

# PROBABILISTIC INFORMATION RETRIEVAL

Many authors divide the classical methods of information retrieval into three categories:

- **Boolean** (based on set theory, and Boolean logic)
- **Vector space** (based on linear algebra)
- **Probabilistic** (based on Bayesian statistics)

Probabilistic methods are one of the oldest but also one of the currently hottest topics in IR.

The *Probabilistic* model of IR (*PIR*)  
is a classical model of IR.

## Basic concept:

The probability that a document is relevant to a query is assumed to depend on

- the terms in the query and
- the terms used to index the document, only.

Given

- a user query  $q$  and
- a document  $d$  in the collection,

the probabilistic model **estimates the probability** that the user will find  $d$  relevant.

## Initial probabilities:

Given a query  $q$  and a document  $d$  the model needs an estimate of the **probability** that the user finds  $d$  relevant. i.e.,  $P(R | d)$ .

## Similarity measure:

$s(q, d)$ , the similarity of  $d$  to  $q$ , is the ratio:

probability that document  $d$  is relevant to  $q$   
probability that document  $d$  is not relevant to  $q$

A mathematical formalisation of the Probabilistic Retrieval:

Given

- documents  $D = \{d_1, \dots, d_j, \dots, d_m\}$
- query  $q$ ,
- terms  $T = \{t_1, \dots, t_i, \dots, t_n\}$ ,
- binary weights matrix

$$\text{matrix } TD = (w_{ij})_{n \times m}, (i = 1, \dots, n, j = 1, \dots, m)$$

where

$$w_{ij} = \begin{cases} 1 & \text{if } t_i \text{ occurs in } d_j \\ 0 & \text{otherwise} \end{cases},$$

Let

$P(R)$ : probability that a randomly selected document is relevant to query  $q$ .

$P(\bar{R})$ : probability that a randomly selected document is not relevant to query  $q$ .

## DEFINITION:

Similarity between  $q$  and  $d_j$ :

$$s(q, d_j) = \frac{P(R|d_j)}{P(\bar{R}|d_j)}$$

where

$P(R|d_j)$  means the probability that the document  $d_j$  is relevant to query  $q$ .

Using Bayes' Rule:

$$s(q, d_j) = \frac{P(R|d_j)}{P(\bar{R}|d_j)} = \frac{\frac{P(d_j|R)P(R)}{P(d_j)}}{\frac{P(d_j|\bar{R})P(\bar{R})}{P(d_j)}} = \frac{P(d_j|R)P(R)}{P(d_j|\bar{R})P(\bar{R})}$$

where

$P(R)$  is called the **prior** probability  
 $P(R|d_j)$  is called the **posterior** probability  
 $P(d_j|R)$  probability that if a relevant document is retrieved, it is  $d_j$ .

because  $P(R)$  and  $P(\overline{R})$  are constant for  $D$ ,

and we are interested in ranking documents on their similarities with  $q$ , we may take:

$$s(q, d_j) = \frac{P(d_j | R) P(R)}{P(d_j | \overline{R}) P(\overline{R})} \approx \frac{P(d_j | R)}{P(d_j | \overline{R})}$$

### **Independence assumption:**

Let

$\mathbf{t} = (t_1, t_2, \dots, t_n)$  be the term vector for document  $d_j$ .  
 $t_i = 1$  if term  $i$  is in the document and 0 otherwise.

We estimate  $P(d_j | R)$  by  $P(\mathbf{t} | R)$

If the index terms are independent

$$\begin{aligned} P(\mathbf{t} | R) &= P(t_1 \cap R) P(t_2 \cap R) \dots P(t_n \cap R) \\ &= P(t_1 | R) P(t_2 | R) \dots P(t_n | R) \\ &= \prod_{i=1}^n P(t_i | R) \end{aligned}$$

Because we assumed that terms are independent of each other, so

$$s(q, d_j) = \frac{P(d_j | R)}{P(d_j | \bar{R})} = \frac{\prod_{i=1}^n P(t_i | R)}{\prod_{i=1}^n P(t_i | \bar{R})}$$

( $d_j$  is a compound event of simultaneous events  $t_i$  which occur in  $d_j$ )

Let

$t_1, \dots, t_i, \dots, t_p$  be the terms occurring in query  $q$

Then,

$s(q, d_j)$  is computed only for query terms:

$$s(q, d_j) = \frac{\prod_{t_i \in q} P(t_i | R)}{\prod_{t_i \in q} P(t_i | \bar{R})}$$

## Retrieval method:

Step 1. Assume  $P(t_i | R) = 0.5$ , and

$$P(t_i | \bar{R}) = \frac{m_i}{m}$$

where  $i = 1, \dots, n$

$m_i$  is number of documents containing term  $t_i$

Step 2. Compute  $s(q, d_j)$ .

Produce ranking.

$V \subseteq D$  denotes retrieved documents.

Step 3. Let  $V_i \subseteq V$  denote documents containing  $t_i \in q$

Revision:

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

$$P(t_i | \bar{R}) = \frac{m_i - |V_i|}{m - |V|}$$

Step 4. Repeat from step 2.