



Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

# Processamento e Recuperação de Informação

## Information Retrieval Models

Departamento de Engenharia Informática  
Instituto Superior Técnico

1º Semestre  
2018/2019



# Outline

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- 1 Generic Document Model
- 2 The Boolean Model
- 3 The Vector Space Model
- 4 The Probabilistic Model
- 5 Comparison of the Classic Models



# Bibliography

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

Ricardo Baeza-Yates, Berthier Ribeiro-Neto, Modern Information Retrieval, 2nd edition. Chapter 3

Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, Introduction to Information Retrieval. Chapters 1, 6, 11 and 12

Bing Liu, Web Data Mining - Exploring Hyperlinks, Contents, and Usage Data. Chapter 6.



# Outline

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- 1 Generic Document Model
- 2 The Boolean Model
- 3 The Vector Space Model
- 4 The Probabilistic Model
- 5 Comparison of the Classic Models



# Retrieval Models

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models





# Index Terms

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

In the classic IR models, documents are represented by **index terms**

- full text/selected keywords
- structure/no structure

Not all terms are equally useful

- index terms can be **weighted**

We assume that terms are **mutually independent**

- this is, of course, a simplification



# Definition of a Document Model

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## Definition

Let  $t$  be the number of index terms in the collection of documents, and let  $k_i$  be a generic index term.

- $K = \{k_1, \dots, k_t\}$  is the set of all index terms.
- A weight  $w_{i,j} \geq 0$  is associated with each index term  $k_i$  of a document  $d_j$ .
- For an index term which does not appear in the document text,  $w_{i,j} = 0$ .
- Each document  $d_j$  is associated a term vector  $\vec{d}_j$ , represented by  $\vec{d}_j = (w_{1,j}, w_{2,j}, \dots, w_{t,j})$ .
- Function  $g_i(\vec{d}_j)$  returns the weight of index term  $k_i$  in vector  $\vec{d}_j$ .



# Outline

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- 1 Generic Document Model
- 2 The Boolean Model
- 3 The Vector Space Model
- 4 The Probabilistic Model
- 5 Comparison of the Classic Models





# Boolean Model Queries

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- Follows Boolean algebra syntax and semantics
- Term weights are binary
  - $w_{i,j} \in \{0, 1\}$
  - $w_{i,j} = 1$  — term present,
  - $w_{i,j} = 0$  — term not present
- Queries are Boolean expressions
  - E.g.,  $q = k_a \wedge (k_b \vee \neg k_c)$
- Documents are considered **relevant** if the query evaluates to 1 (true)



# An Example

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

$d_1$

That government is best which  
governs least

$d_2$

That government is best which  
governs not at all

$d_3$

When men are prepared for it,  
that will be the kind of  
government which they will have



# An Example

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

$d_1$

That government is best which  
governs least

$d_2$

That government is best which  
governs not at all

$d_3$

When men are prepared for it,  
that will be the kind of  
government which they will have

$q = \text{government} \wedge \text{best}$

answer:  $d_1, d_2$



# An Example

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

$d_1$

That government is best which  
governs least

$q = \text{government} \wedge \text{best}$

answer:  $d_1, d_2$

$d_2$

That government is best which  
governs not at all

$q = \text{government} \wedge \text{best} \wedge \neg \text{all}$

answer:  $d_1$

$d_3$

When men are prepared for it,  
that will be the kind of  
government which they will have



# An Example

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

$d_1$

That government is best which  
governs least

$q = \text{government} \wedge \text{best}$

answer:  $d_1, d_2$

$d_2$

That government is best which  
governs not at all

$q = \text{government} \wedge \text{best} \wedge \neg \text{all}$

answer:  $d_1$

$d_3$

When men are prepared for it,  
that will be the kind of  
government which they will have

$q = \text{government} \vee \text{best} \wedge \neg \text{all}$

answer:  $d_1, d_2, d_3$



# Document-Query Similarity

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- Queries can be translated to a disjunction of conjunctive vectors

$$\vec{q} = k_a \wedge (k_b \vee \neg k_c) \Leftrightarrow (1, 1, 1) \vee (1, 1, 0) \vee (1, 0, 0)$$

each tuple corresponds to a vector  $(k_a, k_b, k_c)$

- Similarity of a document to a query is defined as:

$$\text{sim}(d_j, q) = \begin{cases} 1 & \text{if } \exists \vec{q}_c \in \vec{q} \mid \forall_i, g_i(\vec{d}_j) = g_i(\vec{q}_c) \\ 0 & \text{otherwise} \end{cases}$$



# The Boolean Model

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## Why is it good?

- Simple model based on Boolean algebra
- Intuitive concept
- Precise semantics
- Clear formal basis
- Widely adopted by early information systems



# Boolean Model

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## Limitations:

- Retrieval based only on binary decisions
  - More similar to *data retrieval* than *information retrieval*
  - Can retrieve too many, or too little documents
  - Some documents may be more relevant than others
- How do you translate a query to a Boolean expression?
  - Non-expert users may not be able to represent their information needs using Boolean expressions
- Terms are all equally important
  - Index term weighting can bring great improvements in performance





# Outline

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- 1 Generic Document Model
- 2 The Boolean Model
- 3 The Vector Space Model**
- 4 The Probabilistic Model
- 5 Comparison of the Classic Models



# Documents as Vectors

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- Documents are represented as vectors
  - $\vec{d}_j = (w_{1,j}, w_{2,j}, \dots, w_{t,j})$
  - $w_{i,j}$  is the weight of term  $i$  in document  $j$
- Queries are also vectors
  - $\vec{q} = (w_{1,q}, w_{2,q}, \dots, w_{t,q})$
- Vector operations can be used to compare queries  $\times$  documents (or documents  $\times$  documents)



# An Example

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models





# Defining Document Vectors

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

Two questions are still unanswered:

- 1 How do we define term weights?
- 2 How do we compare documents to queries?



# Defining Term Weights — TF

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## Term frequency

Term frequency is a measure of term importance **within a document**

## Definition

Let  $N$  be the total number of documents in the system and  $n_i$  be the number of documents in which term  $k_i$  appears. The **normalized frequency** of a term  $k_i$  in document  $d_j$  is given by:

$$f_{i,j} = \frac{freq_{i,j}}{\max_l freq_{l,j}}$$

where  $freq_{i,j}$  is the number of occurrences of term  $k_i$  in document  $d_j$ .



# Defining Term Weights — IDF

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## (Inverse) Document frequency

Document frequency is a measure of term importance **within a collection**

## Definition

The **inverse document frequency** of a term  $k_i$  is given by:

$$idf_i = \log \frac{N}{n_i}$$



# Defining Term Weights — TF-IDF

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## Definition

The weight of a term  $k_i$  in document  $d_j$  for the vector space model is given by the **tf-idf** formula:

$$w_{i,j} = f_{i,j} \times \log \frac{N}{n_i}$$



# Document Similarity

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- Similarity between documents and queries is a measure of the correlation between their vectors
- Documents/queries that share the same terms, with similar weights, should be more similar
- Thus, as similarity measure, we use the **cosine of the angle between the vectors**

$$\text{sim}(d_j, q) = \frac{\vec{d}_j \cdot \vec{q}}{|\vec{d}_j| \times |\vec{q}|} = \frac{\sum_{i=1}^t w_{i,j} \times w_{i,q}}{\sqrt{\sum_{i=1}^t w_{i,j}^2} \times \sqrt{\sum_{i=1}^t w_{i,q}^2}}$$





# An Example

Processamento  
e Recuperação  
de Informação

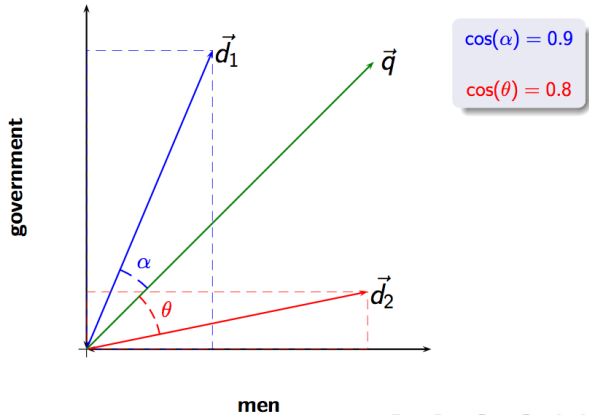
Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models





# The Vector Space Model

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## Why is it so good?

- Simple model, based on linear algebra
- Term weights are not binary
- Allows computing a continuous **degree of similarity** between queries and documents
- Thus, allows **ranking** documents according to their possible relevance



# Improving the VSM (1)

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## The BM25 Model

Consider not only the term frequency and inverse document frequency heuristics, but also the **document length** as a **normalization factor** for the term frequency.

$$TF_{i,j} = \frac{f_{i,j} \times (k_1 + 1)}{f_{i,j} + k_1 \times \left(1 - b + b \frac{|d_j|}{avgdl}\right)}$$

$$IDF_i = \log \frac{N - n_i + 0.5}{n_i + 0.5}$$

$$sim(d_j, q) = \sum_{i \in q} IDF_i \times TF_{i,j}$$



# Improving the VSM (2)

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## Latent Semantic Indexing

- Find a low-rank approximation of the matrix which describes the occurrences of terms in documents
  - Singular Value Decomposition
  - Compare the documents in the low-dimensional space
- The consequence of the rank lowering is that some dimensions are combined (e.g., **mitigates the problem of identifying synonymy**)
- To be detailed latter in the course



# Outline

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- 1 Generic Document Model
- 2 The Boolean Model
- 3 The Vector Space Model
- 4 The Probabilistic Model**
- 5 Comparison of the Classic Models



# The Probabilistic Model

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- Models the IR problem in a probabilistic framework
- Estimates the probability of document  $d_j$  being relevant to the user
- A simple probabilistic model can assume that:
  - 1 the probability depends only on the query and the document
  - 2 there is a subset  $R$  of relevant documents
  - 3 index terms are independent
- A simple probabilistic model can use binary term weights



# Document Query Similarity

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- As a similarity measure, we use the ratio between the probability of finding the relevant documents and the probability of finding the non-relevant documents

$$\text{sim}(d_j, q) = \frac{P(R|\vec{d}_j)}{P(\bar{R}|\vec{d}_j)}$$



# Similarity Probabilities (1)

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## Initial Equation

$$\text{sim}(d_j, q) = \frac{P(R|\vec{d}_j)}{P(\bar{R}|\vec{d}_j)} = \frac{P(\vec{d}_j|R) \times P(R)}{P(\vec{d}_j|\bar{R}) \times P(\bar{R})} \sim \frac{P(\vec{d}_j|R)}{P(\vec{d}_j|\bar{R})}$$

## Assuming term independence...

$$\text{sim}(d_j, q) \sim \frac{(\prod_{g_i(\vec{d}_j)=1} P(k_i|R)) \times (\prod_{g_i(\vec{d}_j)=0} P(\bar{k}_i|R))}{(\prod_{g_i(\vec{d}_j)=1} P(k_i|\bar{R})) \times (\prod_{g_i(\vec{d}_j)=0} P(\bar{k}_i|\bar{R}))}$$

## Tanking logs and removing constant factors...

$$\text{sim}(d_j, q) = \sum_{i=1}^t w_{i,q} \times w_{i,j} \times \left( \log \frac{P(k_i|R)}{1 - P(k_i|R)} + \log \frac{1 - P(k_i|\bar{R})}{P(k_i|\bar{R})} \right)$$





# Similarity Probabilities (2)

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## Blind assumptions

$$\begin{aligned}P(k_i|R) &= 0.5 \\P(k_i|\overline{R}) &= \frac{n_i}{N}\end{aligned}$$

## After document retrieval...

Let  $V$  be the number of returned documents; let  $V_i$  be the number of returned documents containing term  $k_i$ .

$$\begin{aligned}P(k_i|R) &= \frac{V_i}{V} \\P(k_i|\overline{R}) &= \frac{n_i - V_i}{N - V}\end{aligned}$$



# Similarity Probabilities (3)

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

Avoiding small values...

$$P(k_i|R) = \frac{V_i + \frac{n_i}{N}}{V+1}$$
$$P(k_i|\bar{R}) = \frac{n_i - V_i + \frac{n_i}{N}}{N - V + 1}$$



# Problems of this Simple Probabilistic Model

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- There is no accurate estimate for the first run probabilities
  - Index terms are not weighted
  - Terms are assumed mutually independent
- 
- In fact, many different probabilistic retrieval models have been proposed!



# Probabilistic Language Models

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- Another simple probabilistic retrieval model
- Each document  $d$  is treated as (the basis for) a **probabilistic language model**
- Given a query  $q$  rank documents based on  $P(d|q)$

$$P(d|q) = \frac{P(d) \times P(q|d)}{P(q)}$$

- The evidence  $P(q)$  is the same for all documents, so ignore
- $P(d)$  is the prior
  - often treated as the same for all  $d$
  - we can give a higher prior to “high-quality” documents (e.g., those with high PageRank – to be seen latter)
- $P(q|d)$  is likelihood, i.e. the probability of  $q$  given  $d$



# How to compute $P(q|d)$

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- Conditional independence assumption

$$P(q|d) = P(\{t_1, \dots, t_{|q|}\}|d) = \prod_{1 \leq i \leq |q|} P(t_i|d)$$

- $|q|$  is length of  $q$
- $t_k$  is the token occurring at position  $k$  in  $q$

- The above multinomial model is equivalent to:

$$P(q|d) = \prod_{\text{distinct term } t \in q} P(t_k|d)^{TF_{t,q}}$$

- Component  $TF_{t,q}$  is the term frequency of  $t$  in  $q$
- Parameters  $P(t_k|d)$  computed through maximum likelihood estimates

$$P(t_k|d) = \frac{TF_{t_k,d}}{|d|}$$



# LM Retrieval and Naïve Bayes

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

The next class will introduce a simple probabilistic document classifier, known as the **Naïve Bayes** approach

- We want to classify document  $d$ .  
We want to classify a query  $q$
- Human-defined classes: e.g., politics, economics, sports.  
Each document in the collection is a different class
- Assume that  $d$  was produced by the generative model.  
Assume that  $q$  was generated by a generative model
- Which of the classes (= class models) is most likely to have generated the document  $d$ ?  
Which document (=class) is most likely to have generated the query  $q$ ?
- For which class do we have the most evidence?  
For which document (as source for query) do we have the most evidence?



# Outline

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- 1 Generic Document Model
- 2 The Boolean Model
- 3 The Vector Space Model
- 4 The Probabilistic Model
- 5 Comparison of the Classic Models



# What makes these Models Work?

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

## Three main term weighting normalization driving features:

- TF - Term Frequency
- IDF - Inverse Document Frequency
- DL - Document Length





# Comparison of the Classic Models

Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

- Boolean model is considered the weakest
- There is some controversy over which shows better performance: vector space or probabilistic
- However, nowadays, the vector space model is the most widely used



Processamento  
e Recuperação  
de Informação

Generic  
Document  
Model

The Boolean  
Model

The Vector  
Space Model

The  
Probabilistic  
Model

Comparison of  
the Classic  
Models

# Questions?