

# Chapter IR:V

## V. Retrieval Models

- ❑ Overview of Retrieval Models
- ❑ Empirical Models
- ❑ Boolean Retrieval
- ❑ Vector Space Model
- ❑ Probabilistic Models
- ❑ Binary Independence Model
- ❑ Okapi BM25
- ❑ Hidden Variable Models
- ❑ Latent Semantic Indexing
- ❑ Explicit Semantic Analysis
- ❑ Generative Models
- ❑ Language Models
- ❑ Combining Evidence
- ❑ Web Search
- ❑ Learning to Rank

# Probabilistic Models [\[Empirical Models\]](#) [\[Hidden Variable Models\]](#) [\[Generative Models\]](#)

For a pair of query  $q$  and document  $d$ , probabilistic retrieval models estimate the probability  $P(\textit{relevance}(d, q) = 1 \mid \text{data})$  that  $d$  will be considered relevant for  $q$ .

The `data` may include anything available to a retrieval system, including query  $q$ , document  $d$ , user context, user profile, etc.

The relevance function  $\rho(d, q) = P(\textit{relevance}(d, q) = 1 \mid \text{data})$  ranks documents by their probability of being relevant.

Relevance is assumed a binary property of an individual document  $d$  given query  $q$ .

Discriminating factors of probabilistic models:

- ❑ Sample space from which probabilities are derived
- ❑ Estimation method for  $P(\textit{relevance}(d, q) = 1 \mid \text{data})$

# Probabilistic Models

## Probability Ranking Principle [\[Robertson 2009\]](#)

If retrieved documents are ordered by decreasing probability of relevance on the data available, then the system's effectiveness is the best that can be obtained for the data.

# Probabilistic Models

## Probability Ranking Principle [\[Robertson 2009\]](#)

If retrieved documents are ordered by decreasing probability of relevance on the data available, then the system's effectiveness is the best that can be obtained for the data.

The truth about the relevance of a document is unknown to the retrieval system.

It can only reason about the `data` at hand: the query  $q$ , the document  $d$  in question, and all supplementary data.

The probability  $P(\textit{relevance}(d, q) = 1 \mid \textit{data})$  corresponds to the proportion of all users in similar contexts querying for  $q$ , who would consider document  $d$  relevant.

It follows that ranking in decreasing order of that probability is the optimal strategy under uncertain and incomplete knowledge of the true relevance relation.

This maximizes commonly used effectiveness metrics: precision, recall, and others.

Most salient question: How to estimate  $P(\textit{relevance}(d, q) = 1 \mid \textit{data})$ ?

## Remarks:

- ❑ The above rendition of the probability ranking principle is a more concise version of Cooper's original one, published by [\[Robertson 1977\]](#).
- ❑ The probability ranking principle is not universal since counterexamples can be formulated where ranking in order of decreasing relevance probability leads to suboptimal rankings (think of user relevance: someone has seen a document before such that it is maybe less relevant to the user but this might not be part of `data`).

# Binary Independence Model

Retrieval Model  $\mathcal{R} = \langle \mathbf{D}, \mathbf{Q}, \rho \rangle$  [\[Generic Model\]](#) [\[Boolean Retrieval\]](#) [\[VSM\]](#) [\[BM25\]](#) [\[LSI\]](#) [\[ESA\]](#) [\[LM\]](#)

## Document representations $\mathbf{D}$ .

The set of index terms  $T = \{t_1, \dots, t_m\}$  is composed of the word stems of the vocabulary of a document collection.

The representation  $\mathbf{d}$  of a document  $d$  is a function from  $T$  to  $\{0, 1\}$ , where  $\mathbf{d}(t_i) = 1$  is interpreted as “term  $t_i$  present in  $d$ ”, and  $\mathbf{d}(t_i) = 0$  as “term  $t_i$  absent from  $d$ ”.

## Query representations $\mathbf{Q}$ .

A query representation  $\mathbf{q}$  is constructed like a document representation.

## Relevance function $\rho$ .

The relevance function  $\rho(\mathbf{d}, \mathbf{q}) = P(\textit{relevance}(d, q) = 1 \mid \mathbf{d}, \mathbf{q})$

# Binary Independence Model

Retrieval Model  $\mathcal{R} = \langle \mathbf{D}, \mathbf{Q}, \rho \rangle$  [\[Generic Model\]](#) [\[Boolean Retrieval\]](#) [\[VSM\]](#) [\[BM25\]](#) [\[LSI\]](#) [\[ESA\]](#) [\[LM\]](#)

## Document representations $\mathbf{D}$ .

The set of index terms  $T = \{t_1, \dots, t_m\}$  is composed of the word stems of the vocabulary of a document collection.

The representation  $\mathbf{d}$  of a document  $d$  is a function from  $T$  to  $\{0, 1\}$ , where  $\mathbf{d}(t_i) = 1$  is interpreted as “term  $t_i$  present in  $d$ ”, and  $\mathbf{d}(t_i) = 0$  as “term  $t_i$  absent from  $d$ ”.

## Query representations $\mathbf{Q}$ .

A query representation  $\mathbf{q}$  is constructed like a document representation.

## Relevance function $\rho$ .

The relevance function  $\rho(\mathbf{d}, \mathbf{q}) = P(\textit{relevance}(d, q) = 1 \mid \mathbf{d}, \mathbf{q})$

## Remarks:

- ❑ The model is also known as Okapi model, City model, or simply as the probabilistic model.
- ❑ The joint probability space  $(\Omega, \mathcal{P}(\Omega), P)$  underlying the binary independence model is given by the sample space  $\Omega = \{0, 1\} \times \mathcal{P}(T)$ , where  $\mathcal{P}(T)$  denotes the set of all binary document vectors over the set of terms  $T$ .



# Binary Independence Model

## Relevance Function $\rho$ : Derivation

Let  $r$  denote a random variable indicating the true (binary) *relevance*( $d, q$ ), and  $\mathbf{d}, \mathbf{q}$  random variables representing document  $d$  and query  $q$ .

$$P(r = 1 \mid \mathbf{d}, \mathbf{q}) \stackrel{\text{rank}}{=} \frac{P(r = 1 \mid \mathbf{d}, \mathbf{q})}{P(r = 0 \mid \mathbf{d}, \mathbf{q})} \quad (1)$$

(1) **Rank-preserving** replacement of  $P(A)$  by the odds  $\frac{P(A)}{P(\bar{A})}$  in favor of event  $A$ .

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

Let  $r$  denote a random variable indicating the true (binary) *relevance*( $d, q$ ), and  $\mathbf{d}, \mathbf{q}$  random variables representing document  $d$  and query  $q$ .

$$P(r = 1 \mid \mathbf{d}, \mathbf{q}) \stackrel{\text{rank}}{=} \frac{P(r = 1 \mid \mathbf{d}, \mathbf{q})}{P(r = 0 \mid \mathbf{d}, \mathbf{q})} \quad (1)$$

$$= \frac{P(\mathbf{d} \mid r = 1, \mathbf{q}) P(r = 1 \mid \mathbf{q})}{P(\mathbf{d} \mid r = 0, \mathbf{q}) P(r = 0 \mid \mathbf{q})} \quad (2)$$

(1) Rank-preserving replacement of  $P(A)$  by the odds  $\frac{P(A)}{P(\bar{A})}$  in favor of event  $A$ .

(2) Application of Bayes' rule. The common denominator  $P(\mathbf{d} \mid \mathbf{q})$  is canceled.

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

Let  $r$  denote a random variable indicating the true (binary) *relevance*( $d, q$ ), and  $\mathbf{d}, \mathbf{q}$  random variables representing document  $d$  and query  $q$ .

$$P(r = 1 \mid \mathbf{d}, \mathbf{q}) \stackrel{\text{rank}}{=} \frac{P(r = 1 \mid \mathbf{d}, \mathbf{q})}{P(r = 0 \mid \mathbf{d}, \mathbf{q})} \quad (1)$$

$$= \frac{P(\mathbf{d} \mid r = 1, \mathbf{q}) P(r = 1 \mid \mathbf{q})}{P(\mathbf{d} \mid r = 0, \mathbf{q}) P(r = 0 \mid \mathbf{q})} \quad (2)$$

$$\stackrel{\text{rank}}{=} \frac{P(\mathbf{d} \mid r = 1, \mathbf{q})}{P(\mathbf{d} \mid r = 0, \mathbf{q})} \quad (3)$$

(1) Rank-preserving replacement of  $P(A)$  by the odds  $\frac{P(A)}{P(\bar{A})}$  in favor of event  $A$ .

(2) Application of Bayes' rule. The common denominator  $P(\mathbf{d} \mid \mathbf{q})$  is canceled.

(3) Rank-preserving omission of  $\frac{P(r = 1 \mid \mathbf{q})}{P(r = 0 \mid \mathbf{q})}$ ; it does not depend on  $\mathbf{d}$ .

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

Let  $d(t)$  denote a random variable indicating if  $t$  occurs in  $d$ .

$$\frac{P(\mathbf{d} \mid r = 1, \mathbf{q})}{P(\mathbf{d} \mid r = 0, \mathbf{q})} = \prod_{t \in T} \frac{P(d(t) \mid r = 1, \mathbf{q})}{P(d(t) \mid r = 0, \mathbf{q})} \quad (4)$$

(4) Assuming independence between terms.

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

Let  $\mathbf{d}(t)$  denote a random variable indicating if  $t$  occurs in  $d$ .

$$\frac{P(\mathbf{d} \mid r = 1, \mathbf{q})}{P(\mathbf{d} \mid r = 0, \mathbf{q})} = \prod_{t \in T} \frac{P(\mathbf{d}(t) \mid r = 1, \mathbf{q})}{P(\mathbf{d}(t) \mid r = 0, \mathbf{q})} \quad (4)$$

$$= \prod_{t \in \mathbf{d}} \frac{P(\mathbf{d}(t) = 1 \mid r = 1, \mathbf{q})}{P(\mathbf{d}(t) = 1 \mid r = 0, \mathbf{q})} \prod_{t \notin \mathbf{d}} \frac{P(\mathbf{d}(t) = 0 \mid r = 1, \mathbf{q})}{P(\mathbf{d}(t) = 0 \mid r = 0, \mathbf{q})} \quad (5)$$

(4) Assuming independence between terms.

(5) Separation of the the two possible cases for  $\mathbf{d}(t)$ ,  
where  $t \in \mathbf{d}$  means  $t \in T : \mathbf{d}(t) = 1$  and  $t \notin \mathbf{d}$  means  $t \in T : \mathbf{d}(t) = 0$ .

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

Let  $\mathbf{d}(t)$  denote a random variable indicating if  $t$  occurs in  $d$ .

$$\frac{P(\mathbf{d} \mid r = 1, \mathbf{q})}{P(\mathbf{d} \mid r = 0, \mathbf{q})} = \prod_{t \in T} \frac{P(\mathbf{d}(t) \mid r = 1, \mathbf{q})}{P(\mathbf{d}(t) \mid r = 0, \mathbf{q})} \quad (4)$$

$$= \prod_{t \in \mathbf{d}} \frac{P(\mathbf{d}(t) = 1 \mid r = 1, \mathbf{q})}{P(\mathbf{d}(t) = 1 \mid r = 0, \mathbf{q})} \prod_{t \notin \mathbf{d}} \frac{P(\mathbf{d}(t) = 0 \mid r = 1, \mathbf{q})}{P(\mathbf{d}(t) = 0 \mid r = 0, \mathbf{q})} \quad (5)$$

$$= \prod_{t \in \mathbf{d}} \frac{p_t}{s_t} \prod_{t \notin \mathbf{d}} \frac{1 - p_t}{1 - s_t} \quad (6)$$

(4) Assuming independence between terms.

(5) Separation of the the two possible cases for  $\mathbf{d}(t)$ ,  
where  $t \in \mathbf{d}$  means  $t \in T : \mathbf{d}(t) = 1$  and  $t \notin \mathbf{d}$  means  $t \in T : \mathbf{d}(t) = 0$ .

(6) Abbreviation:  $p_t = P(\mathbf{d}(t) = 1 \mid r = 1, \mathbf{q})$  and  $s_t = P(\mathbf{d}(t) = 1 \mid r = 0, \mathbf{q})$ .

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

$$\prod_{t \in \mathbf{d}} \frac{p_t}{s_t} \prod_{t \notin \mathbf{d}} \frac{1 - p_t}{1 - s_t} \stackrel{\text{rank}}{=} \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t}{s_t} \prod_{\substack{t \in \mathbf{q}: \\ t \notin \mathbf{d}}} \frac{1 - p_t}{1 - s_t} \quad (7)$$

- (7) Assumption that if  $t \notin \mathbf{q}$ , then  $p_t = s_t$  (i.e.,  $t$  is equally likely to occur in relevant and irrelevant documents). Renders all corresponding factors idempotent.

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

$$\prod_{t \in \mathbf{d}} \frac{p_t}{s_t} \prod_{t \notin \mathbf{d}} \frac{1 - p_t}{1 - s_t} \stackrel{\text{rank}}{=} \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t}{s_t} \prod_{\substack{t \in \mathbf{q}: \\ t \notin \mathbf{d}}} \frac{1 - p_t}{1 - s_t} \quad (7)$$

$$= \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t}{s_t} \prod_{\substack{t \in \mathbf{q}: \\ t \notin \mathbf{d}}} \frac{1 - p_t}{1 - s_t} \quad (8)$$

- (7) Assumption that if  $t \notin \mathbf{q}$ , then  $p_t = s_t$  (i.e.,  $t$  is equally likely to occur in relevant and irrelevant documents). Renders all corresponding factors idempotent.
- (8) Addition of all missing query terms to the right product and division by the added factors to fulfill the equation.



# Binary Independence Model

## Relevance Function $\rho$ : Derivation

$$\prod_{t \in \mathbf{d}} \frac{p_t}{s_t} \prod_{t \notin \mathbf{d}} \frac{1 - p_t}{1 - s_t} \stackrel{\text{rank}}{=} \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t}{s_t} \prod_{\substack{t \in \mathbf{q}: \\ t \notin \mathbf{d}}} \frac{1 - p_t}{1 - s_t} \quad (7)$$

$$= \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t}{s_t} \prod_{t \in \mathbf{q}} \frac{1 - p_t}{1 - s_t} / \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{1 - p_t}{1 - s_t} \quad (8)$$

- (7) Assumption that if  $t \notin \mathbf{q}$ , then  $p_t = s_t$  (i.e.,  $t$  is equally likely to occur in relevant and irrelevant documents). Renders all corresponding factors idempotent.
- (8) Addition of all missing query terms to the right product and division by the added factors to fulfill the equation.

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

$$\prod_{t \in \mathbf{d}} \frac{p_t}{s_t} \prod_{t \notin \mathbf{d}} \frac{1 - p_t}{1 - s_t} \stackrel{\text{rank}}{=} \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t}{s_t} \prod_{\substack{t \in \mathbf{q}: \\ t \notin \mathbf{d}}} \frac{1 - p_t}{1 - s_t} \quad (7)$$

$$= \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t}{s_t} \prod_{t \in \mathbf{q}} \frac{1 - p_t}{1 - s_t} \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{1 - s_t}{1 - p_t} \quad (8)$$

- (7) Assumption that if  $t \notin \mathbf{q}$ , then  $p_t = s_t$  (i.e.,  $t$  is equally likely to occur in relevant and irrelevant documents). Renders all corresponding factors idempotent.
- (8) Addition of all missing query terms to the right product and division by the added factors to fulfill the equation.

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

$$\prod_{t \in \mathbf{d}} \frac{p_t}{s_t} \prod_{t \notin \mathbf{d}} \frac{1 - p_t}{1 - s_t} \stackrel{\text{rank}}{=} \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t}{s_t} \prod_{\substack{t \in \mathbf{q}: \\ t \notin \mathbf{d}}} \frac{1 - p_t}{1 - s_t} \quad (7)$$

$$= \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t}{s_t} \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{1 - s_t}{1 - p_t} \prod_{t \in \mathbf{q}} \frac{1 - p_t}{1 - s_t} \quad (8)$$

- (7) Assumption that if  $t \notin \mathbf{q}$ , then  $p_t = s_t$  (i.e.,  $t$  is equally likely to occur in relevant and irrelevant documents). Renders all corresponding factors idempotent.
- (8) Addition of all missing query terms to the right product and division by the added factors to fulfill the equation.

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

$$\prod_{t \in \mathbf{d}} \frac{p_t}{s_t} \prod_{t \notin \mathbf{d}} \frac{1 - p_t}{1 - s_t} \stackrel{\text{rank}}{=} \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t}{s_t} \prod_{\substack{t \in \mathbf{q}: \\ t \notin \mathbf{d}}} \frac{1 - p_t}{1 - s_t} \quad (7)$$

$$= \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t(1 - s_t)}{s_t(1 - p_t)} \prod_{t \in \mathbf{q}} \frac{1 - p_t}{1 - s_t} \quad (8)$$

- (7) Assumption that if  $t \notin \mathbf{q}$ , then  $p_t = s_t$  (i.e.,  $t$  is equally likely to occur in relevant and irrelevant documents). Renders all corresponding factors idempotent.
- (8) Addition of all missing query terms to the right product and division by the added factors to fulfill the equation.

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

$$\prod_{t \in \mathbf{d}} \frac{p_t}{s_t} \prod_{t \notin \mathbf{d}} \frac{1 - p_t}{1 - s_t} \stackrel{\text{rank}}{=} \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t}{s_t} \prod_{\substack{t \in \mathbf{q}: \\ t \notin \mathbf{d}}} \frac{1 - p_t}{1 - s_t} \quad (7)$$

$$= \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t(1 - s_t)}{s_t(1 - p_t)} \prod_{t \in \mathbf{q}} \frac{1 - p_t}{1 - s_t} \quad (8)$$

$$\stackrel{\text{rank}}{=} \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t(1 - s_t)}{s_t(1 - p_t)} \quad (9)$$

- (7) Assumption that if  $t \notin \mathbf{q}$ , then  $p_t = s_t$  (i.e.,  $t$  is equally likely to occur in relevant and irrelevant documents). Renders all corresponding factors idempotent.
- (8) Addition of all missing query terms to the right product and division by the added factors to fulfill the equation.
- (9) Rank-preserving omission of the right product; it does not depend on  $\mathbf{d}$ .

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

$$\prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t(1 - s_t)}{s_t(1 - p_t)} \stackrel{\text{rank}}{=} \log \prod_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \frac{p_t(1 - s_t)}{s_t(1 - p_t)} \quad (10)$$

- (10) Rank-preserving logarithmization to allow for computations that do not underflow common floating point number formats.

# Binary Independence Model

## Relevance Function $\rho$ : Derivation

$$\begin{aligned} \prod_{\substack{t \in q: \\ t \in d}} \frac{p_t(1 - s_t)}{s_t(1 - p_t)} &\stackrel{\text{rank}}{=} \log \prod_{\substack{t \in q: \\ t \in d}} \frac{p_t(1 - s_t)}{s_t(1 - p_t)} \\ &= \sum_{\substack{t \in q: \\ t \in d}} \underbrace{\log \frac{p_t(1 - s_t)}{s_t(1 - p_t)}}_{:= \omega^{\text{RSJ}}} \end{aligned} \quad (10)$$

(10) Rank-preserving logarithmization to allow for computations that do not underflow common floating point number formats.

In effect, we accumulate for each term  $t \in q$  the log odds ratio of the **odds in favor** and the **odds against**  $t$  occurring in  $d$  if the document  $d$  is irrelevant to query  $q$ .

This ratio tells us how much more likely it is that  $t$  occurs in  $d$  if  $d$  is relevant to  $q$ .

RSJ  $\sim$  Robertson Spärck-Jones

# Binary Independence Model

## Relevance Function $\rho$ : Estimation

Let  $D$  denote the document collection and  $D_t$  the subset containing term  $t$ .

$$\sum_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \log \frac{p_t(1 - s_t)}{s_t(1 - p_t)} \stackrel{\text{rank}}{=} \sum_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \log \frac{1 - s_t}{s_t} \quad (11)$$

(11) Assumption that a term  $t \in \mathbf{q}$  is equally likely to be present or absent in a random relevant document:  $p_t = 0.5$ . This cancels  $p_t$  and  $1 - p_t$ .



# Binary Independence Model

## Relevance Function $\rho$ : Estimation

Let  $D$  denote the document collection and  $D_t$  the subset containing term  $t$ .

$$\sum_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \log \frac{p_t(1 - s_t)}{s_t(1 - p_t)} \stackrel{\text{rank}}{=} \sum_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \log \frac{1 - s_t}{s_t} \quad (11)$$

$$\stackrel{\text{rank}}{=} \sum_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \log \frac{|D| - |D_t| + 0.5}{|D_t| + 0.5} \quad (12)$$

(11) Assumption that a term  $t \in \mathbf{q}$  is equally likely to be present or absent in a random relevant document:  $p_t = 0.5$ . This cancels  $p_t$  and  $1 - p_t$ .

(12) Maximum likelihood estimation of  $s_t = P(\mathbf{d}(t) = 1 \mid r = 0, \mathbf{q})$ :

$$s_t = \frac{|D_t| + 0.5}{|D| + 1.0},$$

where adding 0.5 (1.0) is used for smoothing (avoiding zeros). Since most documents in  $D$  are irrelevant, this estimation is quite accurate.

## Remarks:

- Adding 0.5 (1.0) 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 relative is referred to as the relative frequency of the event. Estimating the probability as the relative frequency is the maximum likelihood estimate (or MLE), maximum 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  $\alpha$  ( $\beta$ ) to each of the observed counts (totals). These pseudocounts correspond to the use of a uniform distribution over the vocabulary as a Bayesian prior. 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  $\alpha = 0.5, \beta = 1.0$ . 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.

[Manning/Raghavan/Schütze 2008]

# Binary Independence Model

## Relevance Function $\rho$ : Example

$$q = (\mathbf{a} \ \mathbf{c} \ \mathbf{h})$$

Token $t$	<b>a</b>	b	<b>c</b>	d	e	f	g	<b>h</b>
$ D_t $	2	6	2	3	3	1	3	1
$s_t$	0.4	0.9	0.4	0.5	0.5	0.2	0.5	0.2

$$D = \{ \begin{array}{l} d_1 = (\mathbf{a} \ b \ \mathbf{c} \ b \ d), \\ d_2 = (b \ e \ f \ b), \\ d_3 = (b \ g \ \mathbf{c} \ d), \\ d_4 = (b \ d \ e), \\ d_5 = (\mathbf{a} \ b \ e \ g), \\ d_6 = (b \ g \ \mathbf{h}) \end{array} \}$$

$$\begin{aligned} \rho(\mathbf{d}_i, \mathbf{q}) &= \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{a}} + \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{c}} + \underbrace{\log \frac{1 - 0.2}{0.2}}_{t = \mathbf{h}} \\ &= \log 1.5 \quad + \quad \log 1.5 \quad + \quad \log 4 \\ &= 0.1761 \quad + \quad 0.1761 \quad + \quad 0.6021 \\ &= 0.9543 \end{aligned}$$

Document	$\rho(\mathbf{d}_i, \mathbf{q})$
$d_1$	
$d_2$	
$d_3$	
$d_4$	
$d_5$	
$d_6$	

# Binary Independence Model

## Relevance Function $\rho$ : Example

$$q = (\mathbf{a} \ \mathbf{c} \ \mathbf{h})$$

Token $t$	<b>a</b>	b	<b>c</b>	d	e	f	g	<b>h</b>
$ D_t $	2	6	2	3	3	1	3	1
$s_t$	0.4	0.9	0.4	0.5	0.5	0.2	0.5	0.2

$$D = \{ \begin{array}{l} d_1 = (\mathbf{a} \ b \ \mathbf{c} \ b \ d), \\ d_2 = (b \ e \ f \ b), \\ d_3 = (b \ g \ \mathbf{c} \ d), \\ d_4 = (b \ d \ e), \\ d_5 = (\mathbf{a} \ b \ e \ g), \\ d_6 = (b \ g \ \mathbf{h}) \end{array} \}$$

$$\begin{aligned} \rho(\mathbf{d}_1, \mathbf{q}) &= \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{a}} + \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{c}} + \underbrace{\log \frac{1 - 0.2}{0.2}}_{t = \mathbf{h}} \\ &= \log 1.5 + \log 1.5 + \log 4 \\ &= 0.1761 + 0.1761 + 0.6021 \\ &= 0.3522 \end{aligned}$$

Document	$\rho(\mathbf{d}_i, \mathbf{q})$
$d_1$	0.3522
$d_2$	
$d_3$	
$d_4$	
$d_5$	
$d_6$	

# Binary Independence Model

## Relevance Function $\rho$ : Example

$$q = (\mathbf{a} \ \mathbf{c} \ \mathbf{h})$$

Token $t$	<b>a</b>	b	<b>c</b>	d	e	f	g	<b>h</b>
$ D_t $	2	6	2	3	3	1	3	1
$s_t$	0.4	0.9	0.4	0.5	0.5	0.2	0.5	0.2

$$D = \{ \begin{array}{l} d_1 = (\mathbf{a} \ b \ \mathbf{c} \ b \ d), \\ d_2 = (b \ e \ f \ b), \\ d_3 = (b \ g \ \mathbf{c} \ d), \\ d_4 = (b \ d \ e), \\ d_5 = (\mathbf{a} \ b \ e \ g), \\ d_6 = (b \ g \ \mathbf{h}) \end{array} \}$$

$$\begin{aligned} \rho(\mathbf{d}_2, \mathbf{q}) &= \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{a}} + \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{c}} + \underbrace{\log \frac{1 - 0.2}{0.2}}_{t = \mathbf{h}} \\ &= \log 1.5 + \log 1.5 + \log 4 \\ &= 0.1761 + 0.1761 + 0.6021 \\ &= 0 \end{aligned}$$

Document	$\rho(\mathbf{d}_i, \mathbf{q})$
$d_1$	0.3522
$d_2$	0
$d_3$	
$d_4$	
$d_5$	
$d_6$	

# Binary Independence Model

## Relevance Function $\rho$ : Example

$$q = (\mathbf{a} \ \mathbf{c} \ \mathbf{h})$$

Token $t$	<b>a</b>	b	<b>c</b>	d	e	f	g	<b>h</b>
$ D_t $	2	6	2	3	3	1	3	1
$s_t$	0.4	0.9	0.4	0.5	0.5	0.2	0.5	0.2

$$D = \{ \begin{aligned} d_1 &= (\mathbf{a} \ b \ \mathbf{c} \ b \ d), \\ d_2 &= (b \ e \ f \ b), \\ d_3 &= (b \ g \ \mathbf{c} \ d), \\ d_4 &= (b \ d \ e), \\ d_5 &= (\mathbf{a} \ b \ e \ g), \\ d_6 &= (b \ g \ \mathbf{h}) \end{aligned} \}$$

$$\begin{aligned} \rho(\mathbf{d}_3, \mathbf{q}) &= \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{a}} + \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{c}} + \underbrace{\log \frac{1 - 0.2}{0.2}}_{t = \mathbf{h}} \\ &= \log 1.5 + \log 1.5 + \log 4 \\ &= 0.1761 + 0.1761 + 0.6021 \\ &= 0.1761 \end{aligned}$$

Document	$\rho(\mathbf{d}_i, \mathbf{q})$
$d_1$	0.3522
$d_2$	0
$d_3$	0.1761
$d_4$	
$d_5$	
$d_6$	

# Binary Independence Model

## Relevance Function $\rho$ : Example

$$q = (\mathbf{a} \ \mathbf{c} \ \mathbf{h})$$

Token $t$	<b>a</b>	b	<b>c</b>	d	e	f	g	<b>h</b>
$ D_t $	2	6	2	3	3	1	3	1
$s_t$	0.4	0.9	0.4	0.5	0.5	0.2	0.5	0.2

$$D = \{ \begin{array}{l} d_1 = (\mathbf{a} \ b \ \mathbf{c} \ b \ d), \\ d_2 = (b \ e \ f \ b), \\ d_3 = (b \ g \ \mathbf{c} \ d), \\ d_4 = (b \ d \ e), \\ d_5 = (\mathbf{a} \ b \ e \ g), \\ d_6 = (b \ g \ \mathbf{h}) \end{array} \}$$

$$\begin{aligned} \rho(\mathbf{d}_6, \mathbf{q}) &= \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{a}} + \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{c}} + \underbrace{\log \frac{1 - 0.2}{0.2}}_{t = \mathbf{h}} \\ &= \log 1.5 + \log 1.5 + \log 4 \\ &= 0.1761 + 0.1761 + 0.6021 \\ &= 0.6021 \end{aligned}$$

Document	$\rho(\mathbf{d}_i, \mathbf{q})$
$d_1$	0.3522
$d_2$	0
$d_3$	0.1761
$d_4$	0
$d_5$	0.1761
$d_6$	0.6021

# Binary Independence Model

## Relevance Function $\rho$ : Example

$$q = (\mathbf{a} \ \mathbf{c} \ \mathbf{h})$$

Token $t$	<b>a</b>	b	<b>c</b>	d	e	f	g	<b>h</b>
$ D_t $	2	6	2	3	3	1	3	1
$s_t$	0.4	0.9	0.4	0.5	0.5	0.2	0.5	0.2

$$D = \{ \begin{array}{l} d_1 = (\mathbf{a} \ b \ \mathbf{c} \ b \ d), \\ d_2 = (b \ e \ f \ b), \\ d_3 = (b \ g \ \mathbf{c} \ d), \\ d_4 = (b \ d \ e), \\ d_5 = (\mathbf{a} \ b \ e \ g), \\ d_6 = (b \ g \ \mathbf{h}) \end{array} \}$$

$$\begin{aligned} \rho(\mathbf{d}_6, \mathbf{q}) &= \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{a}} + \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{c}} + \underbrace{\log \frac{1 - 0.2}{0.2}}_{t = \mathbf{h}} \\ &= \log 1.5 + \log 1.5 + \log 4 \\ &= 0.1761 + 0.1761 + 0.6021 \\ &= 0.6021 \end{aligned}$$

Ranking	$\rho(\mathbf{d}_i, \mathbf{q})$
$d_6$	0.6021
$d_1$	0.3522
$d_3$	0.1761
$d_5$	0.1761
$d_2$	0
$d_4$	0



# Binary Independence Model

## Relevance Function $\rho$ : Example

$$q = (\mathbf{a} \ \mathbf{c} \ \mathbf{h})$$

Token $t$	<b>a</b>	b	<b>c</b>	d	e	f	g	<b>h</b>
$ D_t $	2	6	2	3	3	1	3	1
$s_t$	0.4	0.9	0.4	0.5	0.5	0.2	0.5	0.2

$$D = \{ \begin{array}{l} d_1 = (\mathbf{a} \ b \ \mathbf{c} \ b \ d), \\ d_2 = (b \ e \ f \ b), \\ d_3 = (b \ g \ \mathbf{c} \ d), \\ d_4 = (b \ d \ e), \\ d_5 = (\mathbf{a} \ b \ e \ g), \\ d_6 = (b \ g \ \mathbf{h}) \end{array} \}$$

$$\begin{aligned} \rho(\mathbf{d}_6, \mathbf{q}) &= \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{a}} + \underbrace{\log \frac{1 - 0.4}{0.4}}_{t = \mathbf{c}} + \underbrace{\log \frac{1 - 0.2}{0.2}}_{t = \mathbf{h}} \\ &= \log 1.5 + \log 1.5 + \log 4 \\ &= 0.1761 + 0.1761 + 0.6021 \\ &= 0.6021 \end{aligned}$$

Ranking	$\rho(\mathbf{d}_i, \mathbf{q})$
$d_6$	0.6021
$d_1$	0.3522
$d_3$	0.1761
$d_5$	0.1761
$d_2$	0
$d_4$	0

Why is  $d_6$  the most relevant document?

# Binary Independence Model

Relevance Function  $\rho$ : Summary [Inverse Document Frequency]

$$\rho(\mathbf{d}, \mathbf{q}) = P(r = 1 \mid \mathbf{d}, \mathbf{q}) \propto \sum_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \underbrace{\log \frac{|D| - |D_t| + 0.5}{|D_t| + 0.5}}_{:= \omega^{\text{BIM}} \approx \text{idf}(t, D)}$$

Assumptions:

1. Binary relevance of a document  $d$  to a query  $q$ , independent of all other documents.
2. Boolean representations  $\mathbf{d}, \mathbf{q}$  of document  $d$  and query  $q$ .
3. Independence of word occurrence in documents.
4. Terms not in query  $\mathbf{q}$  are equally likely to occur in relevant and irrelevant documents.
5. Terms in  $\mathbf{q}$  are equally likely to occur or not to occur in a relevant document  $d$ .
6. The set of irrelevant documents is represented by the entire collection.

## Remarks:

- ❑ It is a personal observation that almost every mathematically inclined graduate student in Information Retrieval attempts to formulate some sort of a non-independent model of IR within the first two to three years of his or her studies. The vast majority of these attempts yield no improvements and remain unpublished. [...] It is natural to wonder why this is the case – the classical model contains an obviously incorrect assumption about the language, and yet most attempts to relax that assumption produce no consistent improvements whatsoever. Contrary to popular belief, word independence is not a necessary assumption in the classical probabilistic model of IR. A necessary and sufficient condition is proportional interdependence [...]: *on average*, all the words in a given document have about as much interdependence under the relevant class as they do under the non-relevant class. [...] the only requirement is that whatever disbalance exists be constant across all documents. If there is anything wrong with the classical model, it is not independence but the assumptions made in the estimation process. [\[Lavrenko 2009\]](#)

# Binary Independence Model

## Query Refinement: Relevance Feedback

Let  $R$  denote a result set for  $q$ , where  $R^+$  and  $R^-$  are subsets of relevant and irrelevant results from relevance feedback, and  $R_t$  the subset containing term  $t$ .

$$\sum_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \log \frac{p_t(1 - s_t)}{s_t(1 - p_t)} \stackrel{\text{rank}}{=} \sum_{\substack{t \in \mathbf{q}: \\ t \in \mathbf{d}}} \log \frac{(|R_t^+| + 0.5)(|R_t^-| - |R^-| - 0.5)}{(|R_t^-| + 0.5)(|R_t^+| - |R^+| - 0.5)} \quad (13)$$

(13) Assumption that partial knowledge of the true relevance function  $r$  is at hand:

$t \in \mathbf{q}$	Relevant	Irrelevant	$\Sigma$
$t \in \mathbf{d}$	$ R_t^+ $	$ R_t^- $	$ R_t $
$t \notin \mathbf{d}$	$ R^+ - R_t^+ $	$ R^- - R_t^- $	$ R - R_t $
$\Sigma$	$ R^+ $	$ R^- $	$ R $

Maximum likelihood estimation of  $p_t = P(\mathbf{d}(t) = 1 \mid r = 1, \mathbf{q})$  and  $s_t = P(\mathbf{d}(t) = 1 \mid r = 0, \mathbf{q})$  with smoothing:

$$p_t = \frac{|R_t^+| + 0.5}{|R^+| + 1.0}; \quad s_t = \frac{|R_t^-| + 0.5}{|R^-| + 1.0}.$$

# Binary Independence Model

## Query Refinement: Relevance Feedback Example

$$q = (\mathbf{b} \ \mathbf{g} \ \mathbf{h})$$

Term $t$	a	<b>b</b>	c	d	e	f	<b>g</b>	<b>h</b>
$ R_t^+ $	1	2	1	1	1	1	0	0
$ R_t^- $	1	3	1	2	2	0	2	0
$p_t$	0.5	0.8	0.5	0.5	0.5	0.5	0.2	0.2
$s_t$	0.4	0.9	0.4	0.6	0.6	0.1	0.6	0.1

$$R = R^+ \cup R^-$$

$$R^+ = \{ d_1 = (a \ \mathbf{b} \ c \ \mathbf{b} \ d), \\ d_2 = (\mathbf{b} \ e \ f \ \mathbf{b}) \}$$

$$R^- = \{ d_3 = (\mathbf{b} \ \mathbf{g} \ c \ d), \\ d_4 = (\mathbf{b} \ d \ e), \\ d_5 = (a \ \mathbf{b} \ e \ \mathbf{g}) \}$$

$$d_6 = (\mathbf{b} \ \mathbf{g} \ \mathbf{h})$$

$$\begin{aligned}
 \rho(\mathbf{d}_6, \mathbf{q}) &= \underbrace{\log \frac{0.8 \cdot (1 - 0.9)}{0.9 \cdot (1 - 0.8)}}_{t = \mathbf{b}} + \underbrace{\log \frac{0.2 \cdot (1 - 0.6)}{0.6 \cdot (1 - 0.2)}}_{t = \mathbf{g}} + \underbrace{\log \frac{0.2 \cdot (1 - 0.1)}{0.1 \cdot (1 - 0.2)}}_{t = \mathbf{h}} \\
 &= \log 0.44 \quad + \quad \log 0.17 \quad + \quad \log 2.25 \\
 &= -0.3522 \quad + \quad -0.7722 \quad + \quad 0.3522 \\
 &= -0.7722
 \end{aligned}$$

# Binary Independence Model

## Query Refinement: Relevance Feedback

Initialization of  $p_t = 0.5$  and  $s_t$  as in the absence of relevance feedback. Updating of  $p_t$  based only on a given  $R$  may be unreliable, especially when  $|R|$  is small. A better update formula for  $p_t$  is as follows:

$$p'_t = \frac{|R_t^+| + \alpha p_t}{|R^+| + \alpha},$$

where  $\alpha$  adjusts the contribution of the previous feedback cycle (e.g.,  $\alpha = |R|$ ).

# Binary Independence Model

## Query Refinement: Relevance Feedback

Initialization of  $p_t = 0.5$  and  $s_t$  as in the absence of relevance feedback. Updating of  $p_t$  based only on a given  $R$  may be unreliable, especially when  $|R|$  is small. A better update formula for  $p_t$  is as follows:

$$p'_t = \frac{|R_t^+| + \alpha p_t}{|R^+| + \alpha},$$

where  $\alpha$  adjusts the contribution of the previous feedback cycle (e.g.,  $\alpha = |R|$ ).

New terms may be added to query  $q$  from relevant documents in  $R^+$ . A term weighting scheme  $\omega'(t)$  ranks terms, determining those that will maximally increase the difference in average score between relevant and non-relevant documents:

$$\omega'(t) = (P(\mathbf{d}(t) = 1 \mid r = 1) - P(\mathbf{d}(t) = 1 \mid r = 0)) \omega(t) \quad (1)$$

$$\approx P(\mathbf{d}(t) = 1 \mid r = 1) \omega(t) \quad (2)$$

$$\approx \frac{|R_t^+|}{|R^+|} \omega(t) \quad (3)$$

$$\stackrel{\text{rank}}{=} |R_t^+| \log \frac{p_t(1 - s_t)}{s_t(1 - p_t)} \quad (4)$$

## Remarks:

### ❑ Assumptions [\[Robertson 2009\]](#):

- (1) Difference in average term weights between relevant and irrelevant documents.
- (2) Omission of the second probability, which is usually much smaller than the first.
- (3) Maximum likelihood estimation of the first probability as proportion of relevant documents  $|R_t^+|$  containing term  $t$  among all relevant documents  $|R^+|$ .
- (4) Multiplication by  $|R^+|$ , which is a constant for a given query  $q$ 's result set  $R$ .

### ❑ The top- $k$ terms are added to $q$ ; say, $k = 10$ .



# Binary Independence Model

## Discussion

### Advantages:

- ❑ grounded in probabilistic theory
- ❑ performs well given some relevance feedback
- ❑ supplies theoretical justification of inverse document frequency

### Disadvantages:

- ❑ in the absence of relevance feedback, only about 50% recall compared to a *tf · idf*-based vector space model
- ❑ does not exploit term frequencies
- ❑ assumptions and rank-preserving simplifications do not generalize to retrieval scenarios other than ad hoc retrieval

# Okapi BM25

Retrieval Model  $\mathcal{R} = \langle \mathbf{D}, \mathbf{Q}, \rho \rangle$  [\[Generic Model\]](#) [\[Boolean Retrieval\]](#) [\[VSM\]](#) [\[BIM\]](#) [\[LSI\]](#) [\[ESA\]](#) [\[LM\]](#)

## Document representations $\mathbf{D}$ .

The set of index terms  $T = \{t_1, \dots, t_m\}$  is typically composed of the word stems of the vocabulary of a document collection, excluding stop words.

The representation  $\mathbf{d}$  of a document  $d$  is a  $|T|$ -dimensional vector, where the  $i$ -th vector component of  $\mathbf{d}$  corresponds to a term weight  $w_i$  of term  $t_i \in T$ , indicating its importance for  $d$ .

The term weighting scheme is the BM25 variant of the Best Match term weighting schemes.

## Query representations $\mathbf{Q}$ .

A query representation  $\mathbf{q}$  is constructed like a document representation.

## Relevance function $\rho$ .

Sum of the term weights computed for the words in a query.

# Okapi BM25

## Background

Empirical evidence suggests that term frequency is an important factor in determining the relevance of a document.

Relaxation of BIM to term frequencies within the 2-Poisson model:

- ❑ Change of the joint sample space to  $\Omega = \{0, 1\} \times \mathbb{N}^{|T|}$ , where  $\mathbb{N}^{|T|}$  denotes the set of document vectors with term frequency weights over the set of terms  $T$ .
- ❑ Starting point is the RSJ weight of the binary independence model:

$$P(r = 1 \mid \mathbf{d}, \mathbf{q}) \propto \omega^{\text{RSJ}} = \log \frac{p_t(1 - s_t)}{s_t(1 - p_t)}$$

- ❑ Estimation of  $p_t = P(\mathbf{d}(t) = tf(t, \mathbf{d}) \mid r = 1, \mathbf{q})$  as mixture of two Poisson distributions, distinguishing “elite” terms that occur unusually frequently from others. An elite term  $t$  encodes whether  $d$  is about the concept underlying  $t$ .
- ❑ Problems: Poisson distribution is a poor fit; too many parameters
- ❑ Approach: empirical approximation of the term weight  $\omega^{\text{RSJ}}$
- ❑ Resulted in the successful term weighting scheme Okapi BM25.

## Remarks:

- ❑ “Okapi” is the name of a retrieval system developed at City University London. “BM” stands for Best Match, and the number 25 refers to the best-performing variant tried. Other variants include BM0, BM1, BM11, and BM15. [Here](#) is an overview of all variants.
- ❑ We assume that each document is generated by filling a certain number of word-positions (fixed length) from a vocabulary of words. Furthermore, we assume a simple multinomial distribution over words, so that for each position each word has a fixed (small) probability of being chosen, independent of what other words have been chosen for other positions. Then it follows that the distribution of *tfs* [term frequencies] for a given word is binomial, which approximates to a Poisson under these conditions.

The eliteness model can be seen as a simple topical model which causes variation in the unigram distributions. The author is assumed first to choose which topics to cover, i.e., which terms to treat as elite and which not. This defines specific probabilities for the unigram model, and the author then fills the word-positions according to this chosen model.

This generative version of the 2-Poisson model (that is, a model for how documents are generated) ties it very closely with the [language models](#).

The model depends fairly crucially on the notion that all documents are of the same (fixed) length.

[[Robertson 2009](#)]

# Okapi BM25

## Term Weighting

$$\omega^{\text{BM25}}(t, d, D) = \omega^{\text{dtf}}(t, d) \cdot \omega^{\text{qtf}}(t, q) \cdot \omega^{\text{BIM}}(t, D)$$

$$\omega^{\text{dtf}}(t, d) = \frac{(k_1 + 1) \cdot \mathbf{tf}(t, d)}{k_1((1 - b) + b \cdot \frac{|d|}{|d|_{\text{avg}}}) + \mathbf{tf}(t, d)}$$

$$\omega^{\text{qtf}}(t, q) = \frac{(k_2 + 1) \cdot \mathbf{tf}(t, q)}{k_2 + \mathbf{tf}(t, q)}$$

$$\omega^{\text{BIM}}(t, D) = \log \frac{|D| - \mathbf{df}(t, D)}{\mathbf{df}(t, D)}$$

# Okapi BM25

## Term Weighting

$$\omega^{\text{BM25}}(t, d, D) = \omega^{\text{dtf}}(t, d) \cdot \omega^{\text{qtf}}(t, q) \cdot \omega^{\text{BIM}}(t, D)$$

$$\omega^{\text{dtf}}(t, d) = \frac{(k_1 + 1) \cdot \textit{tf}(t, d)}{k_1((1 - b) + b \cdot \frac{|d|}{|d|_{\text{avg}}}) + \textit{tf}(t, d)}$$

$$\omega^{\text{qtf}}(t, q) = \frac{(k_2 + 1) \cdot \textit{tf}(t, q)}{k_2 + \textit{tf}(t, q)}$$

$$\omega^{\text{BIM}}(t, D) = \log \frac{|D| - \textit{df}(t, D)}{\textit{df}(t, D)}$$

# Okapi BM25

## Term Weighting

$$\omega^{\text{BM25}}(t, d, D) = \omega^{\text{dtf}}(t, d) \cdot \omega^{\text{qtf}}(t, q) \cdot \omega^{\text{BIM}}(t, D)$$

$$\omega^{\text{dtf}}(t, d) = \frac{(k_1 + 1) \cdot \textcolor{brown}{tf}(t, d)}{k_1((1 - b) + b \cdot \frac{|d|}{|d|_{\text{avg}}}) + \textcolor{brown}{tf}(t, d)}$$

$$\omega^{\text{qtf}}(t, q) = \frac{(k_2 + 1) \cdot \textcolor{brown}{tf}(t, q)}{k_2 + \textcolor{brown}{tf}(t, q)}$$

$$\omega^{\text{BIM}}(t, D) = \log \frac{|D| - \textcolor{brown}{df}(t, D)}{\textcolor{brown}{df}(t, D)}$$

# Okapi BM25

## Term Weighting

$$\omega^{\text{BM25}}(t, d, D) = \omega^{\text{dtf}}(t, d) \cdot \omega^{\text{qtf}}(t, q) \cdot \omega^{\text{BIM}}(t, D)$$

$$\omega^{\text{dtf}}(t, d) = \frac{(k_1 + 1) \cdot \text{tf}(t, d)}{k_1((1 - b) + b \cdot \frac{|d|}{|d|_{\text{avg}}}) + \text{tf}(t, d)}$$

$$\omega^{\text{qtf}}(t, q) = \frac{(k_2 + 1) \cdot \text{tf}(t, q)}{k_2 + \text{tf}(t, q)}$$

$$\omega^{\text{BIM}}(t, D) = \log \frac{|D| - \text{df}(t, D)}{\text{df}(t, D)}$$

### Saturation:

- The eliteness of a term does not grow linearly with its frequency. This is represented by the cumulative distribution function of a Poisson distribution.
- Normalizing term frequency by  $k_1 + \text{tf}(t, d)$  yields a similar function for  $k_1 > 0$ .
- Multiplying by  $(k_1 + 1)$  ensures that the weights are  $\geq 1$ .



# Okapi BM25

## Term Weighting

$$\omega^{\text{BM25}}(t, d, D) = \omega^{\text{dtf}}(t, d) \cdot \omega^{\text{qtf}}(t, q) \cdot \omega^{\text{BIM}}(t, D)$$

$$\omega^{\text{dtf}}(t, d) = \frac{(k_1 + 1) \cdot \text{tf}(t, d)}{k_1((1 - b) + b \cdot \frac{|d|}{|d|_{\text{avg}}}) + \text{tf}(t, d)}$$

$$\omega^{\text{qtf}}(t, q) = \frac{(k_2 + 1) \cdot \text{tf}(t, q)}{k_2 + \text{tf}(t, q)}$$

$$\omega^{\text{BIM}}(t, D) = \log \frac{|D| - \text{df}(t, D)}{\text{df}(t, D)}$$

Document length:

- ❑ Why authors increase document length: verbosity or scope. The former suggests normalization, the latter not; **real documents are mixtures**.
- ❑ Normalization by the average document length  $|d|_{\text{avg}}$  ensures independence of datasets and implementation details.

# Okapi BM25

## Term Weighting

$$\omega^{\text{BM25}}(t, d, D) = \omega^{\text{dtf}}(t, d) \cdot \omega^{\text{qtf}}(t, q) \cdot \omega^{\text{BIM}}(t, D)$$

$$\omega^{\text{dtf}}(t, d) = \frac{(k_1 + 1) \cdot \mathbf{tf}(t, d)}{k_1((1 - b) + b \cdot \frac{|d|}{|d|_{\text{avg}}}) + \mathbf{tf}(t, d)}$$

$$\omega^{\text{qtf}}(t, q) = \frac{(k_2 + 1) \cdot \mathbf{tf}(t, q)}{k_2 + \mathbf{tf}(t, q)}$$

$$\omega^{\text{BIM}}(t, D) = \log \frac{|D| - \mathbf{df}(t, D)}{\mathbf{df}(t, D)}$$

Term frequency weighting for query:

- ❑ Like document term frequency weighting, albeit without length normalization.
- ❑ Can be omitted if queries are generally short.
- ❑ Useful for query by example scenarios.

# Okapi BM25

## Term Weighting

$$\omega^{\text{BM25}}(t, d, D) = \omega^{\text{dtf}}(t, d) \cdot \omega^{\text{qtf}}(t, q) \cdot \omega^{\text{BIM}}(t, D)$$

$$\omega^{\text{dtf}}(t, d) = \frac{(k_1 + 1) \cdot \mathbf{tf}(t, d)}{k_1((1 - b) + b \cdot \frac{|d|}{|d|_{\text{avg}}}) + \mathbf{tf}(t, d)}$$

$$\omega^{\text{qtf}}(t, q) = \frac{(k_2 + 1) \cdot \mathbf{tf}(t, q)}{k_2 + \mathbf{tf}(t, q)}$$

$$\omega^{\text{BIM}}(t, D) = \log \frac{|D| - \mathbf{df}(t, D)}{\mathbf{df}(t, D)}$$

### Parameters:

- ❑  $k_1$  must be optimized against  $D$ ;  $k_1 = 1.2$  is a good value to start with.
- ❑  $k_2$  must be optimized against  $Q$ ; in practice  $0 \leq k_2 \leq 1000$ , the shorter the queries, the less sensitive the overall weight is to  $k_2$ .
- ❑  $b$  must be optimized against  $D$ ;  $b = 0.75$  is a good value to start with.

# Okapi BM25

## Term Weighting

$$\omega^{\text{BM25}}(t, d, D) = \omega^{\text{dtf}}(t, d) \cdot \omega^{\text{qtf}}(t, q) \cdot \omega^{\text{BIM}}(t, D)$$

$$\omega^{\text{dtf}}(t, d) = \frac{(k_1 + 1) \cdot \mathbf{tf}(t, d)}{k_1((1 - b) + b \cdot \frac{|d|}{|d|_{\text{avg}}}) + \mathbf{tf}(t, d)}$$

$$\omega^{\text{qtf}}(t, q) = \frac{(k_2 + 1) \cdot \mathbf{tf}(t, q)}{k_2 + \mathbf{tf}(t, q)}$$

$$\omega^{\text{BIM}}(t, D) = \log \frac{|D| - \mathbf{df}(t, D)}{\mathbf{df}(t, D)}$$

## Extension: BM25F (simple)

If documents have fields of varying importance, they can be weighted as follows:

$$\mathbf{tf}'(t, d) = \sum_{s \in d} k_s \cdot \mathbf{tf}(t, s); \quad |d|' = \sum_{s \in d} k_s \sum_{t \in s} \mathbf{tf}(t, d), \quad |d|'_{\text{avg}} = \frac{1}{|D|} \sum_{d \in D} |d|'$$

where each  $s$  denotes a field of document  $d$ , and  $k_s$  the field-specific weight.

# Okapi BM25

## Term Weighting

$$\omega^{\text{BM25}}(t, d, D) = \omega^{\text{dtf}}(t, d) \cdot \omega^{\text{qtf}}(t, q) \cdot \omega^{\text{BIM}}(t, D)$$

$$\omega^{\text{dtf}}(t, d) = \frac{(k_1 + 1) \cdot \mathbf{tf}(t, d)}{k_1((1 - b) + b \cdot \frac{|d|}{|d|_{\text{avg}}}) + \mathbf{tf}(t, d)}$$

$$\omega^{\text{qtf}}(t, q) = \frac{(k_2 + 1) \cdot \mathbf{tf}(t, q)}{k_2 + \mathbf{tf}(t, q)}$$

$$\omega^{\text{BIM}}(t, D) = \log \frac{|D| - \mathbf{df}(t, D)}{\mathbf{df}(t, D)}$$

## Relevance Function $\rho$

$$\rho(\mathbf{d}, \mathbf{q}) = \sum_{t \in \mathbf{q}} \omega^{\text{BM25}}(t, d, D),$$

where  $D$  is the document collection indexed.

# Okapi BM25

## Discussion

### Advantages:

- ❑ Very good retrieval performance
- ❑ Well tunable to different retrieval scenarios
- ❑ Most terms can be precomputed at indexing time

### Disadvantages:

- ❑ Departure from the theoretic probabilistic foundation → an empirical model