

Introduction to **Information Retrieval**

Probabilistic Information Retrieval

Chris Manning, Pandu Nayak and Prabhakar Raghavan

further edited by Razvan Bunescu

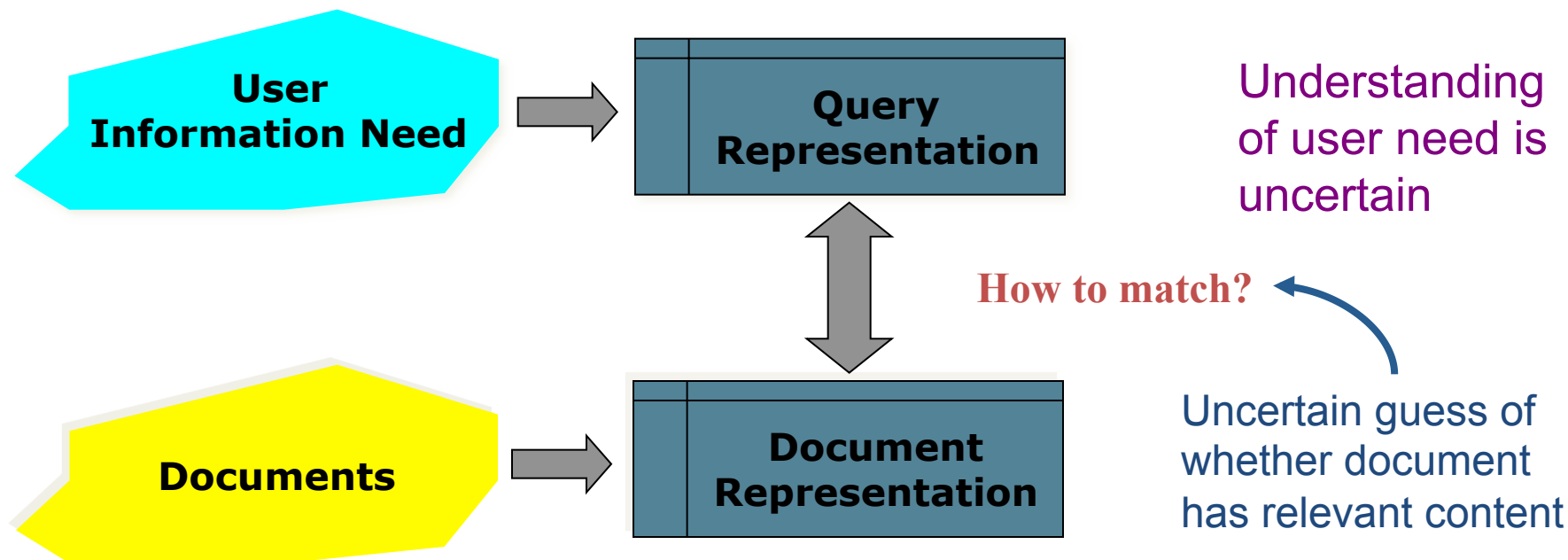
Summary – vector space ranking

- Represent the query as a weighted tf-idf vector
- Represent each document as a weighted tf-idf vector
- Compute the cosine similarity score for the query vector and each document vector
- Rank documents with respect to the query by score
- Return the top K (e.g., $K = 10$) to the user

tf-idf weighting has many variants

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

Why probabilities in IR?



In traditional IR systems, matching between each document and query is attempted in a semantically imprecise space of index terms.

Probabilities provide a principled foundation for uncertain reasoning.
Can we use probabilities to quantify our uncertainties?

Probabilistic IR topics

- Classical probabilistic retrieval model
 - Probability ranking principle, etc.
 - Binary independence model (\approx Naïve Bayes text cat)
 - (Okapi) BM25
- Bayesian networks for text retrieval
- Language model approach to IR
 - An important emphasis in recent work
- *Probabilistic methods are one of the oldest but also one of the currently hottest topics in IR.*
 - *Traditionally: neat ideas, but didn't win on performance*
 - *It may be different now.*

The document ranking problem

- We have a collection of documents
- User issues a query
- A list of documents needs to be returned
- **Ranking method is the core of an IR system:**
 - **In what order do we present documents to the user?**
 - We want the “best” document to be first, second best second, etc....
- **Idea: Rank by probability of relevance of the document w.r.t. information need**
 - $P(R=1 | \text{document}_i, \text{query})$

Recall a few probability basics

- For events A and B :
- Bayes' Rule

$$p(A, B) = p(A \cap B) = p(A | B)p(B) = p(B | A)p(A)$$

$$p(A | B) = \frac{p(B | A)p(A)}{p(B)} = \frac{p(B | A)p(A)}{\sum_{X=A, \bar{A}} p(B | X)p(X)}$$

↑ Posterior
 ↑ Prior

- Odds: $O(A) = \frac{p(A)}{p(\bar{A})} = \frac{p(A)}{1 - p(A)}$

The 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.”

- [1960s/1970s] S. Robertson, W.S. Cooper, M.E. Maron; van Rijsbergen (1979:113); Manning & Schütze (1999:538)

Probability Ranking Principle

Let x represent a document in the collection.

Let R represent **relevance** of a document w.r.t. given (fixed) query and let $\mathbf{R=1}$ represent relevant and $\mathbf{R=0}$ not relevant.

Need to find $p(R=1|x)$ - probability that a document x is **relevant**.

$$p(R = 1 | x) = \frac{p(x | R = 1)p(R = 1)}{p(x)}$$

$p(R=1), p(R=0)$ - prior probability of retrieving a relevant or non-relevant document

$$p(R = 0 | x) = \frac{p(x | R = 0)p(R = 0)}{p(x)}$$

$p(x|R=1), p(x|R=0)$ - probability that if a relevant (not relevant) document is retrieved, it is x .

$$p(R = 0 | x) + p(R = 1 | x) = 1$$

Probability Ranking Principle (PRP)

- Simple case: no selection costs or other utility concerns that would differentially weight errors
- PRP in action: Rank all documents by $p(R=1 | x)$
- Theorem: Using the PRP is optimal, in that it minimizes the loss (Bayes risk) under 1/0 loss
 - Provable if all probabilities correct, etc. [e.g., Ripley 1996]

Probability Ranking Principle

- More complex case: retrieval costs.
 - Let d be a document
 - C – cost of not retrieving a relevant document
 - C' – cost of retrieving a non-relevant document

- Probability Ranking Principle: if

$$C' \cdot p(R = 0 | d) - C \cdot p(R = 1 | d) \leq C' \cdot p(R = 0 | d') - C \cdot p(R = 1 | d')$$

for all d' *not yet retrieved*, then d **is the next document to be retrieved**.

- **We won't further consider cost/utility from now on.**

Probability Ranking Principle

- How do we compute all those probabilities?
 - Do not know exact probabilities, have to use estimates.
 - Binary Independence Model (BIM) – which we discuss next – is the simplest model.
- Questionable assumptions
 - “Relevance” of each document is independent of relevance of other documents.
 - Really, it’s bad to keep on returning **duplicates**
 - Boolean model of relevance.
 - That one has a single step information need:
 - Seeing a range of results might let user refine query.

Probabilistic Retrieval Strategy

- Estimate how terms contribute to relevance:
 - How do things like *tf*, *df*, and *document length* influence your judgments about document relevance?
 - A more nuanced answer is the Okapi formulae [Jones & Robertson].
- Combine to find document relevance probability.
- Order documents by decreasing probability.

Basic concept of probabilistic ranking:

“For a given query, if we know some documents that are relevant, terms that occur in those documents should be given greater weighting in searching for other relevant documents.

By making assumptions about the distribution of terms and applying Bayes Theorem, it is possible to derive weights theoretically.”

Binary Independence Model

- Traditionally used in conjunction with PRP.
- **“Binary” = Boolean**: documents are represented as binary incidence vectors of terms (cf. IIR Chapter 1):
 - $\vec{x} = (x_1, \dots, x_n)$
 - $x_i = 1$ iff term i is present in document x .
- **“Independence”**: terms occur in documents independently
- Different documents can be modeled as the same vector

Binary Independence Model

- Queries: binary term incidence vectors
- Given query q ,
 - for each document d need to compute $p(R | q, d)$.
 - replace with computing $p(R | q, x)$ where x is binary term incidence vector representing d .
 - Interested only in ranking
- Will use odds and Bayes' Rule:

$$O(R | q, \vec{x}) = \frac{p(R = 1 | q, \vec{x})}{p(R = 0 | q, \vec{x})} = \frac{\frac{p(R = 1 | q)p(\vec{x} | R = 1, q)}{p(\vec{x} | q)}}{\frac{p(R = 0 | q)p(\vec{x} | R = 0, q)}{p(\vec{x} | q)}}$$

Binary Independence Model

$$O(R | q, \vec{x}) = \frac{p(R = 1 | q, \vec{x})}{p(R = 0 | q, \vec{x})} = \frac{p(R = 1 | q)}{p(R = 0 | q)} \cdot \frac{p(\vec{x} | R = 1, q)}{p(\vec{x} | R = 0, q)}$$

Constant for a
given query

Needs estimation

- Using **Independence** Assumption:

$$\frac{p(\vec{x} | R = 1, q)}{p(\vec{x} | R = 0, q)} = \prod_{i=1}^n \frac{p(x_i | R = 1, q)}{p(x_i | R = 0, q)}$$

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{i=1}^n \frac{p(x_i | R = 1, q)}{p(x_i | R = 0, q)}$$

Binary Independence Model

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{i=1}^n \frac{p(x_i | R = 1, q)}{p(x_i | R = 0, q)}$$

- Since x_i is either 0 or 1:

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{x_i=1} \frac{p(x_i = 1 | R = 1, q)}{p(x_i = 1 | R = 0, q)} \cdot \prod_{x_i=0} \frac{p(x_i = 0 | R = 1, q)}{p(x_i = 0 | R = 0, q)}$$

- Let $p_i = p(x_i = 1 | R = 1, q)$; $r_i = p(x_i = 1 | R = 0, q)$;

	document	relevant (R=1)	not relevant (R=0)
term present	$x_i = 1$	p_i	r_i
term absent	$x_i = 0$	$(1 - p_i)$	$(1 - r_i)$

Binary Independence Model

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{x_i=1} \frac{p(x_i = 1 | R = 1, q)}{p(x_i = 1 | R = 0, q)} \cdot \prod_{x_i=0} \frac{p(x_i = 0 | R = 1, q)}{p(x_i = 0 | R = 0, q)}$$

$$p_i = p(x_i = 1 | R = 1, q); \quad r_i = p(x_i = 1 | R = 0, q);$$

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{x_i=1} \frac{p_i}{r_i} \cdot \prod_{x_i=0} \frac{1 - p_i}{1 - r_i}$$

- Assume $p_i = r_i$ for terms not occurring in the query ($q_i=0$):

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{\substack{x_i=1 \\ q_i=1}} \frac{p_i}{r_i} \cdot \prod_{\substack{x_i=0 \\ q_i=1}} \frac{(1 - p_i)}{(1 - r_i)}$$

Binary Independence Model

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{x_i=q_i=1} \frac{p_i}{r_i} \cdot \prod_{\substack{x_i=0 \\ q_i=1}} \frac{1-p_i}{1-r_i}$$

All matching terms
Non-matching query terms

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{\substack{x_i=1 \\ q_i=1}} \frac{p_i}{r_i} \cdot \prod_{\substack{x_i=1 \\ q_i=1}} \left(\frac{1-r_i}{1-p_i} \cdot \frac{1-p_i}{1-r_i} \right) \prod_{\substack{x_i=0 \\ q_i=1}} \frac{1-p_i}{1-r_i}$$

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{x_i=q_i=1} \frac{p_i(1-r_i)}{r_i(1-p_i)} \cdot \prod_{q_i=1} \frac{1-p_i}{1-r_i}$$

All matching terms
All query terms

Binary Independence Model

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{x_i=q_i=1} \frac{p_i(1-r_i)}{r_i(1-p_i)} \cdot \prod_{q_i=1} \frac{1-p_i}{1-r_i}$$

Constant for each query

Only quantity to be estimated for rankings

Retrieval Status Value:

$$RSV = \log \prod_{x_i=q_i=1} \frac{p_i(1-r_i)}{r_i(1-p_i)} = \sum_{x_i=q_i=1} \log \frac{p_i(1-r_i)}{r_i(1-p_i)}$$

Binary Independence Model

All boils down to computing RSV.

$$RSV = \log \prod_{x_i=q_i=1} \frac{p_i(1-r_i)}{r_i(1-p_i)} = \sum_{x_i=q_i=1} \log \frac{p_i(1-r_i)}{r_i(1-p_i)}$$

$$RSV = \sum_{x_i=q_i=1} c_i; \quad c_i = \log \frac{p_i(1-r_i)}{r_i(1-p_i)} = \log \frac{p_i}{1-p_i} + \log \frac{1-r_i}{r_i}$$

The c_i are log odds ratios. They function as the term weights in this model

So, how do we compute c_i 's from data ?

Binary Independence Model

- Estimating RSV coefficients in theory
- For each term i look at this table of document counts:

Documents	Relevant	Non-Relevant	Total
$x_i = 1$	s	$df_i - s$	df_i
$x_i = 0$	$S - s$	$N - df_i - S + s$	$N - df_i$
Total	S	$N - S$	N

- Estimates: $p_i \approx \frac{s}{S}$ $r_i \approx \frac{(df_i - s)}{(N - S)}$

$$c_i \approx K(N, df_i, S, s) = \log \frac{s/(S - s)}{(df_i - s)/(N - df_i - S + s)}$$

For now, assume no zero terms. See later lecture.

Estimation – key challenge

- Assumption:
 - non-relevant documents are approximated by the whole collection \Leftrightarrow relevant documents are a very small percentage of the collection
- Then:
 - $df_i - s \approx df_i$ and $N - S \approx N$ and thus $r_i \approx df_i / N$

$$\log \frac{1-r_i}{r_i} \approx \log \frac{N-df_i}{df_i} \approx \log \frac{N}{df_i}$$

IDF!

?

$$RSV = \sum_{x_i=q_i=1} c_i \quad \text{where} \quad c_i = \log \frac{p_i(1-r_i)}{r_i(1-p_i)} = \log \frac{p_i}{1-p_i} + \log \frac{1-r_i}{r_i}$$

Estimation – key challenge

- p_i (probability of occurrence in relevant documents) cannot be approximated as easily
- p_i can be estimated in various ways:
 - from relevant documents if know some:
 - Relevance weighting can be used in a feedback loop.
 - constant (Croft and Harper combination match) – then just get idf weighting of terms (with $p_i=0.5$)

$$RSV = \sum_{x_i=q_i=1} \log \frac{N}{n_i}$$

- proportional to prob. of occurrence in collection
 - Greiff (SIGIR 1998) argues for $1/3 + 2/3 df_i/N$

Probabilistic Relevance Feedback

1. Start with preliminary estimates of relevant probabilities:
 - Use estimates of p_i and r_i from previous section.
2. Use p_i and r_i to retrieve a set V of documents.
3. Interact with the user to refine the estimates:
 - V is partitioned into VR (relevant docs) and VNR (nonrelevant docs).
4. Reestimate p_i and r_i on the basis of these:
 - Use ML estimates (if counts are large enough):
 - $p_i = |VR_i|/|VR|$, or smoothed $p_i = (|VR_i| + 0.5) / (|VR| + 1)$
 - Or combine new information with original guess:

$$p_i^{(k+1)} = \frac{|VR_i| + \kappa p_i^{(k)}}{|VR| + \kappa}$$

Bayesian updating:

- κ is prior weight.
- $\kappa = 0.5$.

5. Repeat from step 2 until user is satisfied.

Pseudo-relevance feedback: $VR = V$

1. Assume initial estimates for p_i and r_i as before.
2. Determine guess of relevant document set:
 - If unsure, a (too) small guess is likely to be best.
 - V is fixed size set of highest ranked documents
3. We need to improve our guesses for p_i and r_i , so:
 - Let V_i be set of documents containing x_i
 - $p_i = (|V_i| + 0.5) / (|V| + 1)$
 - Assume if not retrieved then not relevant
 - $r_i = (df_i - |V_i| + 0.5) / (N - |V| + 1)$
4. Go to step 2. until convergence, then return ranking.

RSV weights vs. VSM weights

$$c_i = \log \frac{p_i(1-r_i)}{r_i(1-p_i)} = \log \frac{p_i}{1-p_i} + \log \frac{1-r_i}{r_i} \approx \log \frac{|V_i|+0.5}{|V|-|V_i|+1} + \log \frac{N}{df_i}$$

The diagram illustrates the relationship between RSV and VSM weights. A central green box containing a question mark has two arrows pointing to the first two log terms of the RSV equation: $\log \frac{p_i}{1-p_i}$ and $\log \frac{1-r_i}{r_i}$. A green box labeled 'IDF' has an arrow pointing to the third log term of the RSV equation: $\log \frac{N}{df_i}$.

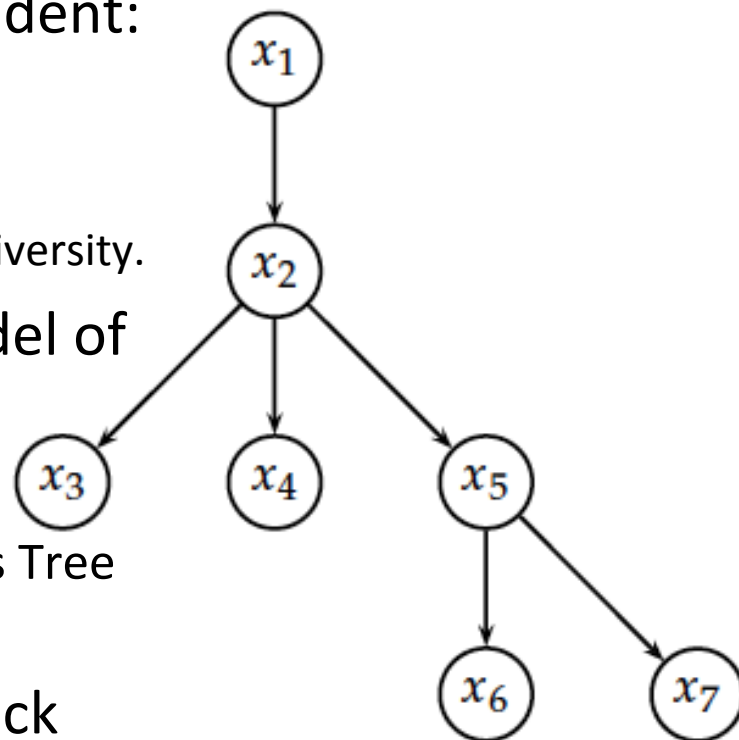
- But things are not quite the same:
 - $p_i / (1 - p_i)$ does not measure term frequency.
 - What does it measure?
 - The two log scaled components are added in RSV, not multiplied as in VSM.

PRP and BIM

- Getting reasonable approximations of probabilities is possible.
- Requires restrictive assumptions:
 - *Term independence*
 - *Terms not in query don't affect the outcome*
 - *Boolean representation of documents/queries/relevance*
 - *Document relevance values are independent*
- Some of these assumptions can be removed.
- Problem: either require partial relevance information or only can derive somewhat inferior term weights.

Removing term independence

- In general, index terms aren't independent:
 - New vs. York.
- Dependencies can be complex:
 - New, York, England, City, Stock, Exchange, University.
- van Rijsbergen (1979) proposed a model of simple tree dependencies:
 - Each term dependent on another.
 - Reinvented by Friedman and Goldszmidt's Tree Augmented Naive Bayes (AAAI 13, 1996).
- In 1970s, estimation problems held back success of this model.



A key limitation of the BIM

- BIM – like much of original IR – was designed for titles or abstracts, and not for modern full text search.
- We want to pay attention to *term frequency* and *document length*, just like in other models we've discussed.
- Goal: be sensitive to these quantities while not adding too many parameters:
 - (Robertson and Zaragoza 2009; Spärck Jones et al. 2000)

Okapi BM25: An Extension to BIM

- BM25 – “Best Match 25” (they had a bunch of tries!):
 - Developed in the context of the Okapi system.
 - Started to be increasingly adopted by other teams during the TREC competitions.
 - It has been used widely and quite successfully across a range of collections and search tasks.
- Goal: be sensitive to these quantities while not adding too many parameters
 - (Robertson and Zaragoza 2009; Spärck Jones et al. 2000)

Retrieval Status Value in BM25

- Similar to the BIM derivation, we have:

$$RSV_d = \sum_{x_i=q_i=1} c_i(tf_{id})$$

- Qualitative properties of $c_i()$:

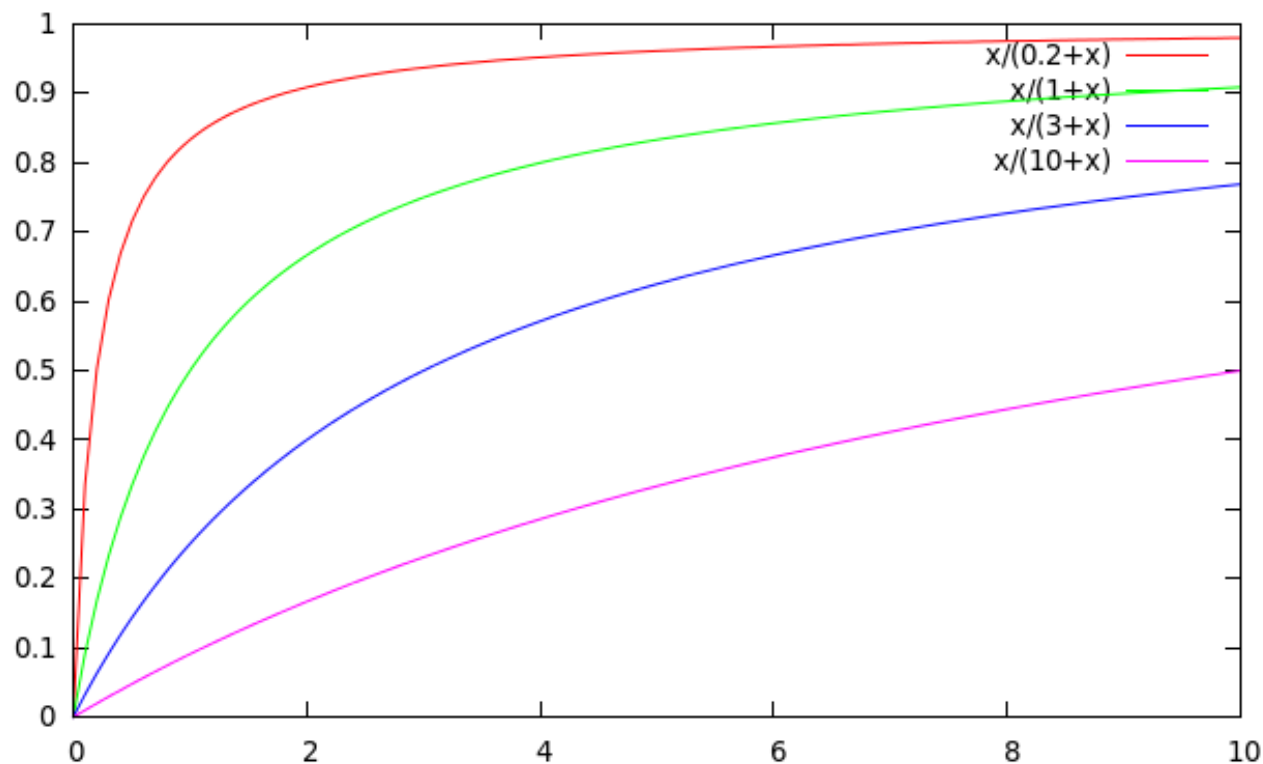
- $c_i(tf_{id})$ increases monotonically with tf_{id}

- $c_i(0)=0$

- $\lim_{tf_{id} \rightarrow \infty} c_i(tf_{id}) = c_i^{BIM}$

=> use the simple parametric curve $c_i(tf) = \frac{tf}{k_1 + tf}$

Saturation function



- For high values of k_1 , increments in tf_i continue to contribute significantly to the score.
- Contributions tail off quickly for low values of k_1

“Early” versions of BM25

- Version 1: using the saturation function:

$$c_i^{BM25v1}(tf_i) = c_i^{BIM} \frac{tf_i}{k_1 + tf_i}$$

- Version 2: BIM simplification to IDF:

$$c_i^{BM25v2}(tf_i) = \log \frac{N}{df_i} \times \frac{(k_1 + 1)tf_i}{k_1 + tf_i}$$

- $(k_1 + 1)$ factor doesn't change ranking, but makes term score 1 when $tf_i = 1$
- Similar to *tf-idf*, but terms scores are bounded.

Document length normalization

- Longer documents are likely to have larger tf_i values.
- Why might documents be longer?
 - Verbosity: suggests observed tf_i too high.
 - Larger scope: suggests observed tf_i may be right.
- A real document collection probably has both effects.
- ... so should apply some kind of normalization

Document length normalization

- Document length:

$$L_d = \sum_{i \in V} tf_{id}$$

- L_{ave} = average document length over collection
- Length normalization component:

$$B = \left((1 - b) + b \frac{L_d}{L_{ave}} \right), \quad 0 \leq b \leq 1$$

- $b = 1 \Rightarrow$ full document length normalization
- $b = 0 \Rightarrow$ no document length normalization

Okapi BM25

- Normalize tf using document length:

$$tf'_i = \frac{tf_i}{B}$$

$$\begin{aligned} c_i^{BM25}(tf_i) &= \log \frac{N}{df_i} \times \frac{(k_1 + 1)tf'_i}{k_1 + tf'_i} \\ &= \log \frac{N}{df_i} \times \frac{(k_1 + 1)tf_i}{k_1((1 - b) + b \frac{L_d}{L_{ave}}) + tf_i} \end{aligned}$$

- BM25 ranking function:

$$RSV_d^{BM25} = \sum_{i \in q} c_i^{BM25}(tf_{id})$$

Okapi BM25

$$RSV_d^{BM25} = \sum_{i \in q} \log \frac{N}{df_i} \cdot \frac{(k_1 + 1)tf_{id}}{k_1 \left((1 - b) + b \frac{L_d}{L_{ave}} \right) + tf_{id}}$$

- k_1 controls term frequency scaling:
 - $k_1 = 0$ is binary model; $k_1 = \text{large}$ is raw term frequency.
- b controls document length normalization:
 - $b = 0$ is no length normalization; $b = 1$ is relative frequency (fully scale by document length).
- Typically, k_1 is set around 1.2–2 and b around 0.75.
- *IIR* sec. 11.4.3 discusses incorporating query term weighting and (pseudo) relevance feedback.

Okapi BM25

- If the query is long (e.g. paragraph), we can use similar weighting for query terms:

$$RSV_d^{BM25} = \sum_{i \in q} \log \frac{N}{df_i} \cdot \frac{(k_1 + 1)tf_{id}}{k_1 \left((1 - b) + b \frac{L_d}{L_{ave}} \right) + tf_{id}} \cdot \frac{(k_2 + 1)tf_{iq}}{k_2 + tf_{iq}}$$

- Length normalization of the query is unnecessary.
- Parameters k_1 , k_3 , and b tuned on separate development collection:
 - or simply set k_1 , k_3 , to values in $[1.2, 2]$, and b around 0.75.

References

- S. E. Robertson and K. Spärck Jones. 1976. Relevance Weighting of Search Terms. *Journal of the American Society for Information Sciences* 27(3): 129–146.
- C. J. van Rijsbergen. 1979. *Information Retrieval*. 2nd ed. London: Butterworths, chapter 6. [Most details of math] <http://www.dcs.gla.ac.uk/Keith/Preface.html>
- N. Fuhr. 1992. Probabilistic Models in Information Retrieval. *The Computer Journal*, 35(3), 243–255. [Easiest read, with BNs]
- S. E. Robertson and H. Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Foundations and Trends in Information Retrieval* 3(4): 333-389.
- K. Spärck Jones, S. Walker, and S. E. Robertson. 2000. A probabilistic model of information retrieval: Development and comparative experiments. Part 1. *Information Processing and Management* 779–808.