

# 11 Probabilistic information retrieval

During the discussion of relevance feedback in Section 9.1.2, we observed that if we have some known relevant and nonrelevant documents, then we can straightforwardly start to estimate the probability of a term  $t$  appearing in a relevant document  $P(t|R = 1)$ , and that this could be the basis of a classifier that decides whether documents are relevant or not. In this chapter, we more systematically introduce this probabilistic approach to information retrieval (IR), which provides a different formal basis for a retrieval model and results in different techniques for setting term weights.

Users start with *information needs*, which they translate into *query representations*. Similarly, there are *documents*, which are converted into *document representations* (the latter differing at least by how text is tokenized, but perhaps containing fundamentally less information, as when a nonpositional index is used). Based on these two representations, a system tries to determine how well documents satisfy information needs. In the Boolean or vector space models of IR, matching is done in a formally defined but semantically imprecise calculus of index terms. Given only a query, an IR system has an uncertain understanding of the information need. Given the query and document representations, a system has an uncertain guess of whether a document has content relevant to the information need. Probability theory provides a principled foundation for such reasoning under uncertainty. This chapter provides one answer as to how to exploit this foundation to estimate how likely it is that a document is relevant to an information need.

There is more than one possible retrieval model with a probabilistic basis. Here, we will introduce probability theory and the probability ranking principle (Sections 11.1–11.2), and then concentrate on the *binary independence model* (Section 11.3), which is the original and still most influential probabilistic retrieval model. Finally, we will introduce related but extended methods that use term counts, including the empirically successful Okapi BM25 weighting scheme, and Bayesian network models for IR (Section 11.4). In Chapter 12, we then present the alternative probabilistic language modeling approach to IR, which has been developed with considerable success in recent years.

## 11.1 Review of basic probability theory

RANDOM  
VARIABLE

We hope that the reader has seen a little basic probability theory previously. We will give a very quick review; some references for further reading appear at the end of the chapter. A variable  $A$  represents an event (a subset of the space of possible outcomes). Equivalently, we can represent the subset via a *random variable*, which is a function from outcomes to real numbers; the subset is the domain over which the random variable  $A$  has a particular value. Often we will not know with certainty whether an event is true in the world. We can ask the probability of the event  $0 \leq P(A) \leq 1$ . For two events  $A$  and  $B$ , the joint event of both events occurring is described by the joint probability  $P(A, B)$ . The conditional probability  $P(A|B)$  expresses the probability of event  $A$  given that event  $B$  occurred. The fundamental relationship between

CHAIN RULE joint and conditional probabilities is given by the *chain rule*:

$$(11.1) \quad P(A, B) = P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

Without making any assumptions, the probability of a joint event equals the probability of one of the events multiplied by the probability of the other event conditioned on knowing the first event happened.

Writing  $P(\bar{A})$  for the complement of an event, we similarly have:

$$(11.2) \quad P(\bar{A}, B) = P(B|\bar{A})P(\bar{A}).$$

PARTITION  
RULE Probability theory also has a *partition rule*, which says that if an event  $B$  can be divided into an exhaustive set of disjoint subcases, then the probability of  $B$  is the sum of the probabilities of the subcases. A special case of this rule gives that:

$$(11.3) \quad P(B) = P(A, B) + P(\bar{A}, B).$$

BAYES' RULE From these we can derive *Bayes' rule* for inverting conditional probabilities:

$$(11.4) \quad P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \left[ \frac{P(B|A)}{\sum_{X \in \{A, \bar{A}\}} P(B|X)P(X)} \right] P(A).$$

This equation can also be thought of as a way of updating probabilities. We start off with an initial estimate of how likely the event  $A$  is when we do

PRIOR  
PROBABILITY not have any other information; this is the *prior probability*  $P(A)$ . Bayes' rule lets us derive a *posterior probability*  $P(A|B)$  after having seen the evidence  $B$ , based on the *likelihood* of  $B$  occurring in the two cases that  $A$  does or does not hold.<sup>1</sup>

<sup>1</sup> The term *likelihood* is just a synonym for *probability*. It is the probability of an event or data according to a model. The term is usually used when people are thinking of holding the data fixed, while varying the model.

## 11.2 The probability ranking principle

203

ODDS Finally, it is often useful to talk about the *odds* of an event, which provide a kind of multiplier for how probabilities change:

$$(11.5) \quad \text{Odds:} \quad O(A) = \frac{P(A)}{P(\bar{A})} = \frac{P(A)}{1 - P(A)}.$$

## 11.2 The probability ranking principle

## 11.2.1 The 1/0 loss case

We assume a ranked retrieval setup as in Section 6.3, where there is a collection of documents, the user issues a query, and an ordered list of documents is returned. We also assume a binary notion of relevance as in Chapter 8. For a query  $q$  and a document  $d$  in the collection, let  $R_{d,q}$  be an indicator random variable that says whether  $d$  is relevant with respect to a given query  $q$ . That is, it takes on a value of 1 when the document is relevant and 0 otherwise. In context we will often write just  $R$  for  $R_{d,q}$ .

Using a probabilistic model, the obvious order in which to present documents to the user is to rank documents by their estimated probability of relevance with respect to the information need:  $P(R = 1|d, q)$ . This is the basis of the *probability ranking principle* (PRP) (van Rijsbergen 1979, 113–114):

PROBABILITY  
RANKING  
PRINCIPLE

If a reference retrieval system's response to each request is a ranking of the documents in the collection in order of decreasing probability of relevance to the user who submitted the request, where the probabilities are estimated as accurately as possible on the basis of whatever data have been made available to the system for this purpose, the overall effectiveness of the system to its user will be the best that is obtainable on the basis of those data.

In the simplest case of the PRP, there are no retrieval costs or other utility concerns that would differentially weight actions or errors. You lose a point for either returning a nonrelevant document or failing to return a relevant document (such a binary situation where you are evaluated on your *accuracy* is called *1/0 loss*). The goal is to return the best possible results as the top  $k$  documents, for any value of  $k$  the user chooses to examine. The PRP then says to simply rank all documents in decreasing order of  $P(R = 1|d, q)$ . If a

BAYES  
OPTIMAL  
DECISION RULE set of retrieval results is to be returned, rather than an ordering, the *Bayes optimal decision rule*, the decision that minimizes the risk of loss, is to simply return documents that are more likely relevant than nonrelevant:

$$(11.6) \quad d \text{ is relevant iff } P(R = 1|d, q) > P(R = 0|d, q).$$

**Theorem 11.1.** *The PRP is optimal, in the sense that it minimizes the expected loss (also known as the Bayes risk) under 1/0 loss.*

BAYES RISK

The proof can be found in Ripley (1996). However, it requires that all probabilities be known correctly. This is never the case in practice. Nevertheless, the PRP still provides a very useful foundation for developing models of IR.

### 11.2.2 The probability ranking principle with retrieval costs

Suppose, instead, that we assume a model of retrieval costs. Let  $C_1$  be the cost of retrieval of a relevant document and  $C_0$  the cost of retrieval of a non-relevant document. Then the PRP says that if for a specific document  $d$  and for all documents  $d'$  not yet retrieved

$$(11.7) \quad C_1 \cdot P(R = 1|d) + C_0 \cdot P(R = 0|d) \leq C_1 \cdot P(R = 1|d') + C_0 \cdot P(R = 0|d')$$

then  $d$  is the next document to be retrieved. Such a model gives a formal framework where we can model differential costs of false positives and false negatives and even system performance issues at the modeling stage, rather than simply at the evaluation stage, as we did in Section 8.6 (page 154). However, we will not further consider loss/utility models in this chapter.

## 11.3 The binary independence model

**BINARY INDEPENDENCE MODEL** The *binary independence model* (BIM) we present in this section is the model that has traditionally been used with the PRP. It introduces some simple assumptions, which make estimating the probability function  $P(R|d, q)$  practical. Here, “binary” is equivalent to Boolean: Documents and queries are both represented as binary term incidence vectors. That is, a document  $d$  is represented by the vector  $\vec{x} = (x_1, \dots, x_M)$  where  $x_t = 1$  if term  $t$  is present in document  $d$  and  $x_t = 0$  if  $t$  is not present in  $d$ . With this representation, many possible documents have the same vector representation. Similarly, we represent  $q$  by the incidence vector  $\vec{q}$  (the distinction between  $q$  and  $\vec{q}$  is less central because commonly  $q$  is in the form of a set of words). “Independence” means that terms are modeled as occurring in documents independently. The model recognizes no association between terms. This assumption is far from correct, but it nevertheless often gives satisfactory results in practice; it is the “naive” assumption of Naive Bayes models, discussed further in Section 13.4 (page 245). Indeed, the BIM is exactly the same as the multivariate Bernoulli Naive Bayes model presented in Section 13.3 (page 243). In a sense, this assumption is equivalent to an assumption of the vector space model, where each term is a dimension that is orthogonal to all other terms.

We will first present a model that assumes that the user has a single-step information need. As discussed in Chapter 9, seeing a range of results might let the user refine their information need. Fortunately, as mentioned there,

## 11.3 The binary independence model

205

it is straightforward to extend the BIM so as to provide a framework for relevance feedback, and we present this model in Section 11.3.4.

To make a probabilistic retrieval strategy precise, we need to estimate how terms in documents contribute to relevance; specifically, we wish to know how term frequency, document frequency, document length, and other statistics that we can compute influence judgments about document relevance, and how they can be reasonably combined to estimate the probability of document relevance. We then order documents by decreasing estimated probability of relevance.

We assume here that the relevance of each document is independent of the relevance of other documents. As we noted in Section 8.5.1 (page 153), this is incorrect; the assumption is especially harmful in practice if it allows a system to return duplicate or near duplicate documents. Under the BIM, we model the probability  $P(R|d, q)$  that a document is relevant via the probability in terms of term incidence vectors  $P(R|\vec{x}, \vec{q})$ . Then, using Bayes rule, we have:

$$(11.8) \quad \begin{aligned} P(R = 1|\vec{x}, \vec{q}) &= \frac{P(\vec{x}|R = 1, \vec{q})P(R = 1|\vec{q})}{P(\vec{x}|\vec{q})} \\ P(R = 0|\vec{x}, \vec{q}) &= \frac{P(\vec{x}|R = 0, \vec{q})P(R = 0|\vec{q})}{P(\vec{x}|\vec{q})}. \end{aligned}$$

Here,  $P(\vec{x}|R = 1, \vec{q})$  and  $P(\vec{x}|R = 0, \vec{q})$  are the probability that if a relevant or nonrelevant, respectively, document is retrieved, then that document's representation is  $\vec{x}$ . You should think of this quantity as defined with respect to a space of possible documents in a domain. How do we compute all these probabilities? We never know the exact probabilities, and so we have to use estimates: Statistics about the actual document collection are used to estimate these probabilities.  $P(R = 1|\vec{q})$  and  $P(R = 0|\vec{q})$  indicate the prior probability of retrieving a relevant or nonrelevant document, respectively, for a query  $\vec{q}$ . Again, if we knew the percentage of relevant documents in the collection, then we could use this number to estimate  $P(R = 1|\vec{q})$  and  $P(R = 0|\vec{q})$ . Because a document is either relevant or nonrelevant to a query, we must have that:

$$(11.9) \quad P(R = 1|\vec{x}, \vec{q}) + P(R = 0|\vec{x}, \vec{q}) = 1.$$

## 11.3.1 Deriving a ranking function for query terms

Given a query  $q$ , we wish to order returned documents by descending  $P(R = 1|d, q)$ . Under the BIM, this is modeled as ordering by  $P(R = 1|\vec{x}, \vec{q})$ . Rather than estimating this probability directly, because we are interested only in the ranking of documents, we work with some other quantities which are easier to compute and which give the same ordering of documents. In

particular, we can rank documents by their odds of relevance (because the odds of relevance is monotonic with the probability of relevance). This makes things easier, because we can ignore the common denominator in (11.8), giving:

$$(11.10) \quad O(R|\vec{x}, \vec{q}) = \frac{P(R=1|\vec{x}, \vec{q})}{P(R=0|\vec{x}, \vec{q})} = \frac{\frac{P(R=1|\vec{q})P(\vec{x}|R=1, \vec{q})}{P(\vec{x}|\vec{q})}}{\frac{P(R=0|\vec{q})P(\vec{x}|R=0, \vec{q})}{P(\vec{x}|\vec{q})}} = \frac{P(R=1|\vec{q})}{P(R=0|\vec{q})} \cdot \frac{P(\vec{x}|R=1, \vec{q})}{P(\vec{x}|R=0, \vec{q})}.$$

The left term in the rightmost expression of Equation (11.10) is a constant for a given query. Because we are only ranking documents, there is thus no need for us to estimate it. The right-hand term does, however, require estimation, and this initially appears to be difficult: How can we accurately estimate the probability of an entire term incidence vector occurring? It is at this point that we make the *Naive Bayes conditional independence assumption* that the presence or absence of a word in a document is independent of the presence or absence of any other word (given the query):

NAIVE BAYES  
ASSUMPTION

$$(11.11) \quad \frac{P(\vec{x}|R=1, \vec{q})}{P(\vec{x}|R=0, \vec{q})} = \prod_{t=1}^M \frac{P(x_t|R=1, \vec{q})}{P(x_t|R=0, \vec{q})}.$$

So:

$$(11.12) \quad O(R|\vec{x}, \vec{q}) = O(R|\vec{q}) \cdot \prod_{t=1}^M \frac{P(x_t|R=1, \vec{q})}{P(x_t|R=0, \vec{q})}.$$

Because each  $x_t$  is either 0 or 1, we can separate the terms to give:

$$(11.13) \quad O(R|\vec{x}, \vec{q}) = O(R|\vec{q}) \cdot \prod_{t:x_t=1} \frac{P(x_t=1|R=1, \vec{q})}{P(x_t=1|R=0, \vec{q})} \cdot \prod_{t:x_t=0} \frac{P(x_t=0|R=1, \vec{q})}{P(x_t=0|R=0, \vec{q})}.$$

Henceforth, let  $p_t = P(x_t=1|R=1, \vec{q})$  be the probability of a term appearing in a document relevant to the query, and  $u_t = P(x_t=1|R=0, \vec{q})$  be the probability of a term appearing in a nonrelevant document. These quantities can be visualized in the following contingency table where the columns add to 1:

(11.14)

	document	relevant ( $R=1$ )	nonrelevant ( $R=0$ )
term present	$x_t = 1$	$p_t$	$u_t$
term absent	$x_t = 0$	$1 - p_t$	$1 - u_t$

Let us make an additional simplifying assumption that terms not occurring in the query are equally likely to occur in relevant and nonrelevant documents; that is, if  $q_t = 0$  then  $p_t = u_t$ . (This assumption can be changed, as when doing relevance feedback in Section 11.3.4.) Then we need only consider terms in the products that appear in the query, and so,

$$(11.15) \quad O(R|\vec{q}, \vec{x}) = O(R|\vec{q}) \cdot \prod_{t:x_t=q_t=1} \frac{p_t}{u_t} \cdot \prod_{t:x_t=0, q_t=1} \frac{1-p_t}{1-u_t}.$$

## 11.3 The binary independence model

207

The left product is over query terms found in the document and the right product is over query terms not found in the document.

We can manipulate this expression by including the query terms found in the document into the right product, but simultaneously dividing through by them in the left product, so the value is unchanged. Then we have:

$$(11.16) \quad O(R|\vec{q}, \vec{x}) = O(R|\vec{q}) \cdot \prod_{t:x_t=q_t=1} \frac{p_t(1-u_t)}{u_t(1-p_t)} \cdot \prod_{t:q_t=1} \frac{1-p_t}{1-u_t}.$$

The left product is still over query terms found in the document, but the right product is now over all query terms. That means that this right product is a constant for a particular query, just like the odds  $O(R|\vec{q})$ . So the only quantity that needs to be estimated to rank documents for relevance to a query is the left product. We can equally rank documents by the logarithm of this term, since log is a monotonic function. The resulting quantity used for ranking is called the *retrieval status value* (RSV) in this model:

RETRIEVAL  
STATUS VALUE

$$(11.17) \quad RSV_d = \log \prod_{t:x_t=q_t=1} \frac{p_t(1-u_t)}{u_t(1-p_t)} = \sum_{t:x_t=q_t=1} \log \frac{p_t(1-u_t)}{u_t(1-p_t)}.$$

So everything comes down to computing the RSV. Define  $c_t$ :

$$(11.18) \quad c_t = \log \frac{p_t(1-u_t)}{u_t(1-p_t)} = \log \frac{p_t}{(1-p_t)} + \log \frac{1-u_t}{u_t}.$$

The  $c_t$  terms are log odds ratios for the terms in the query. We have the odds of the term appearing if the document is relevant ( $p_t/(1-p_t)$ ) and the odds of the term appearing if the document is nonrelevant ( $u_t/(1-u_t)$ ). The *odds ratio* is the ratio of two such odds, and then we finally take the log of that quantity. The value will be 0 if a term has equal odds of appearing in relevant and nonrelevant documents, and positive if it is more likely to appear in relevant documents. The  $c_t$  quantities function as term weights in the model, and the document score for a query is  $RSV_d = \sum_{t:x_t=q_t=1} c_t$ . Operationally, we sum them in accumulators for query terms appearing in documents, just as for the vector space model calculations discussed in Section 7.1 (page 124). We now turn to how we estimate these  $c_t$  quantities for a particular collection and query.

ODDS RATIO

## 11.3.2 Probability estimates in theory

For each term  $t$ , what would these  $c_t$  numbers look like for the whole collection? Equation (11.19) gives a contingency table of counts of documents in the collection, where  $df_t$  is the number of documents that contain term  $t$ :

(11.19)

	documents	relevant	nonrelevant	total
term present	$x_t = 1$	$s$	$df_t - s$	$df_t$
term absent	$x_t = 0$	$S - s$	$(N - df_t) - (S - s)$	$N - df_t$
	total	$S$	$N - S$	$N$

Using this,  $p_t = s/S$  and  $u_t = (df_t - s)/(N - S)$  and

$$(11.20) \quad c_t = K(N, df_t, S, s) = \log \frac{s/(S - s)}{(df_t - s)/((N - df_t) - (S - s))}.$$

To avoid the possibility of zeroes (such as if every or no relevant document has a particular term) it is fairly standard to add  $\frac{1}{2}$  to each of the quantities in the center four terms of (11.19), and then to adjust the marginal counts (the totals) accordingly (so, the bottom right cell totals  $N + 2$ ). Then we have:

$$(11.21) \quad \hat{c}_t = K(N, df_t, S, s) = \log \frac{(s + \frac{1}{2})/(S - s + \frac{1}{2})}{(df_t - s + \frac{1}{2})/(N - df_t - S + s + \frac{1}{2})}.$$

Adding  $\frac{1}{2}$  in this way is a simple form of smoothing. For trials with categorical outcomes (such as noting the presence or absence of a term), one way to estimate the probability of an event from data is simply to count the number of times an event occurred divided by the total number of trials. This is referred to as the *relative frequency* of the event. Estimating the probability as the relative frequency is the *maximum likelihood estimate* (or MLE), because this value makes the observed data maximally likely. However, if we simply use the MLE, then the probability given to events we happened to see is usually too high, whereas other events may be completely unseen and giving them as a probability estimate their relative frequency of 0 is both an underestimate and normally breaks our models; anything multiplied by 0 is 0. Simultaneously decreasing the estimated probability of seen events and increasing the probability of unseen events is referred to as *smoothing*. One simple way of smoothing is to add a number  $\alpha$  to each of the observed counts. These *pseudocounts* correspond to the use of a uniform distribution over the vocabulary as a *Bayesian prior*, following Equation (11.4). We initially assume a uniform distribution over events, where the size of  $\alpha$  denotes the strength of our belief in uniformity, and we then update the probability based on observed events. Because our belief in uniformity is weak, we use  $\alpha = \frac{1}{2}$ . This is a form of *maximum a posteriori* (MAP) estimation, where we choose the most likely point value for probabilities based on the prior and the observed evidence, following Equation (11.4). We will further discuss methods of smoothing estimated counts to give probability models in Section 12.2.2 (page 224); the simple method of adding  $\frac{1}{2}$  to each observed count will do for now.



### 11.3.3 Probability estimates in practice

Under the assumption that relevant documents are a very small percentage of the collection, it is plausible to approximate statistics for nonrelevant documents by statistics from the whole collection. Under this assumption,  $u_t$  (the probability of term occurrence in nonrelevant documents for a query) is  $df_t/N$  and

$$(11.22) \quad \log[(1 - u_t)/u_t] = \log[(N - df_t)/df_t] \approx \log N/df_t$$

In other words, we can provide a theoretical justification for the most frequently used form of idf weighting, which we saw in Section 6.2.1.

The approximation technique in Equation (11.22) cannot easily be extended to relevant documents. The quantity  $p_t$  can be estimated in various ways:

1. We can use the frequency of term occurrence in known relevant documents (if we know some). This is the basis of probabilistic approaches to relevance feedback weighting in a feedback loop, discussed in the next subsection.
2. Croft and Harper (1979) proposed using a constant in their combination match model. For instance, we might assume that  $p_t$  is constant over all terms  $x_t$  in the query and that  $p_t = 0.5$ . This means that each term has even odds of appearing in a relevant document, and so the  $p_t$  and  $(1 - p_t)$  factors cancel out in the expression for  $RSV$ . Such an estimate is weak, but doesn't disagree violently with our hopes for the search terms appearing in many but not all relevant documents. Combining this method with our earlier approximation for  $u_t$ , the document ranking is determined simply by which query terms occur in documents scaled by their idf weighting. For short documents (titles or abstracts) in situations in which iterative searching is undesirable, using this weighting term alone can be quite satisfactory, although in many other circumstances we would like to do better.
3. Greiff (1998) argues that the constant estimate of  $p_t$  in the Croft and Harper (1979) model is theoretically problematic and not observed empirically, and argues that a much better estimate is found by simply estimating  $p_t$  from collection level statistics about the occurrence of  $t$ , as  $p_t = df_t/N$ .

Iterative methods of estimation, which combine some of the above ideas, are discussed in the next subsection.

### 11.3.4 Probabilistic approaches to relevance feedback

We can use (pseudo) relevance feedback (RF), perhaps in an iterative process of estimation, to get a more accurate estimate of  $p_t$ . The probabilistic approach to RF works as follows.

1. Guess initial estimates of  $p_t$  and  $u_t$ . This can be done using the probability estimates of the previous section. For instance, we can assume that  $p_t$  is constant over all  $x_i$  in the query, in particular, perhaps taking  $p_t = \frac{1}{2}$ .
2. Use the current estimates of  $p_t$  and  $u_t$  to determine a best guess at the set of relevant documents  $R = \{d : R_{d,q} = 1\}$ . Use this model to retrieve a set of candidate relevant documents, which we present to the user.
3. We interact with the user to refine the model of  $R$ . We do this by learning from the user relevance judgments for some subset of documents  $V$ . Based on relevance judgments,  $V$  is partitioned into two subsets:  $VR = \{d \in V, R_{d,q} = 1\} \subset R$  and  $VNR = \{d \in V, R_{d,q} = 0\}$ , which is disjoint from  $R$ .
4. We reestimate  $p_t$  and  $u_t$  on the basis of known relevant and nonrelevant documents. If the sets  $VR$  and  $VNR$  are large enough, we may be able to estimate these quantities directly from these documents as maximum likelihood estimates:

$$(11.23) \quad p_t = |VR_t|/|VR|$$

where  $VR_t$  is the set of documents in  $VR$  containing  $x_t$ . In practice, we usually need to smooth these estimates. We can do this by adding  $\frac{1}{2}$  to both the count  $|VR_t|$  and to the number of relevant documents not containing the term, giving:

$$(11.24) \quad p_t = \frac{|VR_t| + \frac{1}{2}}{|VR| + 1}.$$

However, the set of documents judged by the user ( $V$ ) is usually very small, and so the resulting statistical estimate is quite unreliable (noisy), even if the estimate is smoothed. So it is often better to combine the new information with the original guess in a process of Bayesian updating. In this case we have:

$$(11.25) \quad p_t^{(k+1)} = \frac{|VR_t| + \kappa p_t^{(k)}}{|VR| + \kappa}.$$

Here  $p_t^{(k)}$  is the  $k^{\text{th}}$  estimate for  $p_t$  in an iterative updating process and is used as a Bayesian prior in the next iteration with a weighting of  $\kappa$ . Relating this equation back to Equation (11.4) requires a bit more probability theory than we have presented here (we need to use a beta distribution prior, conjugate to the Bernoulli random variable  $X_t$ ). But the form of the resulting equation is quite straightforward: Rather than uniformly distributing pseudocounts, we now distribute a total of  $\kappa$  pseudocounts according to the previous estimate, which acts as the prior distribution. In the absence of other evidence (and assuming that the user is perhaps indicating roughly five relevant or nonrelevant documents) then a value of around  $\kappa = 5$  is perhaps appropriate. That is, the prior is strongly weighted so that the estimate does not change too much from the evidence provided by a very small number of documents.

## 11.3 The binary independence model

211

5. Repeat the above process from Step 2, generating a succession of approximations to  $R$  and hence  $p_t$ , until the user is satisfied.

It is also straightforward to derive a pseudo RF version of this algorithm, where we simply pretend that  $VR = V$ . More briefly:

1. Assume initial estimates for  $p_t$  and  $u_t$  as above.
2. Determine a guess for the size of the relevant document set. If unsure, a conservative (too small) guess is likely to be best. This motivates use of a fixed size set  $V$  of highest ranked documents.
3. Improve our guesses for  $p_t$  and  $u_t$ . We choose from the methods of Equations (11.23) and (11.25) for reestimating  $p_t$ , except now based on the set  $V$  instead of  $VR$ . If we let  $V_t$  be the subset of documents in  $V$  containing  $x_t$  and use add  $\frac{1}{2}$  smoothing, we get:

$$(11.26) \quad p_t = \frac{|V_t| + \frac{1}{2}}{|V| + 1}$$

and if we assume that documents that are non retrieved are nonrelevant then we can update our  $u_t$  estimates as:

$$(11.27) \quad u_t = \frac{df_t - |V_t| + \frac{1}{2}}{N - |V| + 1}.$$

4. Go to Step 2 until the ranking of the returned results converges.

Once we have a real estimate for  $p_t$ , then the  $c_t$  weights used in the  $RSV$  value look almost like a tf-idf value. For instance, using Equation (11.18), Equation (11.22), and Equation (11.26), we have:

$$(11.28) \quad c_t = \log \left[ \frac{p_t}{1 - p_t} \cdot \frac{1 - u_t}{u_t} \right] \approx \log \left[ \frac{|V_t| + \frac{1}{2}}{|V| - |V_t| + 1} \cdot \frac{N}{df_t} \right].$$

But things aren't quite the same:  $p_t/(1 - p_t)$  measures the (estimated) proportion of relevant documents that the term  $t$  occurs in, not term frequency. Moreover, if we apply log identities:

$$(11.29) \quad c_t = \log \frac{|V_t| + \frac{1}{2}}{|V| - |V_t| + 1} + \log \frac{N}{df_t}$$

we see that we are now *adding* the two log-scaled components rather than multiplying them.

**? Exercise 11.1** Work through the derivation of Equation (11.20) from Equations (11.18) and (11.19).

**Exercise 11.2** What are the differences between standard vector space tf-idf weighting and the BIM probabilistic retrieval model (in the case where no document relevance information is available)?

**Exercise 11.3 [★★]** Let  $X_t$  be a random variable indicating whether the term  $t$  appears in a document. Suppose we have  $|R|$  relevant documents in

the document collection and that  $X_t = 1$  in  $s$  of the documents. Take the observed data to be just these observations of  $X_t$  for each document in  $R$ . Show that the MLE for the parameter  $p_t = P(X_t = 1 | R = 1, \vec{q})$ , that is, the value for  $p_t$  which maximizes the probability of the observed data, is  $p_t = s/|R|$ .

**Exercise 11.4** Describe the differences between vector space relevance feedback and probabilistic relevance feedback.

## 11.4 An appraisal and some extensions

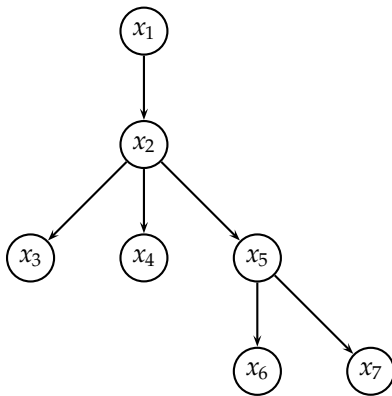
### 11.4.1 An appraisal of probabilistic models

Probabilistic methods are one of the oldest formal models in IR. Already in the 1970s they were held out as an opportunity to place IR on a firmer theoretical footing, and with the resurgence of probabilistic methods in computational linguistics in the 1990s, that hope has returned, and probabilistic methods are again one of the currently hottest topics in IR. Traditionally, probabilistic IR has had neat ideas but the methods have never won on performance. Getting reasonable approximations of the needed probabilities for a probabilistic IR model is possible, but it requires some major assumptions. In the BIM these are:

- a Boolean representation of documents/queries/relevance
- term independence
- terms not in the query don't affect the outcome
- document relevance values are independent

It is perhaps the severity of the modeling assumptions that makes achieving good performance difficult. A general problem seems to be that probabilistic models either require partial relevance information or else only allow for deriving apparently inferior term weighting models.

Things started to change in the 1990s when the BM25 weighting scheme, which we discuss later in this section, showed very good performance, and started to be adopted as a term weighting scheme by many groups. The difference between "vector space" and "probabilistic" IR systems is not that great; in either case, you build an information retrieval scheme in the exact same way that we discussed in Chapter 7. For a probabilistic IR system, it's just that, at the end, you score queries not by cosine similarity and tf-idf in a vector space, but by a slightly different formula motivated by probability theory. Indeed, sometimes people have changed an existing vector-space IR system into an effectively probabilistic system simply by adopted term weighting formulas from probabilistic models. In this section, we briefly present three extensions of the traditional probabilistic model, and in the next



**Figure 11.1** A tree of dependencies between terms. In this graphical model representation, a term  $x_i$  is directly dependent on a term  $x_k$  if there is an arrow  $x_k \rightarrow x_i$ .

chapter, we look at the somewhat different probabilistic language modeling approach to IR.

#### 11.4.2 Tree-structured dependencies between terms

Some of the assumptions of the BIM can be removed. For example, we can remove the assumption that terms are independent. This assumption is very far from true in practice. A case that particularly violates this assumption is term pairs like Hong and Kong, which are strongly dependent. But dependencies can occur in various complex configurations, such as between the set of terms New, York, England, City, Stock, Exchange, and University. van Rijsbergen (1979) proposed a simple, plausible model that allowed a tree structure of term dependencies, as in Figure 11.1. In this model, each term can be directly dependent on only one other term, giving a tree structure of dependencies. When it was invented in the 1970s, estimation problems held back the practical success of this model, but the idea was reinvented as the tree-augmented Naive Bayes model by Friedman and Goldszmidt (1996), who used it with some success on various machine learning data sets.

#### 11.4.3 Okapi BM25: A nonbinary model

BM25 WEIGHTS  
OKAPI  
WEIGHTING

The BIM was originally designed for short catalog records and abstracts of fairly consistent length, and it works reasonably in these contexts, but for modern full-text search collections, it seems clear that a model should pay attention to term frequency and document length, as in Chapter 6. The *BM25 weighting scheme*, often called *Okapi weighting*, after the system in which it was first implemented, was developed as a way of building a probabilistic model sensitive to these quantities while not introducing too many additional

parameters into the model (Spärck Jones et al. 2000). We will not develop the full theory behind the model here, but just present a series of forms that build up to the standard form now used for document scoring. The simplest score for document  $d$  is just idf weighting of the query terms present, as in Equation (11.22):

$$(11.30) \quad RSV_d = \sum_{t \in q} \log \frac{N}{df_t}.$$

Sometimes, an alternative version of idf is used. If we start with the formula in Equation (11.21) but in the absence of relevance feedback information we estimate that  $S = s = 0$ , then we get an alternative idf formulation as follows:

$$(11.31) \quad RSV_d = \sum_{t \in q} \log \frac{N - df_t + \frac{1}{2}}{df_t + \frac{1}{2}}.$$

This variant behaves slightly strangely: If a term occurs in over half the documents in the collection, then this model gives a negative term weight, which is presumably undesirable. But, assuming the use of a stop list, this normally doesn't happen, and the value for each summand can be given a floor of 0.

We can improve on Equation (11.30) by factoring in the frequency of each term and document length:

$$(11.32) \quad RSV_d = \sum_{t \in q} \log \left[ \frac{N}{df_t} \right] \cdot \frac{(k_1 + 1)tf_{td}}{k_1((1 - b) + b \times (L_d/L_{ave})) + tf_{td}}.$$

Here,  $tf_{td}$  is the frequency of term  $t$  in document  $d$ , and  $L_d$  and  $L_{ave}$  are the length of document  $d$  and the average document length for the whole collection. The variable  $k_1$  is a positive tuning parameter that calibrates the document term frequency scaling. A  $k_1$  value of 0 corresponds to a binary model (no term frequency), and a large value corresponds to using raw term frequency.  $b$  is another tuning parameter ( $0 \leq b \leq 1$ ) that determines the scaling by document length:  $b = 1$  corresponds to fully scaling the term weight by the document length, whereas  $b = 0$  corresponds to no length normalization.

If the query is long, then we might also use similar weighting for query terms. This is appropriate if the queries are paragraph-long information needs, but unnecessary for short queries.

$$(11.33) \quad RSV_d = \sum_{t \in q} \left[ \log \frac{N}{df_t} \right] \cdot \frac{(k_1 + 1)tf_{td}}{k_1((1 - b) + b \times (L_d/L_{ave})) + tf_{td}} \cdot \frac{(k_3 + 1)tf_{tq}}{k_3 + tf_{tq}}$$

with  $tf_{tq}$  being the frequency of term  $t$  in the query  $q$ , and  $k_3$  being another positive tuning parameter that this time calibrates term frequency scaling of the query. In the equation presented, there is no length normalization of queries (it is as if  $b = 0$  here). Length normalization of the query is unnecessary because retrieval is being done with respect to a single fixed query. The tuning parameters of these formulas should ideally be set to optimize performance on a development test collection (see page 141). That is, we

## 11.4 An appraisal and some extensions

215

can search for values of these parameters that maximize performance on a separate development test collection (either manually or with optimization methods, such as grid search or something more advanced), and then use these parameters on the actual test collection. In the absence of such optimization, experiments have shown reasonable values are to set  $k_1$  and  $k_3$  to a value between 1.2 and 2 and  $b = 0.75$ .

If we have relevance judgments available, then we can use the full form of (11.21) in place of the approximation  $\log(N/df_t)$  introduced in (11.22):

$$(11.34) \quad RSV_d = \sum_{t \in q} \left[ \log \left[ \frac{(|VR_t| + \frac{1}{2})/(|VNR_t| + \frac{1}{2})}{(df_t - |VR_t| + \frac{1}{2})/(N - df_t - |VR| + |VR_t| + \frac{1}{2})} \right] \right. \\ \left. \times \frac{(k_1 + 1)tf_{td}}{k_1((1 - b) + b(L_d/L_{ave})) + tf_{td}} \times \frac{(k_3 + 1)tf_{tq}}{k_3 + tf_{tq}} \right].$$

Here,  $VR_t$ ,  $NVR_t$ , and  $VR$  are used as in Section 11.3.4. The first part of the expression reflects relevance feedback (or just idf weighting if no relevance information is available), the second implements document term frequency and document length scaling, and the third considers term frequency in the query.

Rather than just providing a term weighting method for terms in a user's query, relevance feedback can also involve augmenting the query (automatically or with manual review) with some (say, ten to twenty) of the top terms in the known-relevant documents as ordered by the relevance factor  $\hat{c}_t$  from Equation (11.21), and the above formula can then be used with such an augmented query  $q$ .

The BM25 term weighting formulas have been used quite widely and quite successfully across a range of collections and search tasks. Especially in the TREC evaluations, they performed well and were widely adopted by many groups. See Spärck Jones et al. (2000) for extensive motivation and discussion of experimental results.

## 11.4.4 Bayesian network approaches to information retrieval

BAYESIAN NETWORKS Turtle and Croft (1989, 1991) introduced into information retrieval the use of *Bayesian networks* (Jensen and Jensen 2001), a form of probabilistic graphical model. We skip the details because fully introducing the formalism of Bayesian networks would require much too much space, but conceptually, Bayesian networks use directed graphs to show probabilistic dependencies between variables, as in Figure 11.1, and have led to the development of sophisticated algorithms for propagating influence so as to allow learning and inference with arbitrary knowledge within arbitrary directed acyclic graphs. Turtle and Croft used a sophisticated network to better model the complex dependencies between a document and a user's information need.

The model decomposes into two parts: a document collection network and a query network. The document collection network is large, but can be pre-computed; it maps from documents to terms to concepts. The concepts are a thesaurus-based expansion of the terms appearing in the document. The query network is relatively small but a new network needs to be built each time a query comes in, and then attached to the document network. The query network maps from query terms, to query subexpressions (built using probabilistic or “noisy” versions of AND and OR operators), to the user’s information need.

The result is a flexible probabilistic network that can generalize various simpler Boolean and probabilistic models. Indeed, this is the primary case of a statistical ranked retrieval model that naturally supports structured query operators. The system allowed efficient large-scale retrieval, and was the basis of the InQuery text retrieval system, built at the University of Massachusetts. This system performed very well in TREC evaluations and for a time was sold commercially. On the other hand, the model still used various approximations and independence assumptions to make parameter estimation and computation possible. There has not been much follow-on work along these lines, but we would note that this model was actually built very early on in the modern era of using Bayesian networks, and there have been many subsequent developments in the theory, and the time is perhaps right for a new generation of Bayesian network-based IR systems.

## 11.5 References and further reading

Longer introductions to probability theory can be found in most introductory probability and statistics books, such as (Grinstead and Snell 1997; Rice 2006; Ross 2006). An introduction to Bayesian utility theory can be found in (Ripley 1996).

The probabilistic approach to IR originated in the United Kingdom in the 1950s. The first major presentation of a probabilistic model is Maron and Kuhns (1960). Robertson and Jones (1976) introduce the main foundations of the BIM and van Rijsbergen (1979) presents in detail the classic BIM probabilistic model. The idea of the PRP is variously attributed to S. E. Robertson, M. E. Maron, and W. S. Cooper (the term *probabilistic ordering principle* is used in Robertson and Jones (1976), but PRP dominates in later work). Fuhr (1992) is a more recent presentation of probabilistic IR, which includes coverage of other approaches such as probabilistic logics and Bayesian networks. Crestani et al. (1998) is another survey. Spärck Jones et al. (2000) is the definitive presentation of probabilistic IR experiments by the “London school,” and Robertson (2005) presents a retrospective on the group’s participation in TREC evaluations, including detailed discussion of the Okapi



*11.5 References and further reading*

217

BM25 scoring function and its development. Robertson et al. (2004) extend BM25 to the case of multiple weighted fields.

The open-source Indri search engine, which is distributed with the Lemur toolkit ([www.lemurproject.org/](http://www.lemurproject.org/)) merges ideas from Bayesian inference networks and statistical language modeling approaches (see Chapter 12), in particular preserving the former's support for structured query operators.