



# Probabilistic Information Retrieval

## Part II: In Depth

*Alexander Dekhtyar*

*Department of Computer Science*

*University of Maryland*

# In this part

---

## ✓ Probability Ranking Principle

- simple case
- case with retrieval costs

## ✓ Binary Independence Retrieval (BIR)

- Estimating the probabilities

## ✓ Binary Independence Indexing (BII)

- dual to BIR

# The Basics

---

- ✓ Bayesian probability formulas

$$p(a | b) p(b) = p(a \cap b) = p(b | a) p(a)$$

$$p(a | b) = \frac{p(b | a) p(a)}{p(b)}$$

$$p(\bar{a} | b) p(b) = p(b | \bar{a}) p(\bar{a})$$

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

# The Basics

---

- Document Relevance:

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

$$p(NR | x) = \frac{p(x | NR) p(NR)}{p(x)}$$

- Note:

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

# Probability Ranking Principle

---

- ✓ Simple case: no selection costs.
- ✓  $x$  is **relevant** iff  $p(R/x) > p(NR/x)$
- ✓ (*Bayes' Decision Rule*)
- ✓ PRP in action: Rank all documents by  $p(R/x)$ .

# Probability Ranking Principle

✓ More complex case: retrieval costs.

–  $C$  - cost of retrieval of relevant document

–  $C'$  - cost of retrieval of non-relevant document

– let  $d$ , be a document

✓ Probability Ranking Principle: if

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

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

A graphic of a spiral-bound notebook with a grey cover and a white page. The spiral binding is on the left side. A horizontal line is drawn across the page, and the text "Next: Binary Independence Model" is written in a black serif font below it.

Next: Binary Independence Model

# Binary Independence Model

---

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



# Binary Independence Model

---

- ✓ Queries: binary vectors of terms
- ✓ Given query  $q$ ,
  - for each document  $d$  need to compute  $p(R/q, d)$ .
  - replace with computing  $p(R/q, x)$  where  $x$  is vector representing  $d$

✓ Interested only in ranking

✓ Will use odds:

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

# Binary Independence Model

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

Constant for each query

Needs estimation

- Using **Independence** Assumption:

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

- So :  $O(R | q, d) = O(R | q) \cdot \prod_{i=1}^n \frac{p(x_i | R, q)}{p(x_i | NR, q)}$

# Binary Independence Model

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

- Since  $x_i$  is either 0 or 1:

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

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

Then...

# Binary Independence Model

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

All matching terms

Non-matching query terms

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

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.

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

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

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

# Binary Independence Model

- Estimating RSV coefficients.
- For each term  $i$  look at the following table:

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

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

Add 0.5 to every expression

# PRP and BIR: The lessons

---

- ✓ Getting reasonable approximations of probabilities is possible.
- ✓ Simple methods work only with restrictive assumptions:
  - *term independence*
  - *terms not in query do not affect the outcome*
  - *boolean representation of documents/queries*
  - *document relevance values are independent*
- ✓ Some of these assumptions can be removed



A graphic of a spiral-bound notebook with a grey cover and a white page. The spiral binding is on the left side. A thin horizontal line is drawn across the page, just above the text.

Next: Binary Independence  
Indexing

# Binary Independence Indexing vs. Binary Independence Retrieval

## •BIR

- ✓ Many Documents, One Query
- ✓ Bayesian Probability:

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

- ✓ Varies: document representation
- ✓ Constant: query (representation)

## •BII

- ✓ One Document, Many Queries
- ✓ Bayesian Probability

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

- ✓ Varies: query
- ✓ Constant: document

# Binary Independence Indexing

---

- ✓ “Learning” from queries
  - More queries: better results

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

- ✓  $p(q/x, R)$  - *probability that if document  $x$  had been deemed relevant, query  $q$  had been asked*
- ✓ The rest of the framework is similar to BIR

# Binary Independence Indexing:

## Key Assumptions

---

- ✓ Term occurrence in queries is *conditionally independent*: 
$$p(\vec{q} | R, \vec{x}) = \prod_{i=1}^n p(q_i | R, \vec{x})$$
- ✓ Relevance of document representation  $\mathbf{x}$  w.r.t. query  $\mathbf{q}$  depends only on the terms present in the query ( $q_i=1$ )
- ✓ For each term  $i$  not used in representation  $\mathbf{x}$  of document  $\mathbf{d}$  ( $x_i=0$ ):  $p(R | z_i, \vec{x}) = p(R | \vec{x})$ 
  - *only positive occurrences of terms count*

# Binary Independence Indexing

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

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

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

constant

Equal to 1