

Modern Information  
Retrieval  
Baeza-Yates and  
Eibeiro-Nelo

or by a variation of this formula. Such term-weighting strategies are called *tf-idf* schemes.

Several variations of the above expression for the weight  $w_{i,q}$  are described in an interesting paper by Salton and Buckley which appeared in 1988 [696]. However, in general, the above expression should provide a good weighting scheme for many collections.

For the query term weights, Salton and Buckley suggest

$$w_{i,q} = \left( 0.5 + \frac{0.5 \text{freq}_{i,q}}{\max_i \text{freq}_{i,q}} \right) \times \log \frac{N}{n_i} \quad (2.4)$$

where  $\text{freq}_{i,q}$  is the raw frequency of the term  $k_i$  in the text of the information request  $q$ .

The main *advantages* of the vector model are: (1) its term-weighting scheme improves retrieval performance; (2) its partial matching strategy allows retrieval of documents that *approximate* the query conditions; and (3) its cosine ranking formula sorts the documents according to their degree of similarity to the query. Theoretically, the vector model has the *disadvantage* that index terms are assumed to be mutually independent (equation 2.3 does *not* account for index term dependencies). However, in practice, consideration of term dependencies might be a disadvantage. Due to the locality of many term dependencies, their indiscriminate application to all the documents in the collection might in fact *hurt* the overall performance.

Despite its simplicity, the vector model is a resilient ranking strategy with general collections. It yields ranked answer sets which are difficult to improve upon without query expansion or relevance feedback (see Chapter 5) within the framework of the vector model. A large variety of alternative ranking methods have been compared to the vector model but the consensus seems to be that, in general, the vector model is either superior or almost as good as the known alternatives. Furthermore, it is simple and fast. For these reasons, the vector model is a popular retrieval model nowadays.

## 2.5.4 Probabilistic Model

In this section, we describe the classic probabilistic model introduced in 1976 by Robertson and Sparck Jones [677] which later became known as the *binary independence retrieval* (BIR) model. Our discussion is intentionally brief and focuses mainly on highlighting the key features of the model. With this purpose in mind, we do not detain ourselves in subtleties regarding the binary independence assumption for the model. The section on bibliographic discussion points to references which cover these details.

The probabilistic model attempts to capture the IR problem within a probabilistic framework. The fundamental idea is as follows. Given a user query, there is a set of documents which contains exactly the relevant documents and

no other. Let us refer to this set of documents as the *ideal answer set*. Given the description of this ideal answer set, we would have no problems in retrieving its documents. Thus, we can think of the querying process as a process of specifying the properties of an ideal answer set (which is analogous to interpreting the IR problem as a problem of clustering). The problem is that we do not know exactly what these properties are. All we know is that there are index terms whose semantics should be used to characterize these properties. Since these properties are not known at query time, an effort has to be made at initially guessing what they could be. This initial guess allows us to generate a preliminary probabilistic description of the ideal answer set which is used to retrieve a first set of documents. An interaction with the user is then initiated with the purpose of improving the probabilistic description of the ideal answer set. Such interaction could proceed as follows.

The user takes a look at the retrieved documents and decides which ones are relevant and which ones are not (in truth, only the first top documents need to be examined). The system then uses this information to refine the description of the ideal answer set. By repeating this process many times, it is expected that such a description will evolve and become closer to the real description of the ideal answer set. Thus, one should always have in mind the need to guess at the beginning the description of the ideal answer set. Furthermore, a conscious effort is made to model this description in probabilistic terms.

The probabilistic model is based on the following fundamental assumption.

*Assumption (Probabilistic Principle)* Given a user query  $q$  and a document  $d_j$  in the collection, the probabilistic model tries to estimate the probability that the user will find the document  $d_j$  interesting (i.e., relevant). The model assumes that this probability of relevance depends on the query and the document representations only. Further, the model assumes that there is a subset of all documents which the user prefers as the answer set for the query  $q$ . Such an *ideal answer set* is labeled  $R$  and should maximize the overall probability of relevance to the user. Documents in the set  $R$  are predicted to be *relevant* to the query. Documents not in this set are predicted to be *non-relevant*.

This assumption is quite troublesome because it does not state explicitly how to compute the probabilities of relevance. In fact, not even the sample space which is to be used for defining such probabilities is given.

Given a query  $q$ , the probabilistic model assigns to each document  $d_j$ , as a measure of its similarity to the query, the ratio  $P(d_j \text{ relevant to } q) / P(d_j \text{ non-relevant to } q)$  which computes the odds of the document  $d_j$  being relevant to the query  $q$ . Taking the odds of relevance as the rank minimizes the probability of an erroneous judgement [282, 785].

**Definition** For the probabilistic model, the index term weight variables are all binary i.e.,  $w_{i,j} \in \{0, 1\}$ ;  $w_{i,q} \in \{0, 1\}$ . A query  $q$  is a subset of index terms. Let  $R$  be the set of documents known (or initially guessed) to be relevant. Let  $\bar{R}$  be the complement of  $R$  (i.e., the set of non-relevant documents). Let  $P(R|\bar{d}_j)$

be the probability that the document  $d_j$  is relevant to the query  $q$  and  $P(\bar{R}|d_j)$  be the probability that  $d_j$  is non-relevant to  $q$ . The similarity  $sim(d_j, q)$  of the document  $d_j$  to the query  $q$  is defined as the ratio

$$sim(d_j, q) = \frac{P(R|\bar{d}_j)}{P(\bar{R}|d_j)} \longrightarrow R(a|b) = \frac{R(b|a)R(a)}{R(b)}$$

Using Bayes' rule,

$$sim(d_j, q) = \frac{P(\bar{d}_j|R) \times P(R)}{P(\bar{d}_j|\bar{R}) \times P(\bar{R})}$$

$P(\bar{d}_j|R)$  stands for the probability of randomly selecting the document  $d_j$  from the set  $R$  of relevant documents. Further,  $P(R)$  stands for the probability that a document randomly selected from the entire collection is relevant. The meanings attached to  $P(\bar{d}_j|\bar{R})$  and  $P(\bar{R})$  are analogous and complementary.

Since  $P(R)$  and  $P(\bar{R})$  are the same for all the documents in the collection, we write,

$$sim(d_j, q) \sim \frac{P(\bar{d}_j|R)}{P(\bar{d}_j|\bar{R})}$$

Assuming independence of index terms,

$$sim(d_j, q) \sim \frac{\prod_{g_i(\bar{d}_j)=1} P(k_i|R) \times \prod_{g_i(\bar{d}_j)=0} P(\bar{k}_i|R)}{\prod_{g_i(\bar{d}_j)=1} P(k_i|\bar{R}) \times \prod_{g_i(\bar{d}_j)=0} P(\bar{k}_i|\bar{R})}$$

$P(k_i|R)$  stands for the probability that the index term  $k_i$  is present in a document randomly selected from the set  $R$ .  $P(\bar{k}_i|R)$  stands for the probability that the index term  $k_i$  is not present in a document randomly selected from the set  $R$ . The probabilities associated with the set  $\bar{R}$  have meanings which are analogous to the ones just described.

Taking logarithms, recalling that  $P(k_i|R) + P(\bar{k}_i|R) = 1$ , and ignoring factors which are constant for all documents in the context of the same query, we can finally write

$$sim(d_j, q) \sim \sum_{i=1}^L u_{i,q} \times u_{i,j} \times \left( \log \frac{P(k_i|R)}{1 - P(k_i|R)} + \log \frac{1 - P(k_i|\bar{R})}{P(k_i|\bar{R})} \right)$$

which is a key expression for ranking computation in the probabilistic model.

Since we do not know the set  $\bar{R}$  at the beginning, it is necessary to devise a method for initially computing the probabilities  $P(k_i|R)$  and  $P(k_i|\bar{R})$ . There are many alternatives for such computation. We discuss a couple of them below.

In the very beginning (i.e., immediately after the query specification), there are no retrieved documents. Thus, one has to make simplifying assumptions such as: (a) assume that  $P(k_i|R)$  is constant for all index terms  $k_i$  (typically, equal to 0.5) and (b) assume that the distribution of index terms among the non-relevant documents can be approximated by the distribution of index terms among all the documents in the collection. These two assumptions yield

$$P(k_i|R) = 0.5$$

$$P(k_i|\bar{R}) = \frac{n_i}{N}$$

where, as already defined,  $n_i$  is the number of documents which contain the index term  $k_i$  and  $N$  is the total number of documents in the collection. Given this initial guess, we can then retrieve documents which contain query terms and provide an initial probabilistic ranking for them. After that, this initial ranking is improved as follows.

Let  $V$  be a subset of the documents initially retrieved and ranked by the probabilistic model. Such a subset can be defined, for instance, as the top  $r$  ranked documents where  $r$  is a previously defined threshold. Further, let  $V_i$  be the subset of  $V$  composed of the documents in  $V$  which contain the index term  $k_i$ . For simplicity, we also use  $V'$  and  $V_i'$  to refer to the number of elements in these sets (it should always be clear when the used variable refers to the set or to the number of elements in it). For improving the probabilistic ranking, we need to improve our guesses for  $P(k_i|R)$  and  $P(k_i|\bar{R})$ . This can be accomplished with the following assumptions: (a) we can approximate  $P(k_i|R)$  by the distribution of the index term  $k_i$  among the documents retrieved so far, and (b) we can approximate  $P(k_i|\bar{R})$  by considering that all the non-retrieved documents are not relevant. With these assumptions, we can write,

$$P(k_i|R) = \frac{V_i}{V}$$

$$P(k_i|\bar{R}) = \frac{n_i - V_i}{N - V}$$

This process can then be repeated recursively. By doing so, we are able to improve on our guesses for the probabilities  $P(k_i|R)$  and  $P(k_i|\bar{R})$  without any assistance from a human subject (contrary to the original idea). However, we can also use assistance from the user for definition of the subset  $V$  as originally conceived.

The last formulas for  $P(k_i|R)$  and  $P(k_i|\bar{R})$  pose problems for small values of  $V$  and  $V_i$  which arise in practice (such as  $V = 1$  and  $V_i = 0$ ). To circumvent these problems, an adjustment factor is often added in which yields

$$P(k_i|R) = \frac{V_i + 0.5}{V + 1}$$

$$P(k_i|\bar{R}) = \frac{n_i - V_i + 0.5}{N - V + 1}$$

Handwritten notes:  $V_i$  is the number of relevant documents,  $n_i$  is the number of relevant documents.

An adjustment factor which is constant and equal to 0.5 is not always satisfactory. An alternative is to take the fraction  $n_i/N$  as the adjustment factor which yields

$$\begin{aligned} P(k_i|R) &= \frac{V_i + \frac{n_i}{N}}{V + 1} \\ P(k_i|\bar{R}) &= \frac{n_i - V_i + \frac{n_i}{N}}{N - V + 1} \end{aligned}$$

This completes our discussion of the probabilistic model.

The main *advantage* of the probabilistic model, in theory, is that documents are ranked in decreasing order of their *probability* of being relevant. The *disadvantages* include: (1) the need to guess the initial separation of documents into relevant and non-relevant sets; (2) the fact that the method does *not* take into account the frequency with which an index term occurs inside a document (i.e., all weights are binary); and (3) the adoption of the independence assumption for index terms. However, as discussed for the vector model, it is not clear that independence of index terms is a bad assumption in practical situations.

### 2.5.5 Brief Comparison of Classic Models

In general, the Boolean model is considered to be the weakest classic method. Its main problem is the inability to recognize partial matches which frequently leads to poor performance. There is some controversy as to whether the probabilistic model outperforms the vector model. Croft performed some experiments and suggested that the probabilistic model provides a better retrieval performance. However, experiments done afterwards by Salton and Buckley refute that claim. Through several different measures, Salton and Buckley showed that the vector model is *expected* to outperform the probabilistic model with general collections. This also seems to be the dominant thought among researchers, practitioners, and the Web community, where the popularity of the vector model runs high.

## 2.6 Alternative Set Theoretic Models

In this section, we discuss two alternative set theoretic models, namely the fuzzy set model and the extended Boolean model.

### 2.6.1 Fuzzy Set Model

Representing documents and queries through sets of keywords yields descriptions which are only partially related to the real semantic contents of the respective documents and queries. As a result, the matching of a document to the query terms is approximate (or vague). This can be modeled by considering that each

query term defines a *fuzzy set* and that each document has a *degree of membership* (usually smaller than 1) in this set. This interpretation of the retrieval process (in terms of concepts from fuzzy theory) is the basic foundation of the various fuzzy set models for information retrieval which have been proposed over the years. Instead of reviewing several of these models here, we focus on a particular one whose description fits well with the models already covered in this chapter. Thus, our discussion is based on the fuzzy set model for information retrieval proposed by Ogawa, Morita, and Kobayashi [616]. Before proceeding, we briefly introduce some fundamental concepts.

#### Fuzzy Set Theory

*Fuzzy set theory* [846] deals with the representation of classes whose boundaries are not well defined. The key idea is to associate a membership function with the elements of the class. This function takes values in the interval  $[0, 1]$  with 0 corresponding to no membership in the class and 1 corresponding to full membership. Membership values between 0 and 1 indicate *marginal* elements of the class. Thus, membership in a fuzzy set is a notion intrinsically *gradual* instead of abrupt (as in conventional Boolean logic).

**Definition** A fuzzy subset  $A$  of a universe of discourse  $U$  is characterized by a membership function  $\mu_A : U \rightarrow [0, 1]$  which associates with each element  $u$  of  $U$  a number  $\mu_A(u)$  in the interval  $[0, 1]$ .

The three most commonly used operations on fuzzy sets are: the complement of a fuzzy set, the union of two or more fuzzy sets, and the intersection of two or more fuzzy sets. They are defined as follows.

**Definition** Let  $U$  be the universe of discourse,  $A$  and  $B$  be two fuzzy subsets of  $U$ , and  $\bar{A}$  be the complement of  $A$  relative to  $U$ . Also, let  $u$  be an element of  $U$ . Then,

$$\begin{aligned} \mu_{\bar{A}}(u) &= 1 - \mu_A(u) \\ \mu_{A \cup B}(u) &= \max(\mu_A(u), \mu_B(u)) \\ \mu_{A \cap B}(u) &= \min(\mu_A(u), \mu_B(u)) \end{aligned}$$

Fuzzy sets are useful for representing vagueness and imprecision and have been applied to various domains. In what follows, we discuss their application to information retrieval.

#### Fuzzy Information Retrieval

As discussed in Chapters 5 and 7, one additional approach to modeling the information retrieval process is to adopt a thesaurus (which defines term re-