

Introduction to

# **CS60092: Information Retrieval**

Sourangshu Bhattacharya

# 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**
  - $P(R=1 \mid \text{document}_i, \text{query})$

# Recall a few probability basics

- For events  $A$  and  $B$ :
- Bayes' Rule

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

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

Prior

Posterior

- 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  $R=1$  represent relevant and  $R=0$  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
      - Spärck Jones / Robertson
- Combine to find document relevance probability
- Order documents by decreasing probability

# Probabilistic Ranking

---

## **Basic concept:**

“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.”

*Van Rijsbergen*

# 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  $p_i = p(x_i = 1 | R = 1, q)$ ;  $r_i = p(x_i = 1 | R = 0, q)$ ;
- Assume, for all terms not occurring in the query ( $q_i = 0$ )  $p_i = r_i$

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

# Binary Independence Model

---

	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

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

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.

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

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

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 our 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
$x_i=1$	$s$	$n-s$	$n$
$x_i=0$	$S-s$	$N-n-S+s$	$N-n$
Total	$S$	$N-S$	$N$

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

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

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

# Estimation – key challenge

---

- If non-relevant documents are approximated by the whole collection, then  $r_i$  (prob. of occurrence in non-relevant documents for query) is  $n/N$  and

$$\log \frac{1-r_i}{r_i} = \log \frac{N-n-S+s}{n-s} \approx \log \frac{N-n}{n} \approx \log \frac{N}{n} = \text{IDF!}$$

# 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. Guess a preliminary probabilistic description of  $R=1$  documents and use it to retrieve a first set of documents
2. Interact with the user to refine the description: learn some definite members with  $R=1$  and  $R=0$
3. Reestimate  $p_i$  and  $r_i$  on the basis of these
  - Or can combine new information with original guess (use Bayesian prior):
4. Repeat, thus generating a succession of approximations to relevant documents

$$p_i^{(2)} = \frac{|V_i| + \kappa p_i^{(1)}}{|V| + \kappa}$$

$\kappa$  is  
prior  
weight

# Iteratively estimating $p_i$ and $r_i$ (= Pseudo-relevance feedback)

---

1. Assume that  $p_i$  is constant over all  $x_i$  in query and  $r_i$  as before
  - $p_i = 0.5$  (even odds) for any given doc
2. Determine guess of relevant document set:
  - $V$  is fixed size set of highest ranked documents on this model
3. We need to improve our guesses for  $p_i$  and  $r_i$ , so
  - Use distribution of  $x_i$  in docs in  $V$ . Let  $V_i$  be set of documents containing  $x_i$ 
    - $p_i = |V_i| / |V|$
  - Assume if not retrieved then not relevant
    - $r_i = (n_i - |V_i|) / (N - |V|)$
4. Go to 2. until converges then return ranking

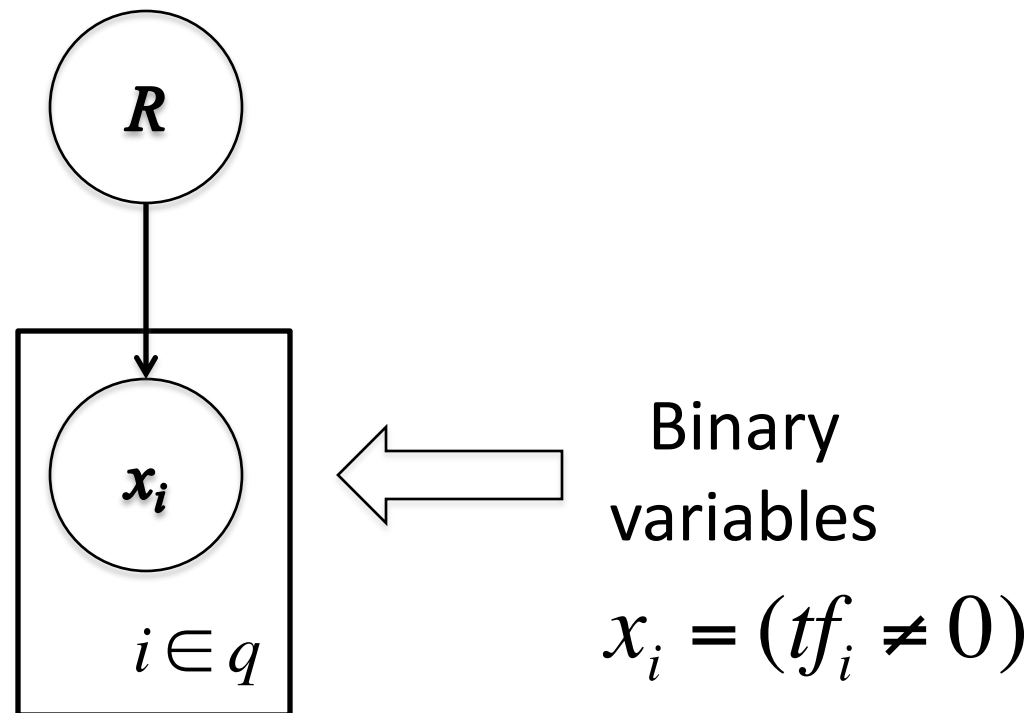


# PRP and BIM

---

- 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

# Graphical model for BIM



# Removing term independence

- In general, index terms aren't independent
- Dependencies can be complex
- van Rijsbergen (1979) proposed model of simple tree dependencies
  - Exactly Friedman and Goldszmidt's Tree Augmented Naive Bayes (AAAI 13, 1996)
- Each term dependent on one other
- In 1970s, estimation problems held back success of this model

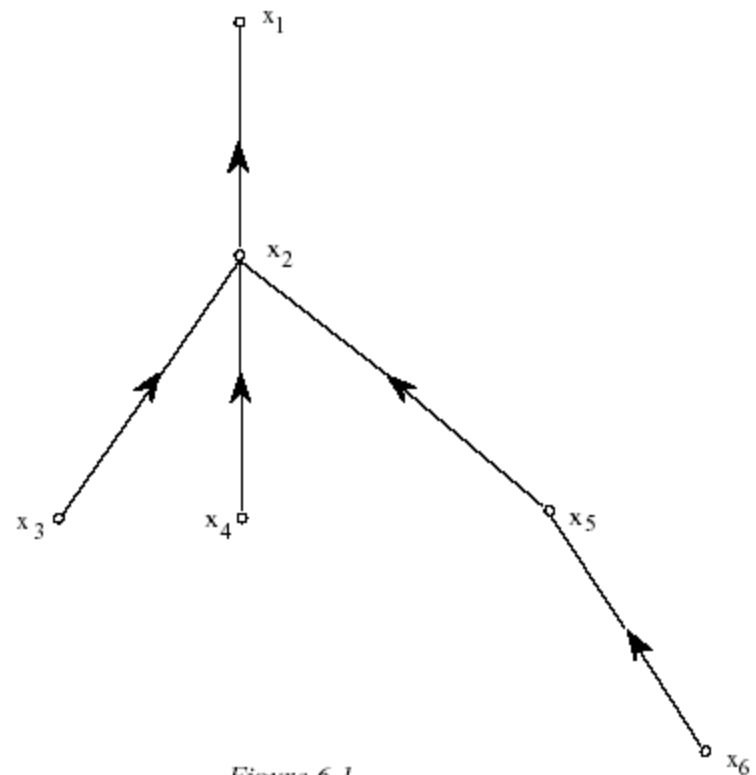


Figure 6.1.

# Summary - BIM

- Boils down to

$$RSV^{BIM} = \sum_{x_i=q_i=1} c_i^{BIM}; \quad c_i^{BIM} = \log \frac{p_i(1-r_i)}{(1-p_i)r_i} \quad \leftarrow \text{Log odds ratio}$$

where

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

- Simplifies to (with constant  $p_i = 0.5$ )

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

# Okapi BM25: A Non-binary Model

---

- The BIM was originally designed for short catalog records of fairly consistent length, and it works reasonably in these contexts
- For modern full-text search collections, a model should pay attention to term frequency and document length
- BestMatch25 (a.k.a **BM25** or **Okapi**) is sensitive to these quantities
- From 1994 until today, BM25 is one of the most widely used and robust retrieval models

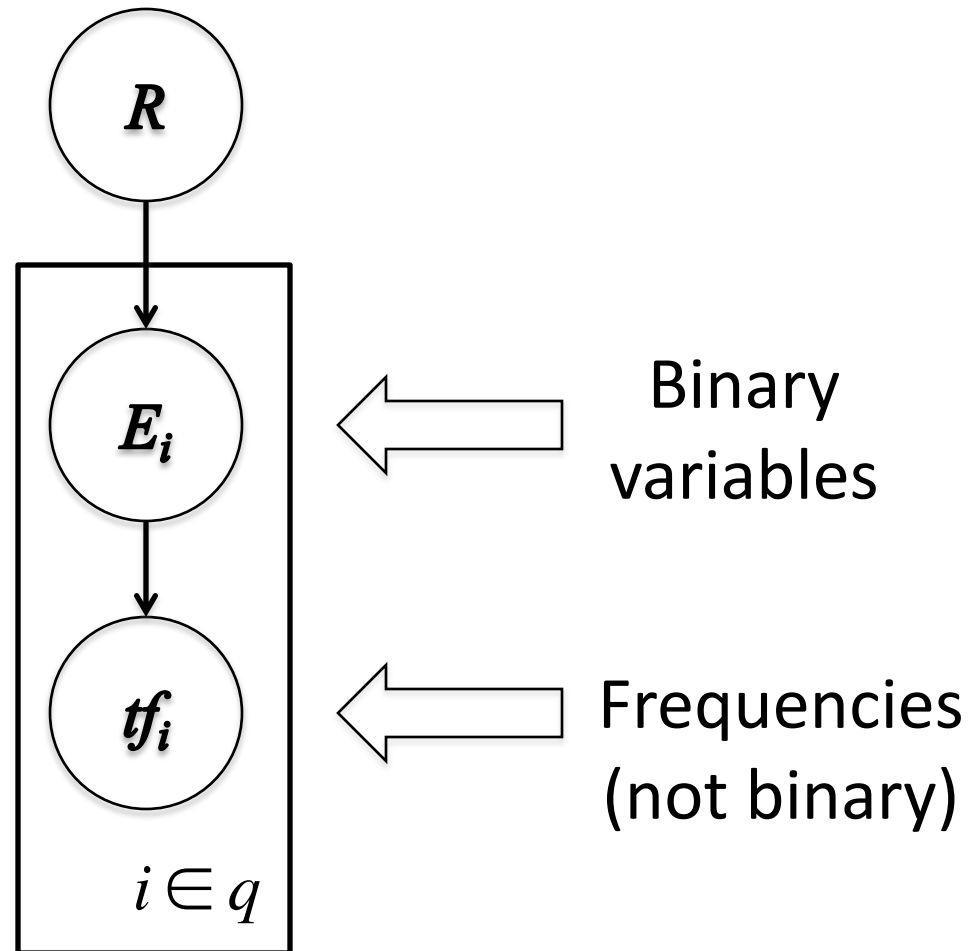
# Elite terms

---

Text from the Wikipedia page on the NFL draft showing **elite terms**

The **National Football League Draft** is an annual event in which the **National Football League (NFL) teams select eligible college football players**. It serves as the **league's** most common source of **player recruitment**. The basic design of the **draft** is that each **team** is given a **position** in the **draft order** in **reverse order** relative to its **record** ...

# Graphical model with eliteness



# Approximation of RSV

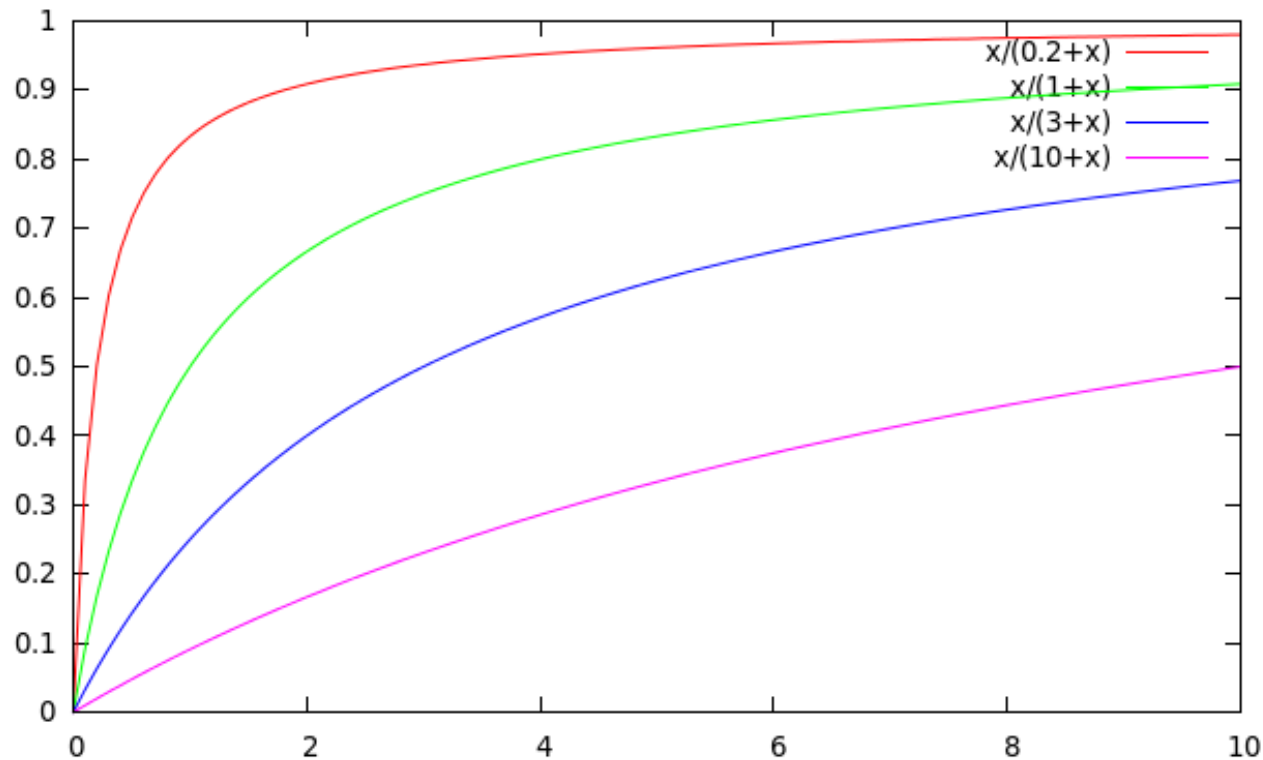
---

- The correction involving  $tf$ :

$$\textit{correction} = \frac{tf}{k_1 + tf}$$



# 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 term 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:

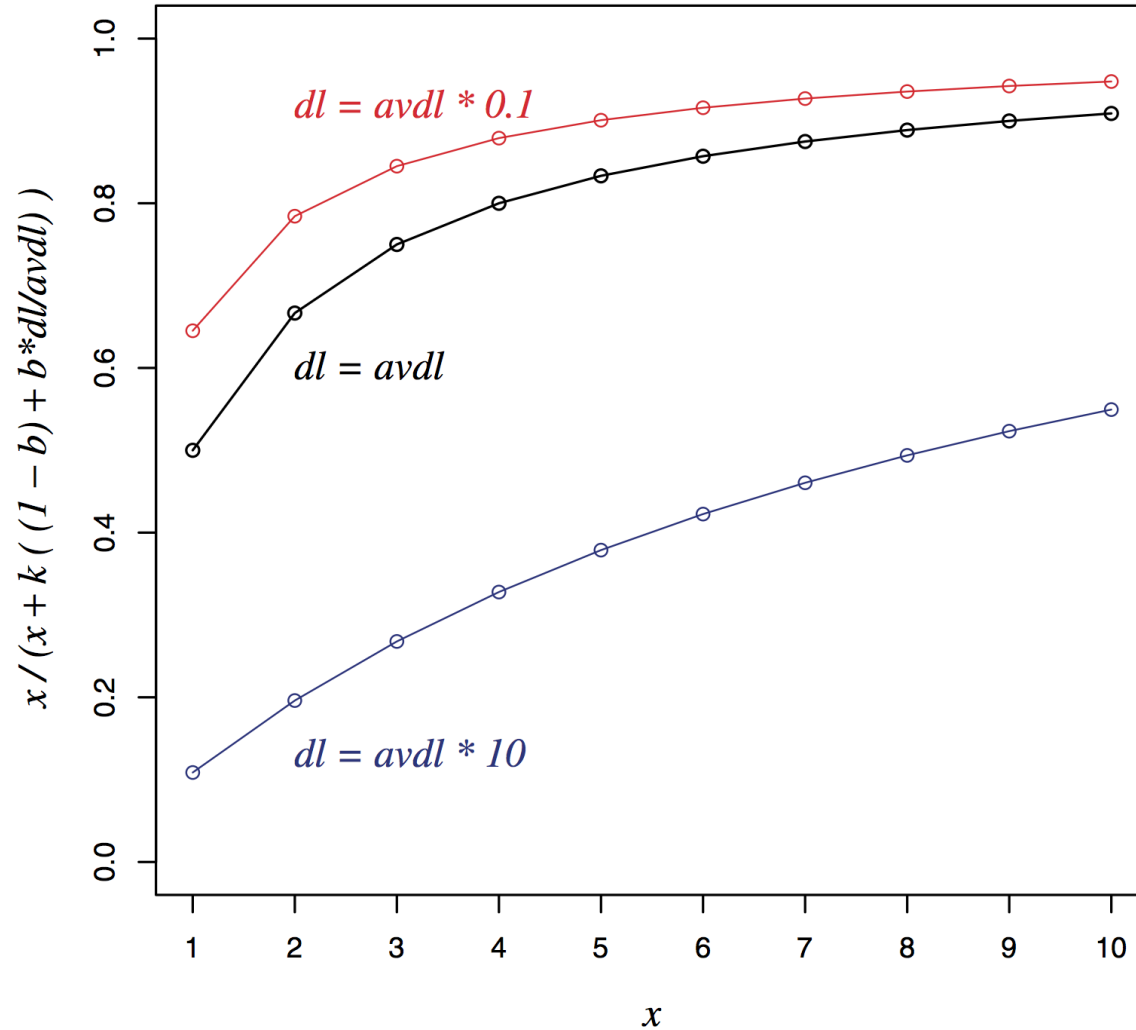
$$dl = \sum_{i \in V} tf_i$$

- *avdl*: Average document length over collection
- Length normalization component

$$B = \left( (1 - b) + b \frac{dl}{avdl} \right), \quad 0 \leq b \leq 1$$

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

# Document length normalization



# Okapi BM25

- Normalize  $tf$  using document length

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

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

- BM25 ranking function

$$RSV^{BM25} = \sum_{i \in q} c_i^{BM25}(tf_i);$$

# Okapi BM25

$$RSV^{BM25} = \sum_{i \in q} \log \frac{N}{df_i} \cdot \frac{(k_1 + 1)tf_i}{k_1((1 - b) + b \frac{dl}{avdl}) + tf_i}$$

- $k_1$  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: A Nonbinary Model

- If the query is long, we might also use similar weighting for query terms

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

- $tf_{tq}$ : term frequency in the query  $q$
- $k_3$ : tuning parameter controlling term frequency scaling of the query
- No length normalisation of queries (because retrieval is being done with respect to a single fixed query)
- The above tuning parameters should ideally be set to optimize performance on a development test collection. In the absence of such optimisation, experiments have shown reasonable values are to set  $k_1$  and  $k_3$  to a value between 1.2 and 2 and  $b = 0.75$