);

→ 3. M2 workshop (next week).

Background Concepts

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	
...							

Figure 1.1 A term-document incidence matrix. Matrix element (t, d) is 1 if the play in column d contains the word in row t , and is 0 otherwise.

Boolean Model

- In the Boolean Retrieval Model queries are represented in the form of a Boolean expression of terms.
- E.g.: [Brutus AND Caesar AND NOT Calpurnia]
 - $\Rightarrow 110100 \text{ AND } 110111 \text{ AND } 101111 = 100100$
- The model views each document as a set of words.
- Bag of Words (BoW) model, where the exact ordering of terms in a document is ignored and only their presence is considered.

Inverted Index

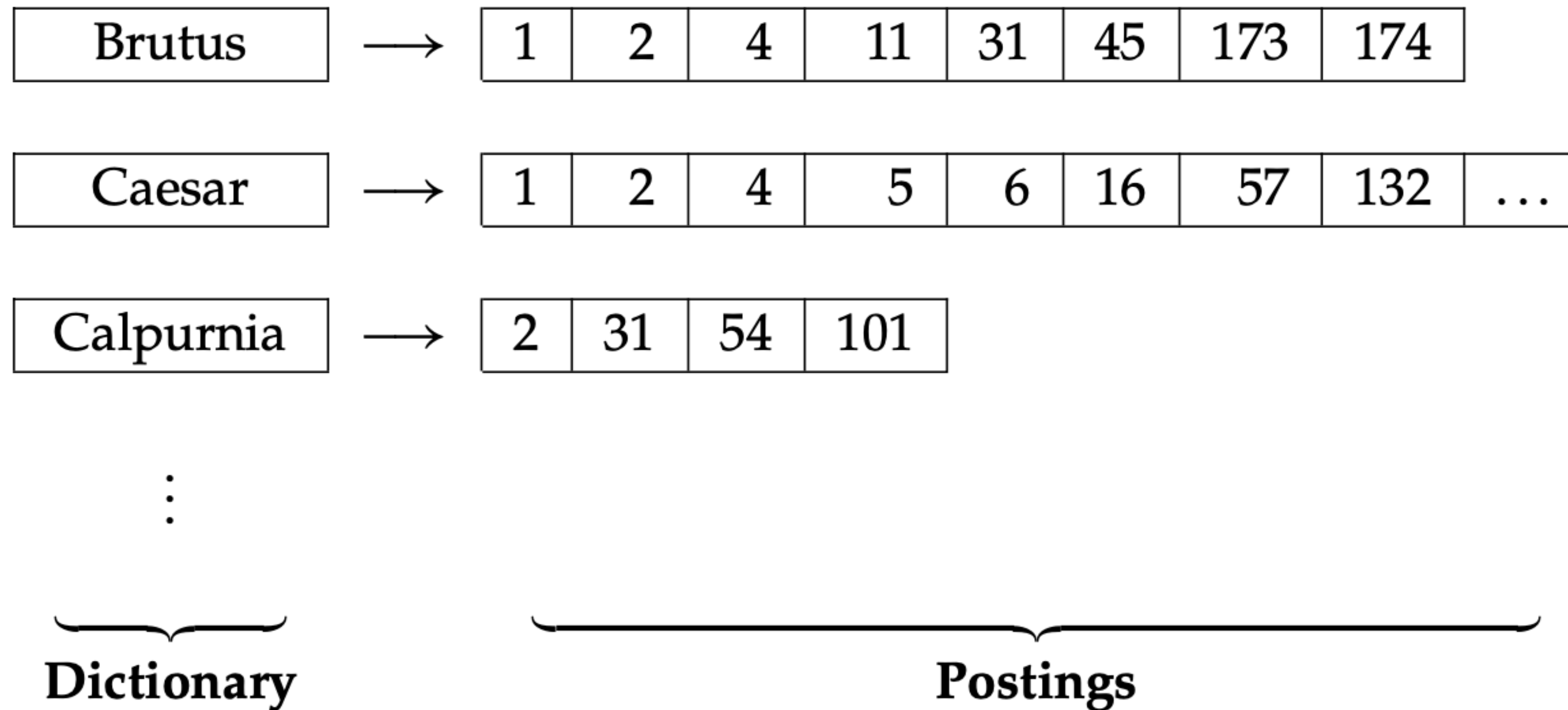


Figure 1.3 The two parts of an inverted index. The dictionary is commonly kept in memory, with pointers to each postings list, which is stored on disk.

Inverted Index (Basic)

Vocabulary	n_i	Occurrences as inverted lists
to	2	[1,4],[2,2]
do	3	[1,2],[3,3],[4,3]
is	1	[1,2]
be	4	[1,2],[2,2],[3,2],[4,2]
or	1	[2,1]
not	1	[2,1]
I	2	[2,2],[3,2]
am	2	[2,2],[3,1]
what	1	[2,1]
think	1	[3,1]
therefore	1	[3,1]
da	1	[4,3]
let	1	[4,2]
it	1	[4,2]

To do is to be.
To be is to do.

d_1

To be or not to be.
I am what I am.

d_2

I think therefore I am.
Do be do be do.

d_3

Do do do, da da da.
Let it be, let it be.

d_4

Inverted Index with Positional Information

Vocabulary	n_i	Occurrences as full inverted lists
to	2	[1,4,[1,4,6,9]], [2,2,[1,5]]
do	3	[1,2,[2,10]], [3,3,[6,8,10]], [4,3,[1,2,3]]
is	1	[1,2,[3,8]]
be	4	[1,2,[5,7]], [2,2,[2,6]], [3,2,[7,9]], [4,2,[9,12]]
or	1	[2,1,[3]]
not	1	[2,1,[4]]
I	2	[2,2,[7,10]], [3,2,[1,4]]
am	2	[2,2,[8,11]], [3,1,[5]]
what	1	[2,1,[9]]
think	1	[3,1,[2]]
therefore	1	[3,1,[3]]
da	1	[4,3,[4,5,6]]
let	1	[4,2,[7,10]]
it	1	[4,2,[8,11]]

To do is to be.
To be is to do.

d_1

To be or not to be.
I am what I am.

d_2

I think therefore I am.
Do be do be do.

d_3

Do do do, da da da.
Let it be, let it be.

d_4

Index Construction

- Choosing the document unit for indexing and the index granularity are important first steps.
- There is a precision / recall tradeoff in this decision.
 - "If the units get too small (e.g. sentences), we are likely to miss important passages because terms were distributed over several mini-documents, whereas if units are too large (e.g. books) we tend to get spurious matches and the relevant information is hard for the user to find."

Tokenization

→ "Given a character sequence and a defined document unit, tokenization is the task of chopping it into pieces, called tokens."

Input: Friends, Romans, Countrymen, lend me your ears;

Output:

Friends	Romans	Countrymen	lend	me	your	ears
---------	--------	------------	------	----	------	------

→ "A token is an instance of a sequence of characters in some particular document that are grouped together as a useful semantic unit for processing."

→ "A type is the class of all tokens containing the same character sequence."

Stop Words

- Stop words are words considered of little value in helping select documents in the retrieval process, e.g. a word that exists in all documents.
- The strategy to determine stop words is by looking the collection frequency of a term (*cf*), i.e. the total number of times a term appears in the document collection.
- A stop list is a (commonly) hand-picked list of terms to be discarded during the indexing process.

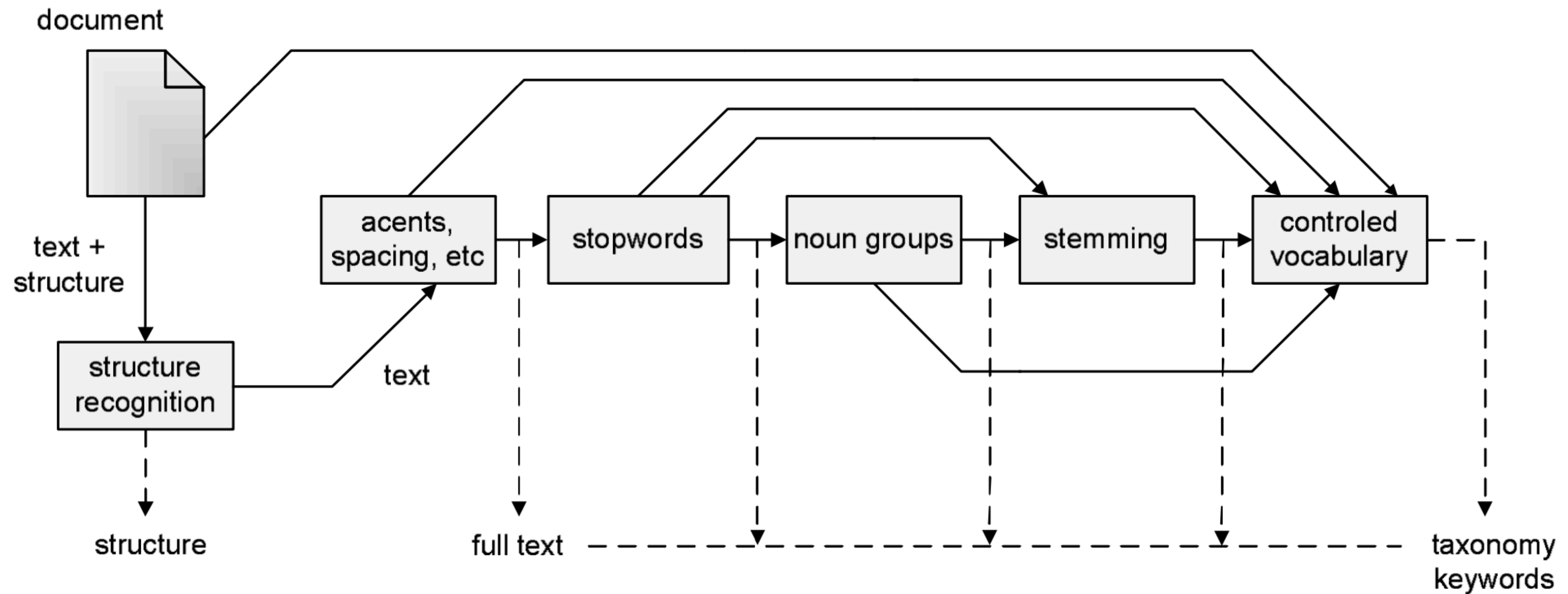
a	an	and	are	as	at	be	by	for	from
has	he	in	is	it	its	of	on	that	the
to	was	were	will	with					

Figure 2.5 A stop list of twenty-five semantically nonselective words that are common in Reuters-RCV1.

Token Normalization

- "Token normalization is the process of canonicalizing tokens so that matches occur despite superficial differences in the character sequences of the tokens."
- Some forms of normalization commonly employed are: accents, capitalization, and stemming and lemmatization for dealing with different forms of a word (e.g. watch, watches, watching).
- Stemming refers to a heuristic process that chops off the ends of words in the hope of reducing inflectional forms.
- Lemmatization refers to the process of reducing inflectional forms by using vocabularies and the morphological analysis of words to find its lemma.

From full text to index terms



Term Weighting

Ranked Retrieval

- With the Boolean model, a document either matches or does not match a query.
- In large document collections this is not feasible, thus it is essential for a search system to rank-order the documents.
- Note that there are scenarios where recall is determinant — i.e., all documents need to be analyzed (e.g., patent search), and thus Boolean search is used.

Parametric and Zone Indexes

- Although we have considered documents to be a simple sequence of terms, most documents have additional structure (e.g. email message). Additionally, metadata is often associated with a document (e.g. date, authors, title).
- Parametric indexes are inverted indexes built for specific parameters, or fields, that support parametric search (e.g. "all documents from author Z containing word Y").
- Zones are a similar concept applied to arbitrary free text (e.g. portion of a document). For example, a document's abstract can be associated to a specific zone index.

Zone Indexes

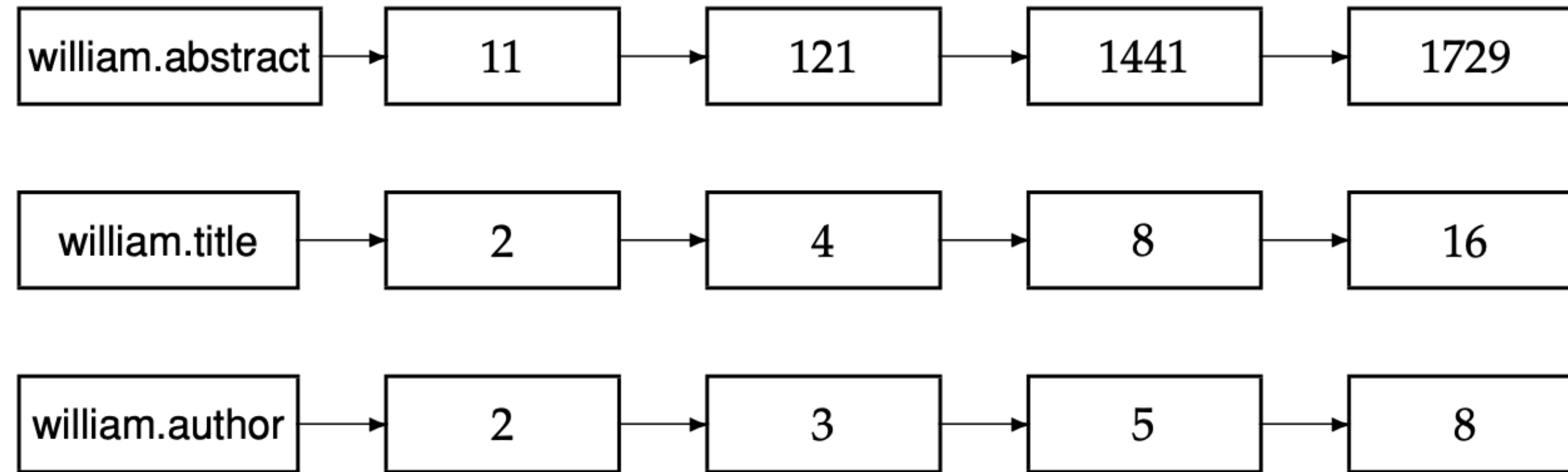


Figure 6.2 Basic zone index; zones are encoded as extensions of dictionary entries.

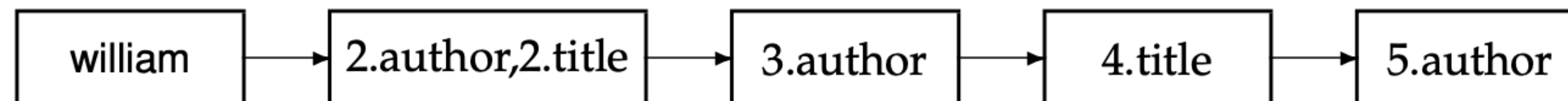


Figure 6.3 Zone index in which the zone is encoded in the postings rather than the dictionary.

Ranked Boolean Retrieval

- Zones (or fields) can be weighted differently to compute each document's relevance simply using a linear combination of zone scores, where each zone of the document contributes a Boolean value.
- Zone weights can be specified by an expert (commonly the end user) but are usually learned by the system based on training examples. This method is known as "machine-learned relevance" or "learning to rank".

Term Frequency (tf)

- The next step is not to consider only the presence (or not) of a query term in a zone, but the number of mentions of the term (i.e., term frequency).
- The simplest form of weighting terms differently is simply to assign the weight of a term to the term's frequency.
- The term frequency of a term in a document is denoted $tf_{t,d}$.
- In the bag of words model, the ordering is ignored but the number of occurrences of each term is key (in contrast with Boolean retrieval).

Document Frequency (df)

- Raw term frequency suffers from a problem: all terms are considered equally important when assessing a query, when in fact some terms are of little use in determining relevance.
- For example, in a collection of thesis dissertations, the term "dissertation" is less likely to be of value since this term probably exists in every document.
- An important measure to incorporate the discriminative power of a term in a collection is the document frequency of a term (df), i.e. the number of documents that contain a term.

Inverse Document Frequency (idf)

- The document frequency of a term is incorporated in the weight of a term by using the concept of inverse document frequency (idf).

$$\text{idf}_t = \log \frac{N}{\text{df}_t}.$$

- Where N is the total number of documents in the collection.
- The rarer the term in a collection, the higher its idf.

tf-idf

- Combining term frequency (tf) with inverse document frequency (idf) results in a classic measure in Information Retrieval, the tf-idf weighting scheme.

$$\mathbf{tf-idf}_{t,d} = \mathbf{tf}_{t,d} \times \mathbf{idf}_t.$$

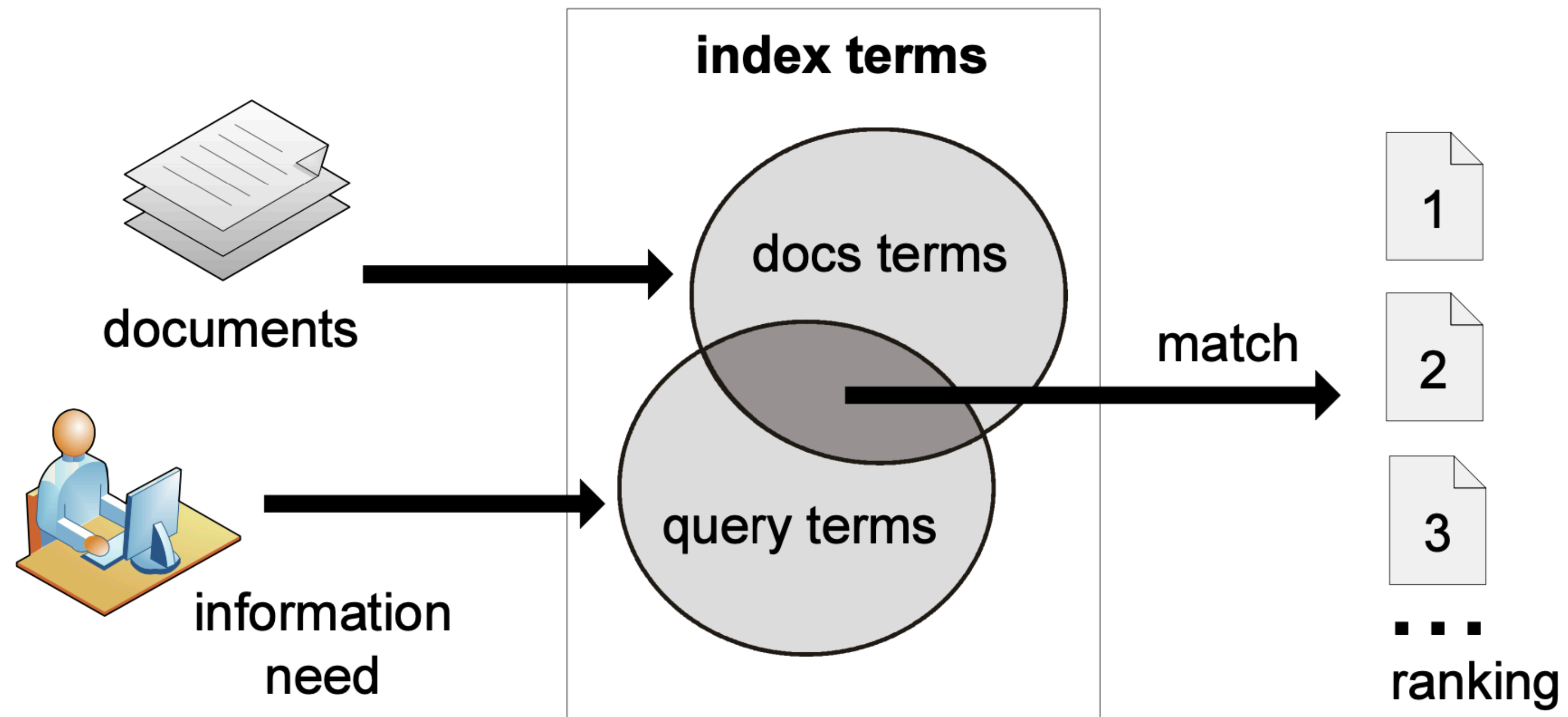
- $\mathbf{tf-idf}_{t,d}$ assigns a term t a weight in a document d that is:
 - highest when t occurs many times within a small number of documents;
 - lower when the term occurs fewer times in a document, or occurs in many documents;
 - lowest when the term occurs in virtually all documents.

Information Retrieval Models

Modeling in IR

- Modeling in IR aims at **producing a ranking function**,
 - i.e. assign scores to documents with regard to a given query.
- The process consists of **two main tasks**:
 - Conception of a logical framework for representing documents and queries.
 - Definition of a ranking function that quantifies the similarities between documents and queries.

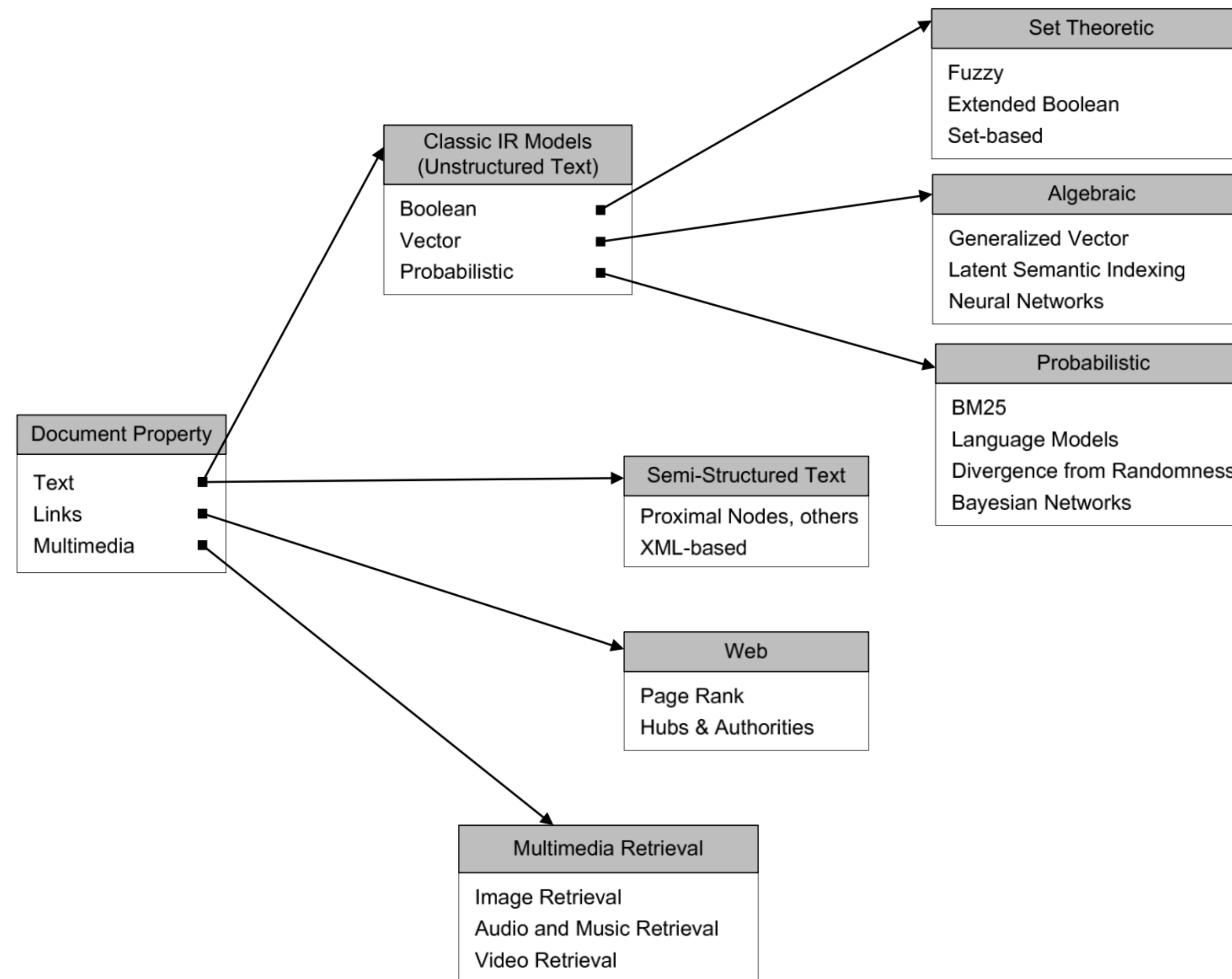
Information Retrieval Process



IR Models

- An IR model defines...
 - A representation (d_i) of the documents (D) in the collection;
 - A representation (q_j) of the queries (Q) submitted by users:
 - A ranking function $R(q_j, d_i)$.

Taxonomy of IR Models



Boolean Model

- Documents are represented in a term-document matrix.
- Queries are specified as boolean expressions.
- The Boolean model predicts if documents are either relevant or non-relevant.
- No ranking is provided.

Vector Space Model

Vector Space Model

- The representation of a set of documents as vectors in a common vector space is known as the vector space model and is fundamental to number Information Retrieval operations.
- In a nutshell, each document is represented as a vector, with a component vector for each dictionary term. tf-idf weights are used as component weights.
- Thus, the set of documents in a collection may be viewed as a set of vectors in a vector space, in which there is one axis for each term.

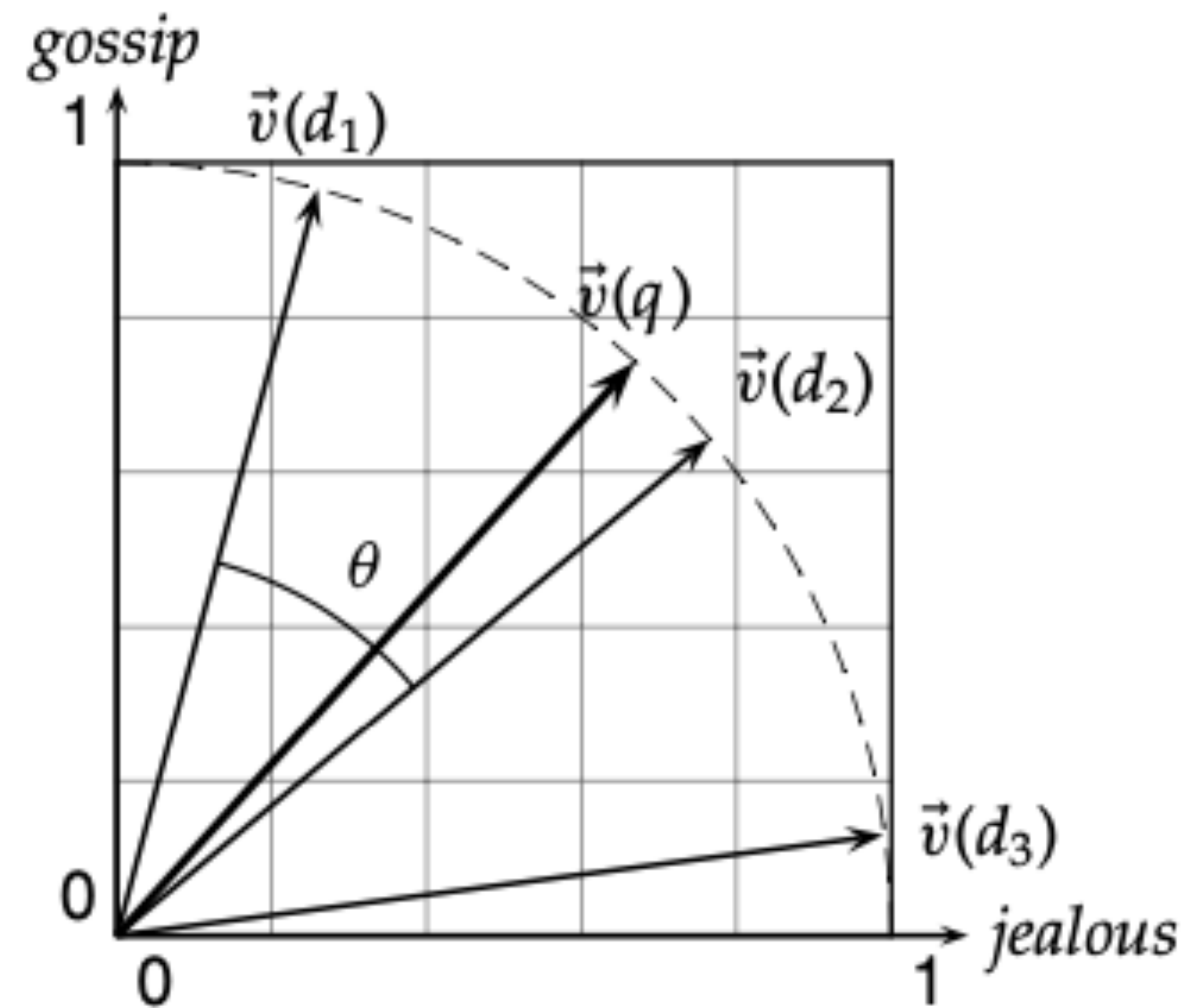


Figure 6.11 Cosine similarity illustrated: $\text{sim}(d_1, d_2) = \cos \theta$.

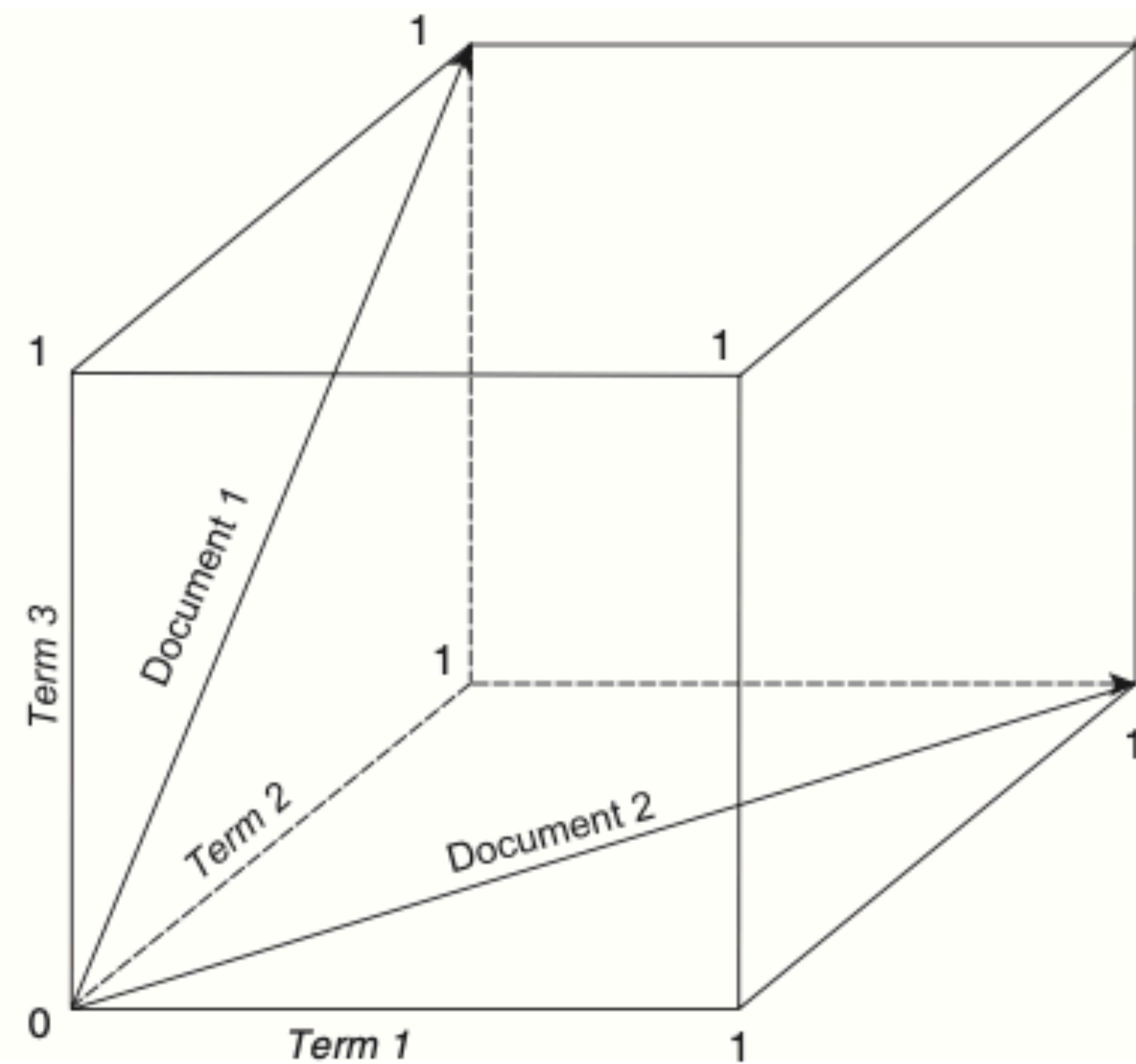


Figure 3.1

Document vectors on the vertices of the *unit hypercube*: each text is represented by a vector whose component, along each term axis, is either 0 or 1. Thus, all vectors must terminate only at one of the corners of the cube. In real life the number of dimensions is far greater than can be illustrated on a two-dimensional page.

Cosine Similarity

- To quantify the similarity between two documents in this vector space, the cosine similarity of the vector representations of the two documents.
- In other words, the similarity between two documents is given by the cosine of the angle between the two vector representations of the documents.
- This approach compensates the effect of document length.

$$\text{sim}(d_1, d_2) = \frac{\vec{V}(d_1) \cdot \vec{V}(d_2)}{|\vec{V}(d_1)| |\vec{V}(d_2)|},$$

Queries as Vectors

- Queries can also be represented as vectors in a n -dimensional space, being n the number of terms in the query. Basically, queries are viewed as very short documents.
- The top ranked results for a given query are thus the documents whose vectors have the highest cosine similarity in comparison with the query vector.

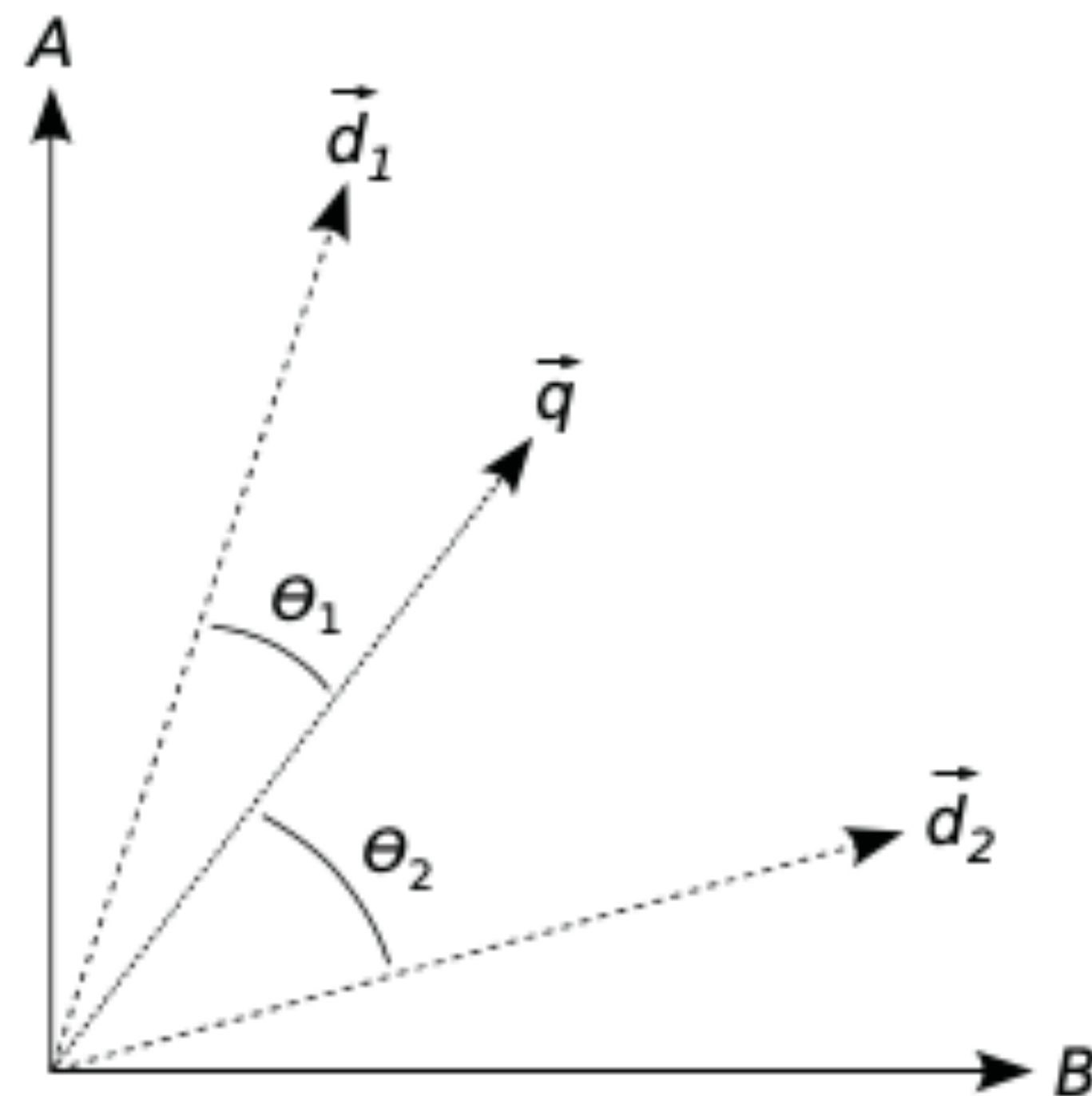


Figure 2.8 Document similarity under the vector space model. Angles are computed between a query vector \vec{q} and two document vectors \vec{d}_1 and \vec{d}_2 . Because $\theta_1 < \theta_2$, d_1 should be ranked higher than d_2 .

Inverse Document Frequency

- Use Google to estimate the inverse document frequency (idf) value of the following keywords:
- portugal, technology, health.
- Again using Google, and the same set of keywords, consider the subset of documents in the .pt domain [site:pt].

$$idf_k = \log \frac{N}{n_k}$$

- Reflect on the differences.

Inverse Document Frequency Solutions

→ What is Google's index size?

→ $\max([the], [and], [is]) \rightarrow \sim 26B$

→ World Wide Web

→ $\text{portugal} - \log(26,000,000,000 / 1,160,000,000) = \log(22.41) \approx 1.35$

→ $\text{technology} - \log(26,000,000,000 / 2,640,000,000) \approx 0.99$

→ $\text{health} - \log(26,000,000,000 / 3,140,000,000) \approx 0.92$

→ .pt domain

→ Update index size [pt site:pt] $\rightarrow \sim 278M$

→ $\text{portugal} - \log(278,000,000 / 99,600,000) = \log(2.79) \approx 0.44$

→ $\text{technology} - \log(278,000,000 / 1,940,000) \approx 2.16$

→ $\text{health} - \log(278,000,000 / 8,260,000) \approx 1.53$

Document Representation

→ Consider the following headlines

- H1: Portugal eyes political balance in presidential election
- H2: Campaigning begins for Portuguese presidential election
- H3: After Portuguese elections, Spain braces for elections

→ Obtain

- An incidence matrix considering the following terms: portugal, election, presidential.
 - A matrix with term frequency weights (raw and log normalized).
 - Inverse document frequency values for each of the three terms.
 - A matrix with tf-idf weights.
- For each headline, obtain a unitary vector representation using raw tf in 1 dimension (portugal), 2 dimensions (portugal, election) and in 3 dimensions (portugal, election, presidential).
- Calculate the cosine similarity between each document and the following query: [portugal election].

Incidence Matrix

	H1	H2	H3
portugal	1	0	0
election	1	1	1
presidential	1	1	0

- H1: Portugal eyes political balance in presidential election
- H2: Campaigning begins for Portuguese presidential election
- H3: After Portuguese elections, Spain braces for elections

Term Frequency (raw and log normalized)

	H1	H2	H3
portugal	1 (1)	0 (0)	0 (0)
election	1 (1)	1 (1)	2 (1.3)
presidential	1 (1)	1 (1)	0 (0)

- H1: Portugal eyes political balance in presidential election
- H2: Campaigning begins for Portuguese presidential election
- H3: After Portuguese elections, Spain braces for elections

Inverse Document Frequency

idf	
portugal election presidential	$\log (3 / 1) \approx 0.48$
	$\log (3 / 3) = 0$
	$\log (3 / 2) \approx 0.18$

- H1: Portugal eyes political balance in presidential election
- H2: Campaigning begins for Portuguese presidential election
- H3: After Portuguese elections, Spain braces for elections

TF-IDF (log normalized tf)

	H1	H2	H3
portugal	$1 \times 0.48 = 0.48$	$0 \times 0.48 = 0$	0
election	0	0	0
presidential	0.18	0.18	0

- H1: Portugal eyes political balance in presidential election
- H2: Campaigning begins for Portuguese presidential election
- H3: After Portuguese elections, Spain braces for elections

Unitary Vector Documents

→ 1-dimension: (portugal)

→ H1: (1)

→ H2: (0)

→ H3: (0)

→ 2-dimension (portugal, election)

→ H1: $(1, 1) > (0.707, 0.707)$

→ H2: $(0, 1) > (0, 1)$

→ H3: $(0, 2) > (0, 1)$

Document Vectors

→ 3-dimension: (portugal, election, presidential)

→ H1: $(1, 1, 1) > (0.577, 0.577, 0.577)$

→ $\|H1\| = \sqrt{1^2 + 1^2 + 1^2} = 1.732$

→ H2: $(0, 1, 1) > (0, 0.707, 0.707)$

→ H3: $(0, 2, 0) > (0, 1, 0)$

TF-IDF (log normalized tf)

→ Q: (portugal, election, presidential) = (1,1,0)

→ Unitary vector: $(1 / \sqrt{(1^2 + 1^2)}) = 0.707$

→ = (0.707, 0.707, 0)

→ $\text{sim}(H1, Q) = (0.577 \times 0.707) + (0.577 \times 0.707) + (0 \times 0.707) = 0.815$

→ $\text{sim}(H2, Q) = (0.707 \times 0) + (0.707 \times 0.707) + (0 \times 0.707) = 0.5$ (60 degrees)

→ $\text{sim}(H3, Q) = 0.707$

→ Ranking for Q: H1, H3, H2

$$\text{Cosine}(D_i, Q) = \frac{\sum_{j=1}^t d_{ij} \cdot q_j}{\sqrt{\sum_{j=1}^t d_{ij}^2 \cdot \sum_{j=1}^t q_j^2}}$$

Extra: similarity between documents

→ H1: (0.577, 0.577, 0.577)

→ H2: (0, 0.707, 0.707)

→ H3: (0, 1, 0)

	H1	H2	H3
H1	1	$(0.577 \times 0.707) \times 2 = 0.815$	0,577
H2	0,815	1	0,707
H3	0,577	0,707	$0 \times 0 + 1 \times 1 + 0 \times 0 = 1$

Language Models for Information Retrieval

Language Models for Information Retrieval

- Central idea: a document is a good match for a query if the document's model is likely to generate the query.
- In the basic language model approach, a probabilistic language model is built for each document in the collection (M_d).
- For a given query, documents are ranked based on the probability of the model generating the query: $P(q|M_d)$.

Language Models

- A language model is a function that puts a probability measure over strings drawn from some vocabulary.
- The sum of all probabilities over a vocabulary for a language model is 1.
- The simplest language model, discards all context information (i.e., nearby words), and estimated the probability of each term independently — this is called an unigram model.
 - In this case, the probability of a sequence of terms (e.g., a query) is simply the product of independent term probabilities:
 - $P_{\text{unigram}}(t_1 t_2 t_3 t_4) = P(t_1) \times P(t_2) \times P(t_3) \times P(t_4)$

Language Models

- Bigram language models condition the probability of each term on the previous item:
 - $P_{\text{bigram}}(t_1 t_2 t_3 t_4) = P(t_1) \times P(t_2|t_1) \times P(t_3|t_2) \times P(t_4|t_3)$
- More complex language models are important in tasks such as speech recognition, spelling correction, or machine translation.
- In Information Retrieval most language-modeling work uses unigram language models. In IR, language models are often estimated from a single model so there is little information to do more.

Example of Language Models

- D1: Portugal eyes political balance in presidential election
- D2: After Portuguese elections, Spain braces for elections
- Unigram Language Models:
 - M_{d1} : portugal: 0.143 (1/7); eye: 0.143; ...
 - M_{d2} : after: 0.143; portugal: 0.143; election: 0.286 (2/7); ...

Retrieval based on Language Models

- Approach for retrieving documents under a language model (LM):
 - 1. Infer a LM for each document.
 - 2. Estimate $P(q|M_{di})$, the probability of generating the query according to each of these document models.
 - 3. Rank the documents according to these probabilities.

Example of Retrieval with Language Models

→ D1: Portugal eyes political balance in presidential election

→ D2: After Portuguese elections, Spain braces for elections

→ Q: [portugal election]

→ $P(q|d_1) = P(\text{portugal}|M_{d1}) \times P(\text{election}|M_{d1}) = 1/7 \times 1/7 = 0.0204$

→ $P(q|d_2) = P(\text{portugal}|M_{d2}) \times P(\text{election}|M_{d2}) = 1/7 \times 2/7 = 0.0408$

Taxonomy of IR Models

