

# CS276A Text Retrieval and Mining

## Lecture 10

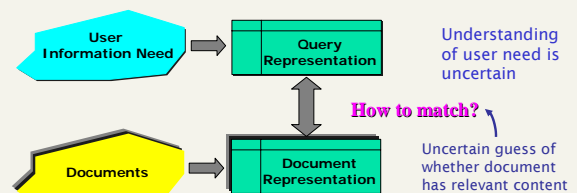
## Recap of the last lecture

- Improving search results
  - Especially for high recall. E.g., searching for *aircraft* so it matches with *plane*; *thermodynamic* with *heat*
- Options for improving results...
  - Global methods
    - Query expansion
      - Thesauri
      - Automatic thesaurus generation
    - Global indirect relevance feedback
  - Local methods
    - Relevance feedback
    - Pseudo relevance feedback

## Probabilistic relevance feedback

- Rather than reweighting in a vector space...
- If user has told us some relevant and some irrelevant documents, then we can proceed to build a probabilistic classifier, such as a Naive Bayes model:
  - $P(t_k|R) = |\mathbf{D}_{rk}| / |\mathbf{D}_r|$
  - $P(t_k|NR) = |\mathbf{D}_{nrk}| / |\mathbf{D}_{nr}|$ 
    - $t_k$  is a term;  $\mathbf{D}_r$  is the set of known relevant documents;  $\mathbf{D}_{rk}$  is the subset that contain  $t_k$ ;  $\mathbf{D}_{nr}$  is the set of known irrelevant documents;  $\mathbf{D}_{nrk}$  is the subset that contain  $t_k$ .

## 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.
- (Naïve) Bayesian Text Categorization
- 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 they've never won 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 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**
  - $P(\text{relevant}|\text{document}, \text{query})$

## Recall a few probability basics

- For events  $a$  and  $b$ :

- Bayes' Rule

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

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

$$p(a | b) = \frac{p(b | a)p(a)}{p(b)} = \frac{p(b | a)p(a)}{\sum_{x=a, \bar{a}} p(b | x)p(x)}$$

← 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  $x$  be a document in the collection.

Let  $R$  represent **relevance** of a document w.r.t. given (fixed)

query and let  $NR$  represent **non-relevance**.  $R = \{0, 1\}$  vs.  $NR/R$

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

$$p(R|x) = \frac{p(x|R)p(R)}{p(x)} \quad p(R), p(NR) \text{ - prior probability of retrieving a (non) relevant document}$$

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

$p(x|R)$ ,  $p(x|NR)$  - probability that if a relevant (non-relevant) document is retrieved, it is  $x$ .

## Probability Ranking Principle (PRP)

- Simple case: no selection costs or other utility concerns that would differentially weight errors
- Bayes' Optimal Decision Rule**
  - $x$  is **relevant** iff  $p(R|x) > p(NR|x)$
- PRP in action: Rank all documents by  $p(R|x)$
- 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 retrieval of relevant document
  - $C'$  - cost of retrieval of non-relevant 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**
- We won't further consider loss/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 Retrieval (BIR) – which we discuss later today – 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 length influence your judgments about document relevance?
    - One answer is the Okapi formulae (S. 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. lecture 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 same vector
- Bernoulli Naive Bayes model (cf. text categorization!)

## 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  $p(R|q, \vec{x})$  where  $\vec{x}$  is binary term incidence vector representing  $d$  Interested only in ranking
- Will use odds and Bayes' Rule:

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

## 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 a given query (points to  $\frac{p(R | q)}{p(NR | q)}$ )

Needs estimation (points to  $\frac{p(\vec{x} | R, q)}{p(\vec{x} | NR, q)}$ )

- 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...

This can be changed (e.g., in relevance feedback)

## Binary Independence Model

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

All matching terms
Non-matching query terms

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

All matching terms
All query terms

## Binary Independence Model

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

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 this table of document counts:

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:  $p_i \approx \frac{s}{S}$      $r_i \approx \frac{(n-s)/(N-S)}{s/(S-s)}$
- For now, assume no zero terms. More next lecture.

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

## 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
  - $\log(1-r_i)/r_i = \log(N-n)/n \approx \log N/n = \text{IDF!}$
- $p_i$  (probability of occurrence in relevant documents) can be estimated in various ways:
  - from relevant documents if know some
    - Relevance weighting can be used in feedback loop
  - constant (Croft and Harper combination match) – then just get idf weighting of terms
  - proportional to prob. of occurrence in collection
    - more accurately, to log of this (Greiff, SIGIR 1998)

## Iteratively estimating $p_i$

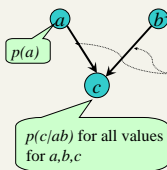
- Assume that  $p_i$  constant over all  $x_i$  in query
  - $p_i = 0.5$  (even odds) for any given doc
- Determine guess of relevant document set:
  - $V$  is fixed size set of highest ranked documents on this model (note: now a bit like tf.idf!)
- We need to improve our guesses for  $p_i$  and  $r_i$ , 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
    - $r_i = (n_i - |V_i|) / (N - |V|)$
- Go to 2. until converges then return ranking

24



## Bayesian Networks

$a, b, c$  - propositions (events).



- Bayesian networks model causal relations between events

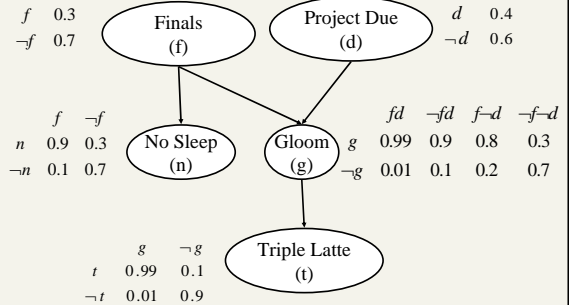
- Inference in Bayesian Nets:

- Given probability distributions for roots and conditional probabilities can compute a priori probability of any instance
- Fixing assumptions (e.g.,  $b$  was observed) will cause recomputation of probabilities

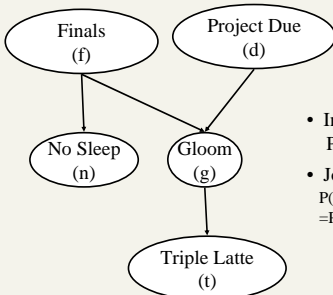
For more information see:

R.G. Cowell, A.P. Dawid, S.L. Lauritzen, and D.J. Spiegelhalter, 1999. *Probabilistic Networks and Expert Systems*. Springer Verlag.  
 J. Pearl, 1988. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan-Kaufman.

## Toy Example



## Independence Assumptions



- Independence assumption:  
 $P(t|g, f) = P(t|g)$
- Joint probability  
 $P(f, d, n, g, t)$   
 $= P(f) P(d) P(n|f) P(g|f, d) P(t|g)$

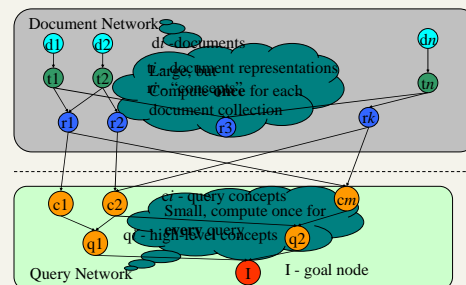
## Chained inference

- Evidence - a node takes on some value
- Inference
  - Compute *belief* (probabilities) of other nodes
    - conditioned on the known evidence
  - Two kinds of inference: *Diagnostic* and *Predictive*
- Computational complexity
  - General network: NP-hard
    - Tree-like networks are easily tractable
    - Much other work on efficient exact and approximate Bayesian network inference
      - Clever dynamic programming
      - Approximate inference ("loopy belief propagation")

## Model for Text Retrieval

- Goal
  - Given a user's information need (evidence), find probability a doc satisfies need
- Retrieval model
  - Model docs in a *document network*
  - Model information need in a *query network*

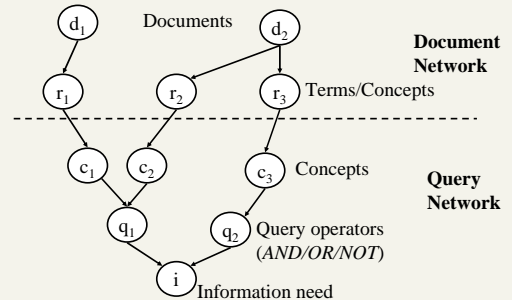
## Bayesian Nets for IR: Idea



## Bayesian Nets for IR

- Construct Document Network (once !)
- For each query
  - Construct best Query Network
  - Attach it to Document Network
  - Find subset of  $d_i$ 's which maximizes the probability value of node  $l$  (best subset).
  - Retrieve these  $d_i$ 's as the answer to query.

## Bayesian nets for text retrieval

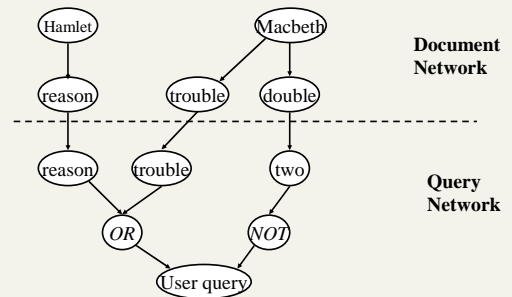


## Link matrices and probabilities

- Prior doc probability  $P(d) = 1/n$
- $P(r|d)$ 
  - within-document term frequency
  - $tf \times idf$  - based
- $P(c|r)$ 
  - 1-to-1
  - thesaurus
- $P(q|c)$ : canonical forms of query operators
  - Always use things like AND and NOT – never store a full CPT\*

\*conditional probability table

## Example: “reason trouble –two”



## Extensions

- Prior probs don't have to be  $1/n$ .
- “User information need” doesn't have to be a query - can be words typed, in docs read, any combination ...
- Phrases, inter-document links
- Link matrices can be modified over time.
  - User feedback.
  - The promise of “personalization”

## Computational details

- Document network built at indexing time
- Query network built/scored at query time
- Representation:
  - Link matrices from docs to any single term are like the postings entry for that term
  - Canonical link matrices are efficient to store and compute
- Attach evidence only at roots of network
  - Can do single pass from roots to leaves

## Bayes Nets in IR

---

- Flexible ways of combining term weights, which can generalize previous approaches
  - Boolean model
  - Binary independence model
  - Probabilistic models with weaker assumptions
- Efficient large-scale implementation
  - InQuery text retrieval system from U Mass
    - Turtle and Croft (1990) [[Commercial version defunct?](#)]
- Need approximations to avoid intractable inference
- Need to estimate all the probabilities by some means (whether more or less ad hoc)
- Much new Bayes net technology yet to be applied?

## Resources

---

- 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]
- F. Crestani, M. Lalmas, C. J. van Rijsbergen, and I. Campbell. 1998. Is This Document Relevant? ... Probably: A Survey of Probabilistic Models in Information Retrieval. *ACM Computing Surveys* 30(4): 528–552.  
<http://www.acm.org/pubs/citations/journals/surveys/1998-30-4/p528-crestani/>  
[Adds very little material that isn't in van Rijsbergen or Fuhr ]

## Resources

---

- H.R. Turtle and W.B. Croft. 1990. Inference Networks for Document Retrieval. *Proc. ACM SIGIR*: 1-24.
- E. Charniak. Bayesian nets without tears. *AI Magazine* 12(4): 50-63 (1991). <http://www.aaai.org/Library/Magazine/Vol12/12-04/vol12-04.html>
- D. Heckerman. 1995. A Tutorial on Learning with Bayesian Networks. Microsoft Technical Report MSR-TR-95-06  
<http://www.research.microsoft.com/~heckerman/>
- N. Fuhr. 2000. Probabilistic Datalog: Implementing Logical Information Retrieval for Advanced Applications. *Journal of the American Society for Information Science* 51(2): 95–110.
- R. K. Belew. 2001. *Finding Out About: A Cognitive Perspective on Search Engine Technology and the WWW*. Cambridge UP 2001.  
MIR 2.5.4, 2.8