# 6 *Scoring, term weighting and the vector space model*

Thus far we have dealt with indexes that support Boolean queries: a document either matches or does not match a query. In the case of large document collections, the resulting number of matching documents can far exceed the number a human user could possibly sift through. Accordingly, it is essential for a search engine to rank-order the documents matching a query. To do this, the search engine computes, for each matching document, a score with respect to the query at hand. In this chapter we initiate the study of assigning a score to a (query, document) pair. This chapter consists of three main ideas.

1. We introduce parametric and zone indexes in Section 6.1, which serve two purposes. First, they allow us to index and retrieve documents by metadata such as the language in which a document is written. Second, they give us a simple means for scoring (and thereby ranking) documents in response to a query.

2. Next, in Section 6.2 we develop the idea of weighting the importance of a term in a document, based on the statistics of occurrence of the term.

3. In Section 6.3 we show that by viewing each document as a vector of such weights, we can compute a score between a query and each document. This view is known as vector space scoring.

Section 6.4 develops several variants of term-weighting for the vector space model. Chapter 7 develops computational aspects of vector space scoring, and related topics.

As we develop these ideas, the notion of a query will assume multiple nuances. In Section 6.1 we consider queries in which specific query terms occur in specified regions of a matching document. Beginning Section 6.2 we will in fact relax the requirement of matching specific regions of a document; instead, we will look at so-called free text queries that simply consist of query terms with no specification on their relative order, importance or where in a document they should be found. The bulk of our study of scoring will be in this latter notion of a query being such a set of terms.
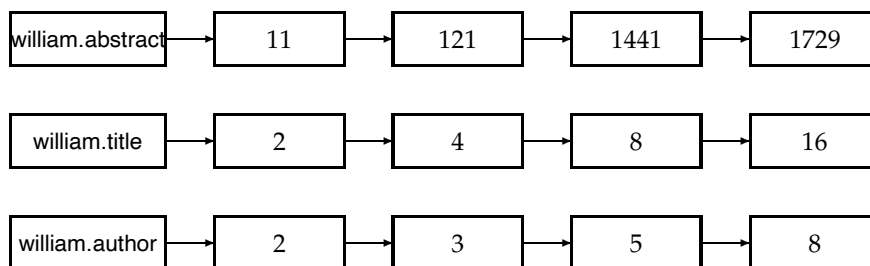
## 6.1    Parametric and zone indexes

METADATA

FIELD

PARAMETRIC INDEX

ZONE

WEIGHTED ZONE
SCORING

We have thus far viewed a document as a sequence of terms. In fact, most documents have additional structure. Digital documents generally encode, in machine-recognizable form, certain *metadata* associated with each document. By metadata, we mean specific forms of data about a document, such as its author(s), title and date of publication. This metadata would generally include *fields* such as the date of creation and the format of the document, as well the author and possibly the title of the document. The possible values of a field should be thought of as finite – for instance, the set of all dates of authorship.

Consider queries of the form "find documents authored by William Shakespeare in 1601, containing the phrase alas poor Yorick". Query processing then consists as usual of postings intersections, except that we may merge postings from standard inverted as well as *parametric indexes*. There is one parametric index for each field (say, date of creation); it allows us to select only the documents matching a date specified in the query. Figure 6.1 illustrates the user's view of such a parametric search. Some of the fields may assume ordered values, such as dates; in the example query above, the year 1601 is one such field value. The search engine may support querying ranges on such ordered values; to this end, a structure like a B-tree may be used for the field's dictionary.

*Zones* are similar to fields, except the contents of a zone can be arbitrary free text. Whereas a field may take on a relatively small set of values, a zone can be thought of as an arbitrary, unbounded amount of text. For instance, document titles and abstracts are generally treated as zones. We may build a separate inverted index for each zone of a document, to support queries such as "find documents with merchant in the title and william in the author list and the phrase gentle rain in the body". This has the effect of building an index that looks like Figure 6.2. Whereas the dictionary for a parametric index comes from a fixed vocabulary (the set of languages, or the set of dates), the dictionary for a zone index must structure whatever vocabulary stems from the text of that zone.

In fact, we can reduce the size of the dictionary by encoding the zone in which a term occurs in the postings. In Figure 6.3 for instance, we show how occurrences of william in the title and author zones of various documents are encoded. Such an encoding is useful when the size of the dictionary is a concern (because we require the dictionary to fit in main memory). But there is another important reason why the encoding of Figure 6.3 is useful: the efficient computation of scores using a technique we will call *weighted zone scoring*.

**Bibliographic Search**
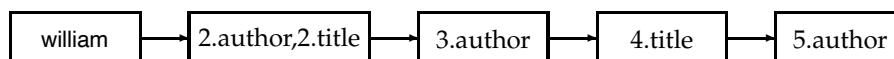
| Search category | Value |
|---|---|
| **Author** | *Example:* Widom, J *or* Garcia-Molina |
| **Title** | Also a part of the title possible |
| **Date of publication** | *Example*: `1997` *or* `<1997` *or* `>1997` limits the search to the documents appeared in, before and after 1997 respectively |
| **Language** | Language the document was written in  English |
| **Project** | ANY |
| **Type** | ANY |
| **Subject group** | ANY |
| **Sorted by** | Date of publication |
| | Start bibliographic search |

Find document via ID

► **Figure 6.1** Parametric search. In this example we have a collection with fields allowing us to select publications by zones such as Author and fields such as Language.

william.abstract → 11 → 121 → 1441 → 1729

william.title → 2 → 4 → 8 → 16

william.author → 2 → 3 → 5 → 8

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

william → 2.author,2.title → 3.author → 4.title → 5.author

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

### 6.1.1 Weighted zone scoring

Thus far in Section 6.1 we have focused on retrieving documents based on Boolean queries on fields and zones. We now turn to a second application of zones and fields.

Given a Boolean query $q$ and a document $d$, weighted zone scoring assigns to the pair $(q, d)$ a score in the interval $[0, 1]$, by computing a linear combination of *zone scores*, where each zone of the document contributes a Boolean value. More specifically, consider a set of documents each of which has $\ell$ zones. Let $g_1, \ldots, g_\ell \in [0, 1]$ such that $\sum_{i=1}^{\ell} g_i = 1$. For $1 \leq i \leq \ell$, let $s_i$ be the Boolean score denoting a match (or absence thereof) between $q$ and the $i$th zone. For instance, the Boolean score from a zone could be 1 if all the query term(s) occur in that zone, and zero otherwise; indeed, it could be any Boolean function that maps the presence of query terms in a zone to $0, 1$. Then, the weighted zone score is defined to be

$$(6.1) \qquad\qquad \sum_{i=1}^{\ell} g_i s_i.$$

RANKED BOOLEAN RETRIEVAL    Weighted zone scoring is sometimes referred to also as *ranked Boolean retrieval*.

**Example 6.1:** Consider the query shakespeare in a collection in which each document has three zones: *author, title* and *body*. The Boolean score function for a zone takes on the value 1 if the query term shakespeare is present in the zone, and zero otherwise. Weighted zone scoring in such a collection would require three weights $g_1, g_2$ and $g_3$, respectively corresponding to the *author, title* and *body* zones. Suppose we set $g_1 = 0.2, g_2 = 0.3$ and $g_3 = 0.5$ (so that the three weights add up to 1); this corresponds to an application in which a match in the *author* zone is least important to the overall score, the *title* zone somewhat more, and the *body* contributes even more.

Thus if the term shakespeare were to appear in the *title* and *body* zones but not the *author* zone of a document, the score of this document would be 0.8.

How do we implement the computation of weighted zone scores? A simple approach would be to compute the score for each document in turn, adding in all the contributions from the various zones. However, we now show how we may compute weighted zone scores directly from inverted indexes. The algorithm of Figure 6.4 treats the case when the query $q$ is a two-term query consisting of query terms $q_1$ and $q_2$, and the Boolean function is AND: 1 if both query terms are present in a zone and 0 otherwise. Following the description of the algorithm, we describe the extension to more complex queries and Boolean functions.

The reader may have noticed the close similarity between this algorithm and that in Figure 1.6. Indeed, they represent the same postings traversal, except that instead of merely adding a document to the set of results for

ZONESCORE($q_1, q_2$)
1   float $scores[N] = [0]$
2   constant $g[\ell]$
3   $p_1 \leftarrow postings(q_1)$
4   $p_2 \leftarrow postings(q_2)$
5   // scores[] is an array with a score entry for each document, initialized to zero.
6   // $p_1$ and $p_2$ are initialized to point to the beginning of their respective postings.
7   //Assume g[] is initialized to the respective zone weights.
8   **while** $p_1 \neq$ NIL and $p_2 \neq$ NIL
9   **do if** $docID(p_1) = docID(p_2)$
10      **then** $scores[docID(p_1)] \leftarrow$ WEIGHTEDZONE($p_1, p_2, g$)
11          $p_1 \leftarrow next(p_1)$
12          $p_2 \leftarrow next(p_2)$
13      **else if** $docID(p_1) < docID(p_2)$
14          **then** $p_1 \leftarrow next(p_1)$
15          **else** $p_2 \leftarrow next(p_2)$
16   **return** $scores$

▶ **Figure 6.4** Algorithm for computing the weighted zone score from two postings lists. Function WEIGHTEDZONE (not shown here) is assumed to compute the inner loop of Equation 6.1.

ACCUMULATOR

a Boolean AND query, we now compute a score for each such document. Some literature refers to the array scores[] above as a set of *accumulators*. The reason for this will be clear as we consider more complex Boolean functions than the AND; thus we may assign a non-zero score to a document even if it does not contain all query terms.

### 6.1.2 Learning weights

MACHINE-LEARNED
RELEVANCE

How do we determine the weights $g_i$ for weighted zone scoring? These weights could be specified by an expert (or, in principle, the user); but increasingly, these weights are "learned" using training examples that have been judged editorially. This latter methodology falls under a general class of approaches to scoring and ranking in information retrieval, known as *machine-learned relevance*. We provide a brief introduction to this topic here because weighted zone scoring presents a clean setting for introducing it; a complete development demands an understanding of machine learning and is deferred to Chapter 15.

1. We are provided with a set of *training examples*, each of which is a tuple consisting of a query $q$ and a document $d$, together with a relevance

judgment for $d$ on $q$. In the simplest form, each relevance judgments is either *Relevant* or *Non-relevant*. More sophisticated implementations of the methodology make use of more nuanced judgments.

2. The weights $g_i$ are then "learned" from these examples, in order that the learned scores approximate the relevance judgments in the training examples.

For weighted zone scoring, the process may be viewed as learning a linear function of the Boolean match scores contributed by the various zones. The expensive component of this methodology is the labor-intensive assembly of user-generated relevance judgments from which to learn the weights, especially in a collection that changes frequently (such as the Web). We now detail a simple example that illustrates how we can reduce the problem of learning the weights $g_i$ to a simple optimization problem.

We now consider a simple case of weighted zone scoring, where each document has a *title* zone and a *body* zone. Given a query $q$ and a document $d$, we use the given Boolean match function to compute Boolean variables $s_T(d, q)$ and $s_B(d, q)$, depending on whether the title (respectively, body) zone of $d$ matches query $q$. For instance, the algorithm in Figure 6.4 uses an AND of the query terms for this Boolean function. We will compute a score between 0 and 1 for each (document, query) pair using $s_T(d, q)$ and $s_B(d, q)$ by using a constant $g \in [0, 1]$, as follows:

$$(6.2) \qquad score(d, q) = g \cdot s_T(d, q) + (1 - g)s_B(d, q).$$

We now describe how to determine the constant $g$ from a set of *training examples*, each of which is a triple of the form $\Phi_j = (d_j, q_j, r(d_j, q_j))$. In each training example, a given training document $d_j$ and a given training query $q_j$ are assessed by a human editor who delivers a relevance judgment $r(d_j, q_j)$ that is either *Relevant* or *Non-relevant*. This is illustrated in Figure 6.5, where seven training examples are shown.

For each training example $\Phi_j$ we have Boolean values $s_T(d_j, q_j)$ and $s_B(d_j, q_j)$ that we use to compute a score from (6.2)

$$(6.3) \qquad score(d_j, q_j) = g \cdot s_T(d_j, q_j) + (1 - g)s_B(d_j, q_j).$$

We now compare this computed score to the human relevance judgment for the same document-query pair $(d_j, q_j)$; to this end, we will quantize each *Relevant* judgment as a 1 and each *Non-relevant* judgment as a 0. Suppose that we define the error of the scoring function with weight $g$ as

$$\varepsilon(g, \Phi_j) = (r(d_j, q_j) - score(d_j, q_j))^2,$$

| Example | DocID | Query | $s_T$ | $s_B$ | Judgment |
|---------|-------|-------|-------|-------|----------|
| $\Phi_1$ | 37 | linux | 1 | 1 | Relevant |
| $\Phi_2$ | 37 | penguin | 0 | 1 | Non-relevant |
| $\Phi_3$ | 238 | system | 0 | 1 | Relevant |
| $\Phi_4$ | 238 | penguin | 0 | 0 | Non-relevant |
| $\Phi_5$ | 1741 | kernel | 1 | 1 | Relevant |
| $\Phi_6$ | 2094 | driver | 0 | 1 | Relevant |
| $\Phi_7$ | 3191 | driver | 1 | 0 | Non-relevant |

▶ **Figure 6.5** An illustration of training examples.

| $s_T$ | $s_B$ | Score |
|-------|-------|-------|
| 0 | 0 | 0 |
| 0 | 1 | $1 - g$ |
| 1 | 0 | $g$ |
| 1 | 1 | 1 |

▶ **Figure 6.6** The four possible combinations of $s_T$ and $s_B$.

where we have quantized the editorial relevance judgment $r(d_j, q_j)$ to 0 or 1. Then, the total error of a set of training examples is given by

$$(6.4) \qquad \sum_j \varepsilon(g, \Phi_j).$$

The problem of learning the constant $g$ from the given training examples then reduces to picking the value of $g$ that minimizes the total error in (6.4).

Picking the best value of $g$ in (6.4) in the formulation of Section 6.1.3 reduces to the problem of minimizing a quadratic function of $g$ over the interval $[0, 1]$. This reduction is detailed in Section 6.1.3.

### 6.1.3 The optimal weight $g$

We begin by noting that for any training example $\Phi_j$ for which $s_T(d_j, q_j) = 0$ and $s_B(d_j, q_j) = 1$, the score computed by Equation (6.2) is $1 - g$. In similar fashion, we may write down the score computed by Equation (6.2) for the three other possible combinations of $s_T(d_j, q_j)$ and $s_B(d_j, q_j)$; this is summarized in Figure 6.6.

Let $n_{01r}$ (respectively, $n_{01n}$) denote the number of training examples for which $s_T(d_j, q_j) = 0$ and $s_B(d_j, q_j) = 1$ and the editorial judgment is *Relevant* (respectively, *Non-relevant*). Then the contribution to the total error in Equation (6.4) from training examples for which $s_T(d_j, q_j) = 0$ and $s_B(d_j, q_j) = 1$

is

$$[1 - (1 - g)]^2 n_{01r} + [0 - (1 - g)]^2 n_{01n}.$$

By writing in similar fashion the error contributions from training examples of the other three combinations of values for $s_T(d_j, q_j)$ and $s_B(d_j, q_j)$ (and extending the notation in the obvious manner), the total error corresponding to Equation (6.4) is

$$(6.5) \qquad (n_{01r} + n_{10n})g^2 + (n_{10r} + n_{01n})(1 - g)^2 + n_{00r} + n_{11n}.$$

By differentiating Equation (6.5) with respect to $g$ and setting the result to zero, it follows that the optimal value of $g$ is

$$(6.6) \qquad \frac{n_{10r} + n_{01n}}{n_{10r} + n_{10n} + n_{01r} + n_{01n}}.$$

**?**

**Exercise 6.1**

When using weighted zone scoring, is it necessary for all zones to use the same Boolean match function?

**Exercise 6.2**

In Example 6.1 above with weights $g_1 = 0.2, g_2 = 0.31$ and $g_3 = 0.49$, what are all the distinct score values a document may get?

**Exercise 6.3**

Rewrite the algorithm in Figure 6.4 to the case of more than two query terms.

**Exercise 6.4**

Write pseudocode for the function WeightedZone for the case of two postings lists in Figure 6.4.

**Exercise 6.5**

Apply Equation 6.6 to the sample training set in Figure 6.5 to estimate the best value of $g$ for this sample.

**Exercise 6.6**

For the value of $g$ estimated in Exercise 6.5, compute the weighted zone score for each (query, document) example. How do these scores relate to the relevance judgments in Figure 6.5 (quantized to 0/1)?

**Exercise 6.7**

Why does the expression for $g$ in (6.6) not involve training examples in which $s_T(d_t, q_t)$ and $s_B(d_t, q_t)$ have the same value?

## 6.2 Term frequency and weighting

Thus far, scoring has hinged on whether or not a query term is present in a zone within a document. We take the next logical step: a document or zone that mentions a query term more often has more to do with that query and therefore should receive a higher score. To motivate this, we recall the notion of a free text query introduced in Section 1.4: a query in which the terms of the query are typed freeform into the search interface, without any connecting search operators (such as Boolean operators). This query style, which is extremely popular on the web, views the query as simply a set of words. A plausible scoring mechanism then is to compute a score that is the sum, over the query terms, of the match scores between each query term and the document.

Towards this end, we assign to each term in a document a *weight* for that term, that depends on the number of occurrences of the term in the document. We would like to compute a score between a query term $t$ and a document $d$, based on the weight of $t$ in $d$. The simplest approach is to assign the weight to be equal to the number of occurrences of term $t$ in document $d$. This weighting scheme is referred to as *term frequency* and is denoted $\mathrm{tf}_{t,d}$, with the subscripts denoting the term and the document in order.

TERM FREQUENCY

For a document $d$, the set of weights determined by the tf weights above (or indeed any weighting function that maps the number of occurrences of $t$ in $d$ to a positive real value) may be viewed as a quantitative digest of that document. In this view of a document, known in the literature as the *bag of words model*, the exact ordering of the terms in a document is ignored but the number of occurrences of each term is material (in contrast to Boolean retrieval). We only retain information on the number of occurrences of each term. Thus, the document "Mary is quicker than John" is, in this view, identical to the document "John is quicker than Mary". Nevertheless, it seems intuitive that two documents with similar bag of words representations are similar in content. We will develop this intuition further in Section 6.3.

BAG OF WORDS

Before doing so we first study the question: are all words in a document equally important? Clearly not; in Section 2.2.2 (page 27) we looked at the idea of *stop words* – words that we decide not to index at all, and therefore do not contribute in any way to retrieval and scoring.

### 6.2.1 Inverse document frequency

Raw term frequency as above suffers from a critical problem: all terms are considered equally important when it comes to assessing relevancy on a query. In fact certain terms have little or no discriminating power in determining relevance. For instance, a collection of documents on the auto industry is likely to have the term auto in almost every document. To this

| Word      | cf    | df   |
|-----------|-------|------|
| try       | 10422 | 8760 |
| insurance | 10440 | 3997 |

▶ **Figure 6.7** Collection frequency (cf) and document frequency (df) behave differently, as in this example from the Reuters collection.

end, we introduce a mechanism for attenuating the effect of terms that occur too often in the collection to be meaningful for relevance determination. An immediate idea is to scale down the term weights of terms with high *collection frequency*, defined to be the total number of occurrences of a term in the collection. The idea would be to reduce the tf weight of a term by a factor that grows with its collection frequency.

DOCUMENT
FREQUENCY

Instead, it is more commonplace to use for this purpose the *document frequency* $\mathrm{df}_t$, defined to be the number of documents in the collection that contain a term $t$. This is because in trying to discriminate between documents for the purpose of scoring it is better to use a document-level statistic (such as the number of documents containing a term) than to use a collection-wide statistic for the term. The reason to prefer df to cf is illustrated in Figure 6.7, where a simple example shows that collection frequency (cf) and document frequency (df) can behave rather differently. In particular, the cf values for both try and insurance are roughly equal, but their df values differ significantly. Intuitively, we want the few documents that contain insurance to get a higher boost for a query on insurance than the many documents containing try get from a query on try.

How is the document frequency df of a term used to scale its weight? Denoting as usual the total number of documents in a collection by $N$, we define

INVERSE DOCUMENT
FREQUENCY

the *inverse document frequency* (idf) of a term $t$ as follows:

(6.7)
$$\mathrm{idf}_t = \log \frac{N}{\mathrm{df}_t}.$$

Thus the idf of a rare term is high, whereas the idf of a frequent term is likely to be low. Figure 6.8 gives an example of idf's in the Reuters collection of 806,791 documents; in this example logarithms are to the base 10. In fact, as we will see in Exercise 6.12, the precise base of the logarithm is not material to ranking. We will give on page 227 a justification of the particular form in Equation (6.7).

### 6.2.2 Tf-idf weighting

We now combine the definitions of term frequency and inverse document frequency, to produce a composite weight for each term in each document.

| term | $df_t$ | $idf_t$ |
|------|-------:|--------:|
| car | 18,165 | 1.65 |
| auto | 6723 | 2.08 |
| insurance | 19,241 | 1.62 |
| best | 25,235 | 1.5 |

▶ **Figure 6.8** Example of idf values. Here we give the idf's of terms with various frequencies in the Reuters collection of 806,791 documents.

TF-IDF    The *tf-idf* weighting scheme assigns to term $t$ a weight in document $d$ given by

(6.8)
$$\text{tf-idf}_{t,d} = \text{tf}_{t,d} \times \text{idf}_t.$$

In other words, $\text{tf-idf}_{t,d}$ assigns to term $t$ a weight in document $d$ that is

1. highest when $t$ occurs many times within a small number of documents (thus lending high discriminating power to those documents);

2. lower when the term occurs fewer times in a document, or occurs in many documents (thus offering a less pronounced relevance signal);

3. lowest when the term occurs in virtually all documents.

DOCUMENT VECTOR    At this point, we may view each document as a *vector* with one component corresponding to each term in the dictionary, together with a weight for each component that is given by (6.8). For dictionary terms that do not occur in a document, this weight is zero. This vector form will prove to be crucial to scoring and ranking; we will develop these ideas in Section 6.3. As a first step, we introduce the *overlap score measure*: the score of a document $d$ is the sum, over all query terms, of the number of times each of the query terms occurs in $d$. We can refine this idea so that we add up not the number of occurrences of each query term $t$ in $d$, but instead the tf-idf weight of each term in $d$.

(6.9)
$$\text{Score}(q,d) = \sum_{t \in q} \text{tf-idf}_{t,d}.$$

In Section 6.3 we will develop a more rigorous form of Equation (6.9).

**?** **Exercise 6.8**

Why is the idf of a term always finite?

**Exercise 6.9**

What is the idf of a term that occurs in every document? Compare this with the use of stop word lists.

|           | Doc1 | Doc2 | Doc3 |
|-----------|------|------|------|
| car       | 27   | 4    | 24   |
| auto      | 3    | 33   | 0    |
| insurance | 0    | 33   | 29   |
| best      | 14   | 0    | 17   |

▶ **Figure 6.9**   Table of tf values for Exercise 6.10.

**Exercise 6.10**

Consider the table of term frequencies for 3 documents denoted Doc1, Doc2, Doc3 in Figure 6.9. Compute the tf-idf weights for the terms car, auto, insurance, best, for each document, using the idf values from Figure 6.8.

**Exercise 6.11**

Can the tf-idf weight of a term in a document exceed 1?

**Exercise 6.12**

How does the base of the logarithm in (6.7) affect the score calculation in (6.9)? How does the base of the logarithm affect the relative scores of two documents on a given query?

**Exercise 6.13**

If the logarithm in (6.7) is computed base 2, suggest a simple approximation to the idf of a term.
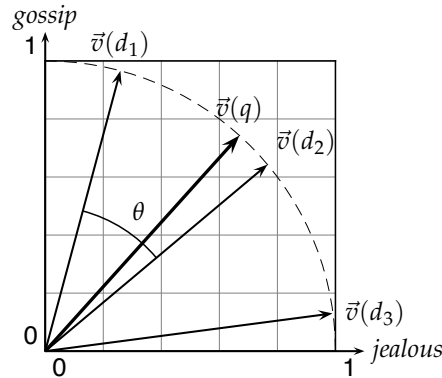
## 6.3  The vector space model for scoring

In Section 6.2 (page 117) we developed the notion of a document vector that captures the relative importance of the terms in a document. The representation of a set of documents as vectors in a common vector space is known as VECTOR SPACE MODEL the *vector space model* and is fundamental to a host of information retrieval operations ranging from scoring documents on a query, document classification and document clustering. We first develop the basic ideas underlying vector space scoring; a pivotal step in this development is the view (Section 6.3.2) of queries as vectors in the same vector space as the document collection.

### 6.3.1  Dot products

We denote by $\vec{V}(d)$ the vector derived from document $d$, with one component in the vector for each dictionary term. Unless otherwise specified, the reader may assume that the components are computed using the tf-idf weighting scheme, although the particular weighting scheme is immaterial to the discussion that follows. The set of documents in a collection then may be viewed as a set of vectors in a vector space, in which there is one axis for

► **Figure 6.10** Cosine similarity illustrated. $\text{sim}(d_1, d_2) = \cos \theta$.

each term. This representation loses the relative ordering of the terms in each document; recall our example from Section 6.2 (page 117), where we pointed out that the documents *Mary is quicker than John* and *John is quicker than Mary* are identical in such a *bag of words* representation.

How do we quantify the similarity between two documents in this vector space? A first attempt might consider the magnitude of the vector difference between two document vectors. This measure suffers from a drawback: two documents with very similar content can have a significant vector difference simply because one is much longer than the other. Thus the relative distributions of terms may be identical in the two documents, but the absolute term frequencies of one may be far larger.

To compensate for the effect of document length, the standard way of quantifying the similarity between two documents $d_1$ and $d_2$ is to compute the *cosine similarity* of their vector representations $\vec{V}(d_1)$ and $\vec{V}(d_2)$

COSINE SIMILARITY

(6.10)
$$\text{sim}(d_1, d_2) = \frac{\vec{V}(d_1) \cdot \vec{V}(d_2)}{|\vec{V}(d_1)||\vec{V}(d_2)|},$$

DOT PRODUCT

where the numerator represents the *dot product* (also known as the *inner product*) of the vectors $\vec{V}(d_1)$ and $\vec{V}(d_2)$, while the denominator is the product of their *Euclidean lengths*. The dot product $\vec{x} \cdot \vec{y}$ of two vectors is defined as $\sum_{i=1}^{M} x_i y_i$. Let $\vec{V}(d)$ denote the document vector for $d$, with $M$ components $\vec{V}_1(d) \dots \vec{V}_M(d)$. The Euclidean length of $d$ is defined to be $\sqrt{\sum_{i=1}^{M} \vec{V}_i^2(d)}$.

EUCLIDEAN LENGTH

LENGTH-NORMALIZATION

The effect of the denominator of Equation (6.10) is thus to *length-normalize* the vectors $\vec{V}(d_1)$ and $\vec{V}(d_2)$ to unit vectors $\vec{v}(d_1) = \vec{V}(d_1)/|\vec{V}(d_1)|$ and

|           | Doc1 | Doc2 | Doc3 |
|-----------|------|------|------|
| car       | 0.88 | 0.09 | 0.58 |
| auto      | 0.10 | 0.71 | 0    |
| insurance | 0    | 0.71 | 0.70 |
| best      | 0.46 | 0    | 0.41 |

▶ **Figure 6.11**  Euclidean normalized tf values for documents in Figure 6.9.

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

▶ **Figure 6.12**   Term frequencies in three novels.  The novels are Austen's *Sense and Sensibility*, *Pride and Prejudice* and Brontë's *Wuthering Heights.*

$\vec{v}(d_2) = \vec{V}(d_2)/|\vec{V}(d_2)|$. We can then rewrite (6.10) as

(6.11)                                    $$\text{sim}(d_1, d_2) = \vec{v}(d_1) \cdot \vec{v}(d_2).$$

**Example 6.2:**    Consider the documents in Figure 6.9.  We now apply Euclidean normalization to the tf values from the table, for each of the three documents in the table. The quantity $\sqrt{\sum_{i=1}^{M} \vec{V}_i^2(d)}$ has the values 30.56, 46.84 and 41.30 respectively for Doc1, Doc2 and Doc3.  The resulting Euclidean normalized tf values for these documents are shown in Figure 6.11.

Thus, (6.11) can be viewed as the dot product of the normalized versions of the two document vectors. This measure is the cosine of the angle $\theta$ between the two vectors, shown in Figure 6.10. What use is the similarity measure $\text{sim}(d_1, d_2)$? Given a document $d$ (potentially one of the $d_i$ in the collection), consider searching for the documents in the collection most similar to $d$. Such a search is useful in a system where a user may identify a document and seek others like it – a feature available in the results lists of search engines as a *more like this* feature. We reduce the problem of finding the document(s) most similar to $d$ to that of finding the $d_i$ with the highest dot products (sim values) $\vec{v}(d) \cdot \vec{v}(d_i)$. We could do this by computing the dot products between $\vec{v}(d)$ and each of $\vec{v}(d_1), \ldots, \vec{v}(d_N)$, then picking off the highest resulting sim values.

**Example 6.3:**    Figure 6.12 shows the number of occurrences of three terms (affection, jealous and gossip) in each of the following three novels: Jane Austen's *Sense and Sensibility* (SaS) and *Pride and Prejudice* (PaP) and Emily Brontë's *Wuthering Heights* (WH).

| term | SaS | PaP | WH |
|------|-----|-----|-----|
| affection | 0.996 | 0.993 | 0.847 |
| jealous | 0.087 | 0.120 | 0.466 |
| gossip | 0.017 | 0 | 0.254 |

▶ **Figure 6.13** Term vectors for the three novels of Figure 6.12. These are based on raw term frequency only and are normalized as if these were the only terms in the collection. (Since affection and jealous occur in all three documents, their tf-idf weight would be 0 in most formulations.)

Of course, there are many other terms occurring in each of these novels. In this example we represent each of these novels as a unit vector in three dimensions, corresponding to these three terms (only); we use raw term frequencies here, with no idf multiplier. The resulting weights are as shown in Figure 6.13.

Now consider the cosine similarities between pairs of the resulting three-dimensional vectors. A simple computation shows that $\text{sim}(\vec{v}(SAS), \vec{v}(PAP))$ is 0.999, whereas $\text{sim}(\vec{v}(SAS), \vec{v}(WH))$ is 0.888; thus, the two books authored by Austen (SaS and PaP) are considerably closer to each other than to Brontë's *Wuthering Heights*. In fact, the similarity between the first two is almost perfect (when restricted to the three terms we consider). Here we have considered tf weights, but we could of course use other term weight functions.

TERM-DOCUMENT MATRIX Viewing a collection of $N$ documents as a collection of vectors leads to a natural view of a collection as a *term-document matrix*: this is an $M \times N$ matrix whose rows represent the $M$ terms (dimensions) of the $N$ columns, each of which corresponds to a document. As always, the terms being indexed could be stemmed before indexing; for instance, jealous and jealousy would under stemming be considered as a single dimension. This matrix view will prove to be useful in Chapter 18.

### 6.3.2 Queries as vectors

There is a far more compelling reason to represent documents as vectors: we can also view a *query* as a vector. Consider the query $q = $ jealous gossip. This query turns into the unit vector $\vec{v}(q) = (0, 0.707, 0.707)$ on the three coordinates of Figures 6.12 and 6.13. The key idea now: to assign to each document $d$ a score equal to the dot product

$$\vec{v}(q) \cdot \vec{v}(d).$$

In the example of Figure 6.13, *Wuthering Heights* is the top-scoring document for this query with a score of 0.509, with *Pride and Prejudice* a distant second with a score of 0.085, and *Sense and Sensibility* last with a score of 0.074. This simple example is somewhat misleading: the number of dimen-

sions in practice will be far larger than three: it will equal the vocabulary size $M$.

To summarize, by viewing a query as a "bag of words", we are able to treat it as a very short document. As a consequence, we can use the cosine similarity between the query vector and a document vector as a measure of the score of the document for that query. The resulting scores can then be used to select the top-scoring documents for a query. Thus we have

(6.12)
$$\text{score}(q,d) = \frac{\vec{V}(q) \cdot \vec{V}(d)}{|\vec{V}(q)||\vec{V}(d)|}.$$

A document may have a high cosine score for a query even if it does not contain all query terms. Note that the preceding discussion does not hinge on any specific weighting of terms in the document vector, although for the present we may think of them as either tf or tf-idf weights. In fact, a number of weighting schemes are possible for query as well as document vectors, as illustrated in Example 6.4 and developed further in Section 6.4.

Computing the cosine similarities between the query vector and each document vector in the collection, sorting the resulting scores and selecting the top $K$ documents can be expensive — a single similarity computation can entail a dot product in tens of thousands of dimensions, demanding tens of thousands of arithmetic operations. In Section 7.1 we study how to use an inverted index for this purpose, followed by a series of heuristics for improving on this.

**Example 6.4:** We now consider the query best car insurance on a fictitious collection with $N = 1,000,000$ documents where the document frequencies of auto, best, car and insurance are respectively 5000, 50000, 10000 and 1000.

| term | query | | | | document | | | product |
|------|----|----|----|----------|----|----|----------|----|
|      | tf | df | idf | $w_{t,q}$ | tf | wf | $w_{t,d}$ |    |
| auto | 0 | 5000 | 2.3 | 0 | 1 | 1 | 0.41 | 0 |
| best | 1 | 50000 | 1.3 | 1.3 | 0 | 0 | 0 | 0 |
| car | 1 | 10000 | 2.0 | 2.0 | 1 | 1 | 0.41 | 0.82 |
| insurance | 1 | 1000 | 3.0 | 3.0 | 2 | 2 | 0.82 | 2.46 |

In this example the weight of a term in the query is simply the idf (and zero for a term not in the query, such as auto); this is reflected in the column header $w_{t,q}$ (the entry for auto is zero because the query does not contain the term auto). For documents, we use tf weighting with no use of idf but with Euclidean normalization. The former is shown under the column headed wf, while the latter is shown under the column headed $w_{t,d}$. Invoking (6.9) now gives a net score of $0 + 0 + 0.82 + 2.46 = 3.28$.

### 6.3.3 Computing vector scores

In a typical setting we have a collection of documents each represented by a vector, a free text query represented by a vector, and a positive integer $K$. We

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

▶ **Figure 6.14** The basic algorithm for computing vector space scores.

seek the $K$ documents of the collection with the highest vector space scores on the given query. We now initiate the study of determining the $K$ documents with the highest vector space scores for a query. Typically, we seek these $K$ top documents in ordered by decreasing score; for instance many search engines use $K = 10$ to retrieve and rank-order the first page of the ten best results. Here we give the basic algorithm for this computation; we develop a fuller treatment of efficient techniques and approximations in Chapter 7.

Figure 6.14 gives the basic algorithm for computing vector space scores. The array Length holds the lengths (normalization factors) for each of the $N$ documents, whereas the array Scores holds the scores for each of the documents. When the scores are finally computed in Step 9, all that remains in Step 10 is to pick off the $K$ documents with the highest scores.

The outermost loop beginning Step 3 repeats the updating of Scores, iterating over each query term $t$ in turn. In Step 5 we calculate the weight in the query vector for term $t$. Steps 6-8 update the score of each document by adding in the contribution from term $t$. This process of adding in contribu-
TERM-AT-A-TIME tions one query term at a time is sometimes known as *term-at-a-time* scoring or accumulation, and the $N$ elements of the array *Scores* are therefore known
ACCUMULATOR as *accumulators*. For this purpose, it would appear necessary to store, with each postings entry, the weight wf$_{t,d}$ of term $t$ in document $d$ (we have thus far used either tf or tf-idf for this weight, but leave open the possibility of other functions to be developed in Section 6.4). In fact this is wasteful, since storing this weight may require a floating point number. Two ideas help alleviate this space problem. First, if we are using inverse document frequency, we need not precompute idf$_t$; it suffices to store $N/\text{df}_t$ at the head of the postings for $t$. Second, we store the term frequency tf$_{t,d}$ for each postings entry. Finally, Step 12 extracts the top $K$ scores – this requires a priority queue

data structure, often implemented using a heap. Such a heap takes no more than $2N$ comparisons to construct, following which each of the $K$ top scores can be extracted from the heap at a cost of $O(\log N)$ comparisons.

Note that the general algorithm of Figure 6.14 does not prescribe a specific implementation of how we traverse the postings lists of the various query terms; we may traverse them one term at a time as in the loop beginning at Step 3, or we could in fact traverse them concurrently as in Figure 1.6. In such a concurrent postings traversal we compute the scores of one document at a time, so that it is sometimes called *document-at-a-time* scoring. We will say more about this in Section 7.1.5.

DOCUMENT-AT-A-TIME

**?**

**Exercise 6.14**

If we were to stem jealous and jealousy to a common stem before setting up the vector space, detail how the definitions of tf and idf should be modified.

**Exercise 6.15**

Recall the tf-idf weights computed in Exercise 6.10. Compute the Euclidean normalized document vectors for each of the documents, where each vector has four components, one for each of the four terms.

**Exercise 6.16**

Verify that the sum of the squares of the components of each of the document vectors in Exercise 6.15 is 1 (to within rounding error). Why is this the case?

**Exercise 6.17**

With term weights as computed in Exercise 6.15, rank the three documents by computed score for the query car insurance, for each of the following cases of term weighting in the query:

1. The weight of a term is 1 if present in the query, 0 otherwise.

2. Euclidean normalized idf.

## 6.4  Variant tf-idf functions

For assigning a weight for each term in each document, a number of alternatives to tf and tf-idf have been considered. We discuss some of the principal ones here; a more complete development is deferred to Chapter 11. We will summarize these alternatives in Section 6.4.3 (page 128).

### 6.4.1  Sublinear tf scaling

It seems unlikely that twenty occurrences of a term in a document truly carry twenty times the significance of a single occurrence. Accordingly, there has been considerable research into variants of term frequency that go beyond counting the number of occurrences of a term. A common modification is

to use instead the logarithm of the term frequency, which assigns a weight given by

(6.13)
$$\text{wf}_{t,d} = \begin{cases} 1 + \log \text{tf}_{t,d} & \text{if } \text{tf}_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}.$$

In this form, we may replace tf by some other function wf as in (6.13), to obtain:

(6.14)
$$\text{wf-idf}_{t,d} = \text{wf}_{t,d} \times \text{idf}_t.$$

Equation (6.9) can then be modified by replacing tf-idf by wf-idf as defined in (6.14).

### 6.4.2 Maximum tf normalization

One well-studied technique is to normalize the tf weights of all terms occurring in a document by the maximum tf in that document. For each document $d$, let $\text{tf}_{\max}(d) = \max_{\tau \in d} \text{tf}_{\tau,d}$, where $\tau$ ranges over all terms in $d$. Then, we compute a normalized term frequency for each term $t$ in document $d$ by

(6.15)
$$\text{ntf}_{t,d} = a + (1-a)\frac{\text{tf}_{t,d}}{\text{tf}_{\max}(d)},$$

SMOOTHING where $a$ is a value between 0 and 1 and is generally set to 0.4, although some early work used the value 0.5. The term $a$ in (6.15) is a *smoothing* term whose role is to damp the contribution of the second term – which may be viewed as a scaling down of tf by the largest tf value in $d$. We will encounter smoothing further in Chapter 13 when discussing classification; the basic idea is to avoid a large swing in $\text{ntf}_{t,d}$ from modest changes in $\text{tf}_{t,d}$ (say from 1 to 2). The main idea of maximum tf normalization is to mitigate the following anomaly: we observe higher term frequencies in longer documents, merely because longer documents tend to repeat the same words over and over again. To appreciate this, consider the following extreme example: supposed we were to take a document $d$ and create a new document $d'$ by simply appending a copy of $d$ to itself. While $d'$ should be no more relevant to any query than $d$ is, the use of (6.9) would assign it twice as high a score as $d$. Replacing tf-idf$_{t,d}$ in (6.9) by ntf-idf$_{t,d}$ eliminates the anomaly in this example. Maximum tf normalization does suffer from the following issues:

1. The method is unstable in the following sense: a change in the stop word list can dramatically alter term weightings (and therefore ranking). Thus, it is hard to tune.

2. A document may contain an outlier term with an unusually large number of occurrences of that term, not representative of the content of that document.

| Term frequency | | Document frequency | | Normalization | |
|---|---|---|---|---|---|
| n (natural) | $\text{tf}_{t,d}$ | n (no) | 1 | n (none) | 1 |
| l (logarithm) | $1 + \log(\text{tf}_{t,d})$ | t (idf) | $\log \frac{N}{\text{df}_t}$ | c (cosine) | $\frac{1}{\sqrt{w_1^2 + w_2^2 + \ldots + w_M^2}}$ |
| a (augmented) | $0.5 + \frac{0.5 \times \text{tf}_{t,d}}{\max_t(\text{tf}_{t,d})}$ | p (prob idf) | $\max\{0, \log \frac{N - \text{df}_t}{\text{df}_t}\}$ | u (pivoted unique) | $1/u$ (Section 6.4.4) |
| b (boolean) | $\begin{cases} 1 & \text{if } \text{tf}_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$ | | | b (byte size) | $1/CharLength^{\alpha}, \alpha < 1$ |
| L (log ave) | $\frac{1 + \log(\text{tf}_{t,d})}{1 + \log(\text{ave}_{t \in d}(\text{tf}_{t,d}))}$ | | | | |

▶ **Figure 6.15** SMART notation for tf-idf variants. Here *CharLength* is the number of characters in the document.

3. More generally, a document in which the most frequent term appears roughly as often as many other terms should be treated differently from one with a more skewed distribution.

### 6.4.3  Document and query weighting schemes

Equation (6.12) is fundamental to information retrieval systems that use any form of vector space scoring. Variations from one vector space scoring method to another hinge on the specific choices of weights in the vectors $\vec{V}(d)$ and $\vec{V}(q)$. Figure 6.15 lists some of the principal weighting schemes in use for each of $\vec{V}(d)$ and $\vec{V}(q)$, together with a mnemonic for representing a specific combination of weights; this system of mnemonics is sometimes called SMART notation, following the authors of an early text retrieval system. The mnemonic for representing a combination of weights takes the form *ddd.qqq* where the first triplet gives the term weighting of the document vector, while the second triplet gives the weighting in the query vector. The first letter in each triplet specifies the term frequency component of the weighting, the second the document frequency component, and the third the form of normalization used. It is quite common to apply different normalization functions to $\vec{V}(d)$ and $\vec{V}(q)$. For example, a very standard weighting scheme is *lnc.ltc*, where the document vector has log-weighted term frequency, no idf (for both effectiveness and efficiency reasons), and cosine normalization, while the query vector uses log-weighted term frequency, idf weighting, and cosine normalization.
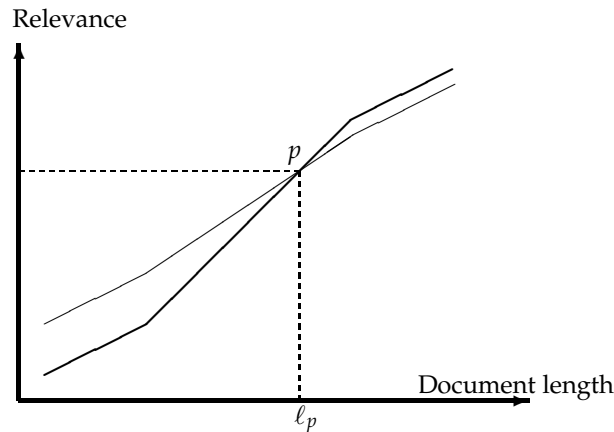
### 6.4.4 Pivoted normalized document length

In Section 6.3.1 we normalized each document vector by the Euclidean length of the vector, so that all document vectors turned into unit vectors. In doing so, we eliminated all information on the length of the original document; this masks some subtleties about longer documents. First, longer documents will – as a result of containing more terms – have higher tf values. Second, longer documents contain more distinct terms. These factors can conspire to raise the scores of longer documents, which (at least for some information needs) is unnatural. Longer documents can broadly be lumped into two categories: (1) *verbose* documents that essentially repeat the same content – in these, the length of the document does not alter the relative weights of different terms; (2) documents covering multiple different topics, in which the search terms probably match small segments of the document but not all of it – in this case, the relative weights of terms are quite different from a single short document that matches the query terms. Compensating for this phenomenon is a form of document length normalization that is independent of term and document frequencies. To this end, we introduce a form of normalizing the vector representations of documents in the collection, so that the resulting "normalized" documents are not necessarily of unit length. Then, when we compute the dot product score between a (unit) query vector and such a normalized document, the score is skewed to account for the effect of document length on relevance. This form of compensation for document length is known as *pivoted document length normalization*.

PIVOTED DOCUMENT LENGTH NORMALIZATION

Consider a document collection together with an ensemble of queries for that collection. Suppose that we were given, for each query $q$ and for each document $d$, a Boolean judgment of whether or not $d$ is relevant to the query $q$; in Chapter 8 we will see how to procure such a set of relevance judgments for a query ensemble and a document collection. Given this set of relevance judgments, we may compute a *probability of relevance* as a function of document length, averaged over all queries in the ensemble. The resulting plot may look like the curve drawn in thick lines in Figure 6.16. To compute this curve, we bucket documents by length and compute the fraction of relevant documents in each bucket, then plot this fraction against the median document length of each bucket. (Thus even though the "curve" in Figure 6.16 appears to be continuous, it is in fact a histogram of discrete buckets of document length.)

On the other hand, the curve in thin lines shows what might happen with the same documents and query ensemble if we were to use relevance as prescribed by cosine normalization Equation (6.12) – thus, cosine normalization has a tendency to distort the computed relevance vis-à-vis the true relevance, at the expense of longer documents. The thin and thick curves crossover at a point $p$ corresponding to document length $\ell_p$, which we refer to as the *pivot*
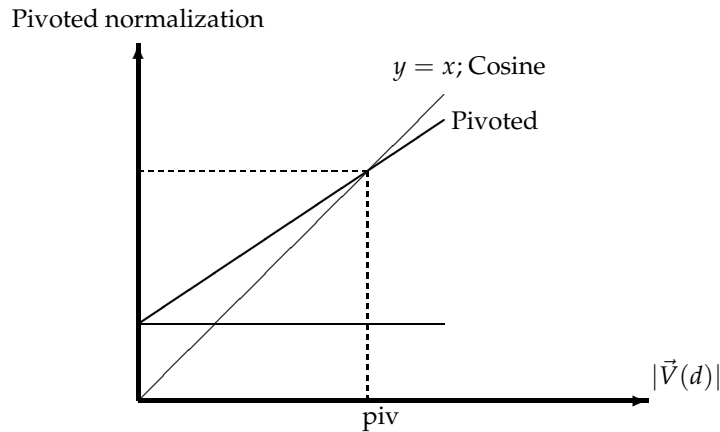
Relevance



► **Figure 6.16**   Pivoted document length normalization.

*length*; dashed lines mark this point on the $x-$ and $y-$ axes. The idea of pivoted document length normalization would then be to "rotate" the cosine normalization curve counter-clockwise about $p$ so that it more closely matches thick line representing the relevance vs. document length curve. As mentioned at the beginning of this section, we do so by using in Equation (6.12) a normalization factor for each document vector $\vec{V}(d)$ that is not the Euclidean length of that vector, but instead one that is larger than the Euclidean length for documents of length less than $\ell_p$, and smaller for longer documents.

To this end, we first note that the normalizing term for $\vec{V}(d)$ in the denominator of Equation (6.12) is its Euclidean length, denoted $|\vec{V}(d)|$. In the simplest implementation of pivoted document length normalization, we use a normalization factor in the denominator that is linear in $|\vec{V}(d)|$, but one of slope $< 1$ as in Figure 6.17. In this figure, the $x-$ axis represents $|\vec{V}(d)|$, while the $y-$axis represents possible normalization factors we can use. The thin line $y = x$ depicts the use of cosine normalization. Notice the following aspects of the thick line representing pivoted length normalization:

1. It is linear in the document length and has the form

(6.16)                                 $$a|\vec{V}(d)| + (1-a)\mathrm{piv},$$

Online edition (c) 2009 Cambridge UP

Pivoted normalization



► **Figure 6.17**   Implementing pivoted document length normalization by linear scaling.

where piv is the cosine normalization value at which the two curves intersect.

2. Its slope is $a < 1$ and (3) it crosses the $y = x$ line at piv.

It has been argued that in practice, Equation (6.16) is well approximated by

$$au_d + (1 - a)\text{piv},$$

where $u_d$ is the number of unique terms in document $d$.

   Of course, pivoted document length normalization is not appropriate for all applications. For instance, in a collection of answers to frequently asked questions (say, at a customer service website), relevance may have little to do with document length. In other cases the dependency may be more complex than can be accounted for by a simple linear pivoted normalization. In such cases, document length can be used as a feature in the machine learning based scoring approach of Section 6.1.2.

**Exercise 6.18**

EUCLIDEAN DISTANCE   One measure of the similarity of two vectors is the *Euclidean distance* (or $L_2$ distance) between them:

$$|\vec{x} - \vec{y}| = \sqrt{\sum_{i=1}^{M}(x_i - y_i)^2}$$

Online edition (c) 2009 Cambridge UP

| word | query | | | | | document | | | |
| | tf | wf | df | idf | $q_i = $ wf-idf | tf | wf | $d_i = $ normalized wf | $q_i \cdot d_i$ |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| digital | | | 10,000 | | | | | | |
| video | | | 100,000 | | | | | | |
| cameras | | | 50,000 | | | | | | |

▶ **Table 6.1**  Cosine computation for Exercise 6.19.

Given a query $q$ and documents $d_1, d_2, \ldots$, we may rank the documents $d_i$ in order of increasing Euclidean distance from $q$. Show that if $q$ and the $d_i$ are all normalized to unit vectors, then the rank ordering produced by Euclidean distance is identical to that produced by cosine similarities.

**Exercise 6.19**

Compute the vector space similarity between the query "digital cameras" and the document "digital cameras and video cameras" by filling out the empty columns in Table 6.1. Assume $N = 10,000,000$, logarithmic term weighting (wf columns) for query and document, idf weighting for the query only and cosine normalization for the document only. Treat and as a stop word. Enter term counts in the tf columns. What is the final similarity score?

**Exercise 6.20**

Show that for the query affection, the relative ordering of the scores of the three documents in Figure 6.13 is the reverse of the ordering of the scores for the query jealous gossip.

**Exercise 6.21**

In turning a query into a unit vector in Figure 6.13, we assigned equal weights to each of the query terms. What other principled approaches are plausible?

**Exercise 6.22**

Consider the case of a query term that is not in the set of $M$ indexed terms; thus our standard construction of the query vector results in $\vec{V}(q)$ not being in the vector space created from the collection. How would one adapt the vector space representation to handle this case?

**Exercise 6.23**

Refer to the tf and idf values for four terms and three documents in Exercise 6.10. Compute the two top scoring documents on the query best car insurance for each of the following weighing schemes: (i) `nnn.atc`; (ii) `ntc.atc`.

**Exercise 6.24**

Suppose that the word coyote does not occur in the collection used in Exercises 6.10 and 6.23. How would one compute `ntc.atc` scores for the query coyote insurance?

## 6.5   References and further reading

Chapter 7 develops the computational aspects of vector space scoring. Luhn (1957; 1958) describes some of the earliest reported applications of term weighting. His paper dwells on the importance of medium frequency terms (terms that are neither too commonplace nor too rare) and may be thought of as anticipating tf-idf and related weighting schemes. Spärck Jones (1972) builds on this intuition through detailed experiments showing the use of inverse document frequency in term weighting. A series of extensions and theoretical justifications of idf are due to Salton and Buckley (1987) Robertson and Jones (1976), Croft and Harper (1979) and Papineni (2001). Robertson maintains a web page (http://www.soi.city.ac.uk/˜ser/idf.html) containing the history of idf, including soft copies of early papers that predated electronic versions of journal article. Singhal et al. (1996a) develop pivoted document length normalization. Probabilistic language models (Chapter 11) develop weighting techniques that are more nuanced than tf-idf; the reader will find this development in Section 11.4.3.

We observed that by assigning a weight for each term in a document, a document may be viewed as a vector of term weights, one for each term in the collection. The SMART information retrieval system at Cornell (Salton 1971b) due to Salton and colleagues was perhaps the first to view a document as a vector of weights. The basic computation of cosine scores as described in Section 6.3.3 is due to Zobel and Moffat (2006). The two query evaluation strategies term-at-a-time and document-at-a-time are discussed by Turtle and Flood (1995).

The SMART notation for tf-idf term weighting schemes in Figure 6.15 is presented in (Salton and Buckley 1988, Singhal et al. 1995; 1996b). Not all versions of the notation are consistent; we most closely follow (Singhal et al. 1996b). A more detailed and exhaustive notation was developed in Moffat and Zobel (1998), considering a larger palette of schemes for term and document frequency weighting. Beyond the notation, Moffat and Zobel (1998) sought to set up a space of feasible weighting functions through which hill-climbing approaches could be used to begin with weighting schemes that performed well, then make local improvements to identify the best combinations. However, they report that such hill-climbing methods failed to lead to any conclusions on the best weighting schemes.