# Retrieval Models Overview

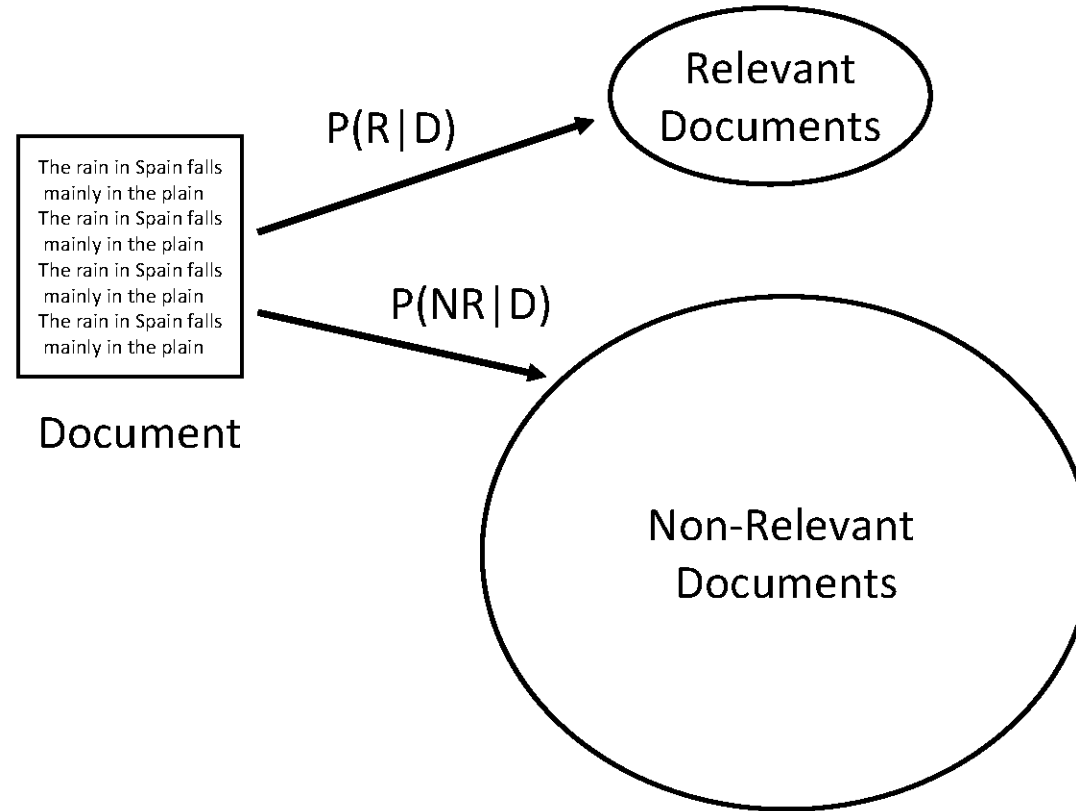- Older models
  - Boolean retrieval
  - Vector Space model
- <span style="color:red">Probabilistic Models</span>
  - <span style="color:red">Language models</span>
  - <span style="color:red">BM25</span>
- Combining evidence
  - Inference networks
  - Learning to Rank

# Probability Ranking Principle

- Robertson (1977)
  - "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."

# IR as Classification

The rain in Spain falls
mainly in the plain
The rain in Spain falls
mainly in the plain
The rain in Spain falls
mainly in the plain
The rain in Spain falls
mainly in the plain

Document

$P(R|D)$

$P(NR|D)$

Relevant
Documents

Non-Relevant
Documents

# Bayes Classifier

- Bayes Decision Rule
  - A document *D* is relevant if *P*(*R*|*D*) > *P*(*NR*|*D*)

- Estimating probabilities
  - Use Bayes Rule

  $$P(R|D) = \frac{P(D|R)P(R)}{P(D)} \qquad P(NR\,|\,D) = \frac{P(D\,|\,NR)P(NR)}{P(D)}$$

  - Classify a document as relevant if

  $$\frac{P(D\,|\,R)P(R)}{P(D)} > \frac{P(D\,|\,NR)P(NR)}{P(D)}$$

    - lhs is *likelihood ratio*

  $$\frac{P(D|R)}{P(D|NR)} > \frac{P(NR)}{P(R)}$$

# Estimating P(D|R)

- Assume independence

$$P(D|R) = \prod_{i=1}^{t} P(d_i|R)$$

- *Binary independence model*
  - document represented by a vector of binary features indicating term occurrence (or non-occurrence)
  - Assume:
    - $p_i$ is probability that term i occurs (i.e., has value 1) in relevant document
    - $s_i$ is probability of occurrence in non-relevant document

# Binary Independence Model

$$\frac{P(D|R)}{P(D|NR)} = \prod_{i:d_i=1} \frac{p_i}{s_i} \cdot \prod_{i:d_i=0} \frac{1-p_i}{1-s_i}$$

$$= \prod_{i:d_i=1} \frac{p_i}{s_i} \cdot \left( \prod_{i:d_i=1} \frac{1-s_i}{1-p_i} \cdot \prod_{i:d_i=1} \frac{1-p_i}{1-s_i} \right) \cdot \prod_{i:d_i=0} \frac{1-p_i}{1-s_i}$$

$$= \prod_{i:d_i=1} \frac{p_i(1-s_i)}{s_i(1-p_i)} \cdot \prod_i \frac{1-p_i}{1-s_i}$$

# Binary Independence Model

- Scoring function is

$$\sum_{i:d_i=1} \log \frac{p_i(1-s_i)}{s_i(1-p_i)}$$

- Query provides information about relevant documents

- If we assume $p_i$ constant, $s_i$ approximated by entire collection, get *idf*-like weight

$$\log \frac{0.5(1-\frac{n_i}{N})}{\frac{n_i}{N}(1-0.5)} = \log \frac{N-n_i}{n_i}$$

# Contingency Table

|  | Relevant | Non-relevant | Total |
|---|---|---|---|
| $d_i = 1$ | $r_i$ | $n_i - r_i$ | $n_i$ |
| $d_i = 0$ | $R - r_i$ | $N - n_i - R + r_i$ | $N - n_i$ |
| Total | $R$ | $N - R$ | $N$ |

$$p_i = (r_i + 0.5)/(R + 1)$$

$$s_i = (n_i - r_i + 0.5)/(N - R + 1)$$

Gives scoring function:

$$\sum_{i:d_i=q_i=1} \log \frac{(r_i+0.5)/(R-r_i+0.5)}{(n_i-r_i+0.5)/(N-n_i-R+r_i+0.5)}$$

# BM25

- Popular and effective ranking algorithm based on binary independence model
  - adds document and query term weights

$$\sum_{i \in Q} \log \frac{(r_i+0.5)/(R-r_i+0.5)}{(n_i-r_i+0.5)/(N-n_i-R+r_i+0.5)} \cdot \frac{(k_1+1)f_i}{K+f_i} \cdot \frac{(k_2+1)qf_i}{k_2+qf_i}$$

  - $k_1$, $k_2$ and $K$ are parameters whose values are set empirically
$$K = k_1((1-b)+b \cdot \frac{dl}{avdl})$$
  - $dl$ is document length
  - $avdl$ is average document length
  - Typical value for $k_1$ is 1.2, $k_2$ varies from 0 to 1000, b = 0.75

# BM25 Example

- Query with two terms, "president lincoln", ($qf = 1$)
- No relevance information ($r$ and $R$ are zero)
- $N$ = 500,000 documents
- *"president"* occurs in 40,000 documents ($n_1 = 40,000$)
- *"lincoln"* occurs in 300 documents ($n_2 = 300$)
- "president" occurs 15 times in doc ($f_1 = 15$)
- *"lincoln"* occurs 25 times ($f_2 = 25$)
- document length is 90% of the average length ($dl/avdl = .9$)
- $k_1 = 1.2$, $b = 0.75$, and $k_2 = 100$
- $K = 1.2 \cdot (0.25 + 0.75 \cdot 0.9) = 1.11$

# BM25 Example

$$BM25\ (Q,D) = \sum_{i \in Q} \log \frac{(r_i+0.5)/(R-r_i+0.5)}{(n_i-r_i+0.5)/(N-n_i-R+r_i+0.5)} \cdot \frac{(k_1+1)f_i}{K+f_i} \cdot \frac{(k_2+1)qf_i}{k_2+qf_i}$$

$$BM25(Q,D) =$$

$$\log \frac{(0+0.5)/(0-0+0.5)}{(40000-0+0.5)/(500000-40000-0+0+0.5)}$$

$$\times \frac{(1.2+1)15}{1.11+15} \times \frac{(100+1)1}{100+1}$$

$$+\log \frac{(0+0.5)/(0-0+0.5)}{(300-0+0.5)/(500000-300-0+0+0.5)}$$

$$\times \frac{(1.2+1)25}{1.11+25} \times \frac{(100+1)1}{100+1}$$

$$= \log 460000.5/40000.5 \cdot 33/16.11 \cdot 101/101$$

$$+\log 499700.5/300.5 \cdot 55/26.11 \cdot 101/101$$

$$= 2.44 \cdot 2.05 \cdot 1 + 7.42 \cdot 2.11 \cdot 1$$

$$= 5.00 + 15.66 = 20.66$$

# BM25 Example

- Effect of term frequencies

| Frequency of "president" | Frequency of "lincoln" | BM25 score |
|---|---|---|
| 15 | 25 | 20.66 |
| 15 | 1 | 12.74 |
| 15 | 0 | 5.00 |
| 1 | 25 | 18.2 |
| 0 | 25 | 15.66 |

# Summary

- Probabilistic models for information retrieval
  - Language models
  - BM25

- More advanced models available
  - Learning to rank
  - Neural/dense retrieval models