

1 Um estudo sobre o uso do modelo probabilístico no desenvolvimento de
2 uma máquina de busca

3

4 Stephany Castro da SILVA¹;

5 stephanycastro.es@gmail.com;

6 Aurélio Andrade de Menezes JÚNIOR¹.

7 aurelio.menezesjr@gmail.com

8

9

10

11 ¹Universidade Federal do Amazonas (UFAM), Instituto de Ciências Exatas e Tecnologia
12 (ICET), Itacoatiara, AM, Brasil.

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28 UM ESTUDO SOBRE O USO DO MODELO PROBABILÍSTICO NO
29 DESENVOLVIMENTO DE UMA MÁQUINA DE BUSCA

30

31 **Resumo**

32 Este trabalho propôs a criação de uma base de dados que contivesse de forma
33 digitalizada as devidas Resoluções existentes e propôs também a realização de um
34 estudo sobre o uso do Modelo BM25 no desenvolvimento de uma Máquina de Busca
35 voltada para resoluções institucionais utilizadas pela Universidade Federal do
36 Amazonas, como um meio de minimizar o problema de acesso e recuperação de tais
37 instrumentos normativos.

38 **Palavras-chave:** Modelo BM25; Resoluções; Máquina de Busca.

39 **Abstract**

40 This paper proposes the creation of a database that contained digitized so appropriate
41 existing resolutions and also proposed to conduct a study on using the BM25 model in
42 developing a focused search machine for institucional resolutions used by the Federal
43 University of Amazonas as a means of minimizing the problem of access and retrieval
44 of such normative tools.

45 **Keywords:** Model BM25; Resolutions; Search engine.

46

47

48

49

50

51 **Introdução**

52 Com o advento da Web e de sua grande utilização, a quantidade de banco de
53 dados existentes em instituições, tais como, escolas, empresas e associações, cresceu
54 rapidamente. Isso aumentou a quantidade de dados existentes e de informações trocadas
55 de forma eletrônica. De acordo com (Júnior, 2012), a Web tem revolucionado tanto o
56 acesso a informações pessoais quanto o gerenciamento do conhecimento em
57 instituições, assim como bancos de dados tem sido usado com grande frequência no
58 nosso dia a dia nas empresas para armazenamento e controle de informações.

59 Entretanto, se hoje temos a facilidade de acessar grandes quantidades de
60 informações, por outro lado, também temos um problema causado por isso: Como
61 encontrar uma determinada informação dentro de grandes quantidades de Informações.
62 Para resolver esse problema utiliza-se a Recuperação da Informação - RI.

63 RI é uma área responsável pela representação, armazenamento, organização e
64 acesso a itens de informação (Garcia, 2002). Além disso, é uma ciência de pesquisa
65 relacionada à busca por informações que podem estar contidas em determinados
66 documentos, busca pelos documentos desejados pelo usuário e busca em banco de
67 dados, este último podendo ser relacionais e isolados, ou até mesmo interligado a Web.
68 Esses documentos podem conter qualquer tipo de mídia, tais como, texto, imagem e
69 som. Na maioria das vezes, são compostos de texto em linguagem natural ou de
70 informação textual associada a outros tipos de dados.

71 Dentro desse contexto, há as Máquinas de busca que existem para facilitar o
72 acesso a informações relevantes de acordo com o desejo do usuário. Máquinas de busca
73 são projetadas para operarem ambientes onde a quantidade de conteúdo disponível

74 supera a capacidade do usuário de acessá-lo de forma eficiente durante a pesquisa
75 (Júnior, 2012). Esses sistemas são implementados a partir de Modelos em Recuperação
76 de Informação que são o núcleo de qualquer sistema de RI. Tais modelos são utilizados
77 para representar características semânticas dos elementos envolvidos nos sistemas,
78 sendo os três tipos de modelos para Recuperação de Informação que possuem o seu uso
79 mais difundido: o modelo booleano, o modelo vetorial e o modelo probabilístico.

80 O modelo booleano é composto por conjuntos de documentos e operações
81 clássicas da teoria de conjuntos. O modelo vetorial é composto por documentos e
82 consultas representados como vetores em um espaço n dimensional e operações de
83 álgebra linear aplicáveis aos vetores (Oliveira, 2010). O modelo probabilístico define as
84 representações para documentos e consultas baseado na teoria das probabilidades. A
85 base desse modelo está no princípio da ordenação probabilística: onde dada uma
86 consulta e um documento, tenta-se estimar a probabilidade do usuário considerar o
87 documento relevante à consulta (Oliveira, 2010).

88 Baseado em tais modelos clássicos, outros modelos, ou funções específicas
89 foram sendo criados, tais como, a função BM25 que foi desenvolvido baseado em
90 modelos probabilísticos. Tal função de classificação é utilizada para atribuir relevância
91 aos documentos existentes em uma Máquina de Busca, nele a busca de uma informação
92 é baseada na probabilidade de um documento ser relevante para a consulta (Martins,
93 2009).

94 A função BM25 é baseada no conceito de utilização de uma “bolsa” de palavras
95 para classificar um conjunto de documentos de acordo com a sua relevância para uma
96 determinada consulta, independentemente da inter-relação entre os termos da consulta
97 em um documento. O desempenho de um sistema de RI é avaliado de acordo com a sua

98 capacidade em recuperar o maior número de itens relevantes, ao mesmo tempo em que
99 filtra ao máximo os itens irrelevantes (Oliveira, 2010).

100 Sendo o objetivo de uma máquina de busca facilitar o acesso a informação
101 relevante, conhecer a sua eficiência é importante não só para os pesquisadores de RI,
102 mas também para quem usa estes sistemas. Pesquisadores e usuários precisam ter
103 maneiras efetivas para saber quão bons são os sistemas para uma dada tarefa e como
104 estes podem ser melhorados.

105 Algumas medidas mais comuns para avaliar o desempenho de um sistema
106 computacional são tempo e espaço. Se o tempo de resposta do sistema for pequeno e o
107 tamanho do espaço de memória utilizado também, o sistema pode ser considerado
108 excelente. Entretanto, quando se trata de sistemas cujo objetivo é recuperar informações
109 outras medidas devem ser utilizadas.

110 Sendo assim, para determinar se um sistema de busca é bom ou não, utiliza-se
111 Métodos de Avaliação específicos, que verificam um sistema de RI através da
112 comparação das respostas geradas por este e o conjunto ideal de respostas esperadas.
113 Isto fornece uma estimativa da qualidade do algoritmo de recuperação de informação
114 avaliado. As métricas usuais para a avaliação do resultado de um sistema de recuperação
115 de informação são: revocação e precisão. Para que estas medidas sejam relevantes, é
116 necessário conhecer bem o conteúdo dos documentos da coleção (Oliveira, 2010).
117 Revocação mede a proporção de documentos relevantes que foram retornados como
118 resultado a uma consulta do usuário e a Precisão mede quantos documentos relevantes
119 foram recuperados

120

121 **Material e Métodos**

122 O projeto foi elaborado de acordo com as etapas previstas no cronograma. As
123 atividades realizadas foram referentes à:

124 **Realizar levantamento bibliográfico de trabalhos relacionados a métodos**
125 **de avaliação.** Foram estudadas e analisadas as técnicas e aplicações de um método de
126 avaliação em uma Máquina de Busca. Assim, foi possível conhecer a fundo esses
127 determinados tipos de métodos, verificar a aplicação de um método de avaliação de
128 desempenho na Máquina de Busca, para aplicar tais métodos na prática para a validação
129 da Máquina de Busca criada.

130 **Analisar o comportamento do Modelo BM25.** Foi avaliada o funcionamento
131 do Modelo BM25, por meio de estudos dos conceitos referentes a esse tipo de modelo e
132 verificado a sua aplicação no desenvolvimento de uma Máquinas de Busca.

133 **Realizar levantamento de resoluções da UFAM que possam ser**
134 **digitalizadas para a criação de uma base de dados e analisar os conceitos referentes**
135 **à implementação de uma Máquina de Busca.** Foram levantadas resoluções que
136 podiam ser digitalizadas para a criação de uma base de dados. Após isso, foram
137 estudados os conceitos sobre Máquina de Busca e sobre o funcionamento de tais
138 sistemas de RI, para por em prática e implementar o Modelo BM25 para aplicá-lo em
139 uma Máquina de Busca.

140 **Estudar e analisar as técnicas e aplicações de um método de avaliação em**
141 **uma Máquina de Busca.** Foram realizados estudos referentes a esses determinados
142 tipos de métodos. Além disso, foi verificado métodos de avaliação de desempenho para
143 serem aplicadas na validação da Máquina de Busca desenvolvida.

144 **Resultados**

145 Durante os primeiros meses de execução do projeto, foi realizado levantamento
146 bibliográfico de trabalhos e artigos relacionados a área do estudo. Os assuntos
147 levantados foram referentes à Recuperação da Informação, dentro da qual, foram
148 realizados estudos abordando o Modelo Probabilístico (modelo base para a construção
149 do Modelo BM25), a qual o mesmo, foi utilizado durante o decorrer do trabalho, mais
150 precisamente para a avaliação do seu uso no desenvolvimento de uma Máquina de
151 Busca.

152 Sendo um dos objetivos do projeto a implementação do Modelo BM25 como
153 uma forma de facilitar e agilizar o processo de consulta de um conjunto de resoluções
154 usadas pela Universidade Federal do Amazonas – UFAM, fez-se necessário a
155 construção de uma base de dados contendo resoluções, obtidas a partir da página oficial
156 da Instituição. Os textos obtidos foram posteriormente todos convertidos para o formato
157 .txt, dado a maior facilidade em aplicar as buscas em documentos com este formato.

158 Com a compreensão dos principais assuntos relacionados ao Modelo BM25,
159 foi possível criar uma versão do mesmo, capaz de executar os principais passos do
160 processo de recuperação da informação que são: consulta, indexação e pesquisa. A
161 consulta consiste na especificação de um conjunto de termos que representam a
162 necessidade de informação do usuário, a indexação envolve a criação de estruturas de
163 dados associadas aos documentos de uma coleção, onde se encontra a possível
164 informação que foi requerida pelo usuário. A pesquisa envolve o processo de recuperar
165 os documentos de acordo com a consulta do usuário. A ordenação dos documentos
166 recuperados dos mais relevantes para o menos relevantes, é feito usando Modelos como
167 o BM25.

168 Como pode ser observado na **Figura 1** e **Figura 2**, a primeira etapa de
169 funcionamento da estrutura da máquina de busca, consiste na Busca informada pelo
170 usuário, ou seja, é a entrada dos termos para a pesquisa das resoluções na base onde
171 estão as resoluções levantadas em etapas anteriores do projeto. A segunda parte da
172 busca é feita aplicando graus (valores) de relevância na coleção como pode ser
173 observado na **Figura 3**, que permitem posteriormente o Modelo BM25, classificar cada
174 documento relevante para a consulta e por fim realizar a ordenação de acordo com a
175 relevância. Assim, quanto maior o grau de relevância do documento, mais importante
176 ele é para o usuário.

177 As métricas utilizadas para avaliação da eficiência da Máquina de Busca
178 desenvolvida por este projeto, foram a de precisão e revocação. Tais métodos foram
179 selecionados a partir dos estudos realizados sobre o assunto, onde identificou-se que
180 ambas são métricas comumente utilizadas dentro da área de Recuperação da
181 Informação, como forma de avaliar rankings de busca. O processo de avaliação da
182 estratégia utilizado para essa avaliação, seguiu os seguintes passos:

183 1. Criação da lista de consultas com as palavras-chave: a criação da lista
184 iniciou-se com a escolha de 11 palavras – chave, as mesmas foram selecionados após a
185 realização de análises na base de dados onde se identificou as palavras que mais se
186 repetiam durante o decorrer do texto.

187 2. Identificação dos documentos mais relevantes: foram identificados os
188 documentos relevantes e foi atribuído um peso de 1 a 8, onde 8 indica o documento
189 mais relevante e 1 indica o documento menos relevante a uma dada consulta.

190 3. Execução da busca no sistema de RI desenvolvido: utilizando as palavras-
191 chave da lista de consulta pré-definidas, o algoritmo foi executado e obteve-se a lista de
192 resultado por ele gerado.

193 4. Análise dos resultados: o resultado da busca do Modelo BM25 foi
194 comparado com os documentos relevantes da lista de consulta.

195

196

197 **Discussões**

198 Este trabalho propôs a criação de uma base de dados que contivesse de forma
199 digitalizada as devidas Resoluções existentes e propôs também a realização de um
200 estudo sobre o uso do Modelo BM25 no desenvolvimento de uma Máquina de Busca
201 voltada para resoluções institucionais utilizadas pela Universidade Federal do
202 Amazonas, como um meio de minimizar o problema de acesso e recuperação de tais
203 instrumentos normativos.

204 A aplicação das métricas precisão e revocação como meio de avaliar a máquina
205 de busca desenvolvida utilizando o Modelo BM25, observou-se que a média de precisão
206 do algoritmo de 0,710, com a revocação de 0,043. Esses resultados permitiram verificar
207 que quanto maior a precisão do algoritmo de busca, menor será a sua revocação, além
208 disso, também foi possível perceber que o Modelo BM25 concentra seus resultados
209 mais relevantes próximos do topo.

210 Como resultado, pudemos identificar que a abordagem aplicada neste trabalho
211 gerou precisão de 71% e revocação 4,3%, além de apresentar tempo médio de resposta
212 de 58 milissegundos. Um fator que pode ter influenciado no resultado obtido é que
213 como a maioria das resoluções disponíveis para acesso via web estava no formato de

214 imagem, a conversão dos mesmos para texto pode ter causado uma certa perda de dados
215 importantes para a busca.

216 Entretanto, os resultados do experimento foram satisfatórios e serviram para
217 demonstrar a relevância deste trabalho e além disso, permitiu a obtenção de
218 conhecimentos importantes na área da Recuperação da Informação, principalmente a
219 assuntos referentes ao Modelo BM25 e aos Métodos de Avaliação.

220

221

222 **Agradecimentos**

223 Os autores agradecem a FAPEAM (Fundação de Amparo a Pesquisa do
224 Amazonas) pelo apoio financeiro, e a todos aqueles que contribuíram e que continuarão
225 contribuindo de alguma forma para o sucesso deste trabalho.

226

227

228

229

230

231

232

233

234

235

236

237

238 **Bibliografia Citada**

239 Ferreira, Luciene Costa. 2009. *Repositório Temático Digital dos*
240 *Instrumentos Normativos da UFBE: Acesso e Recuperação Da Informação*
241 *Através da Plataforma DSPACE*. Universidade Federal da Paraíba.

242 Garcia, Tiago Freire. 2002. *Proposta de uma Máquina de Busca Eficiente para*
243 *Documetos na Web Usando Lógica Fuzzy*. Monografia. Universidade Federal de Lavras.
244 3p.

245 Júnior, Aurélio Andrade de Menezes. 2012. *Um Método para Busca de Competências*
246 *a Partir de currículos Lattes*. Dissertação de Mestrado. Universidade Federal do
247 Amazonas. Instituto de Computação. Programa de Pós-Graduação em Informática. 12p.

248 Oliveira, Renan Rodrigues. 2010. *Recuperação Contextualizada de Documentos*
249 *Integrados pelo Protocolo OAI-PMH*. Dissertação de Mestrado. Universidade Federal
250 de Goiás. Instituto de Informática. Programa de Pós-Graduação do Instituto de
251 Informática da Universidade Federal de Goiás. 44-45p.

252

253

254

255

256

257

258

259

260

261

262

263

264

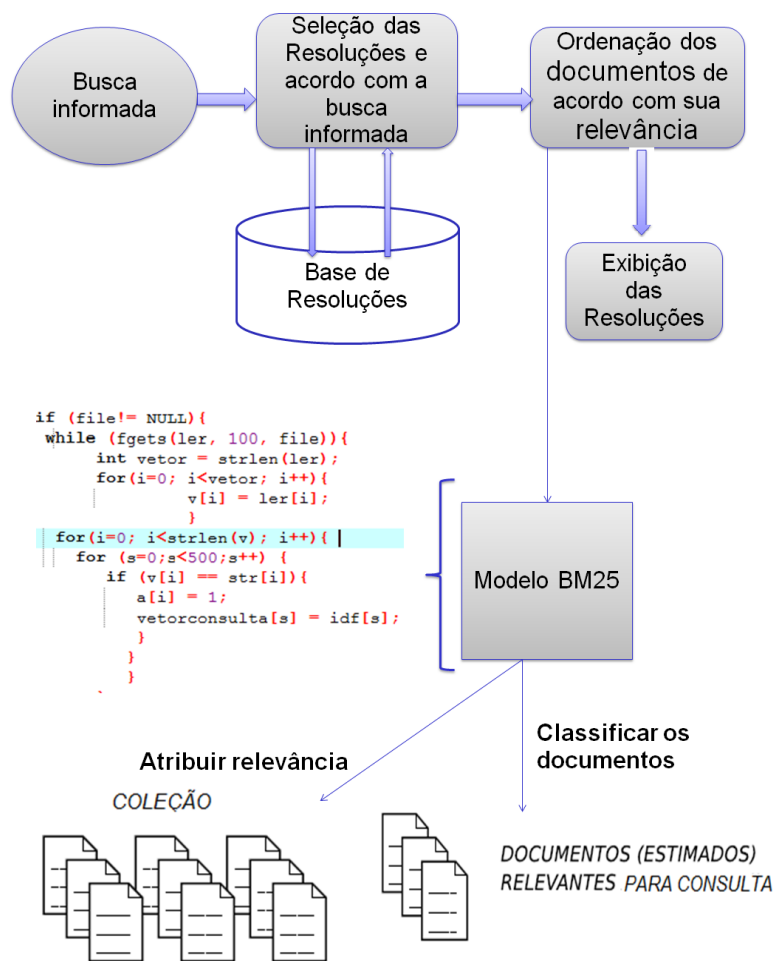
265

266

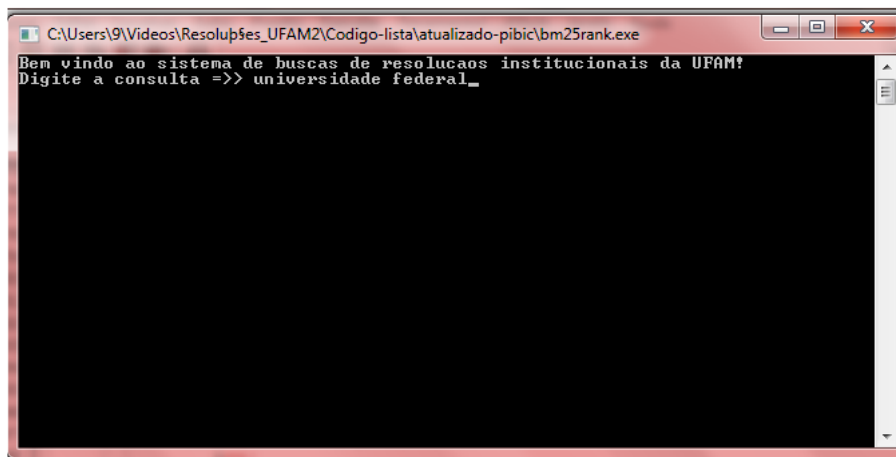
267

268

269

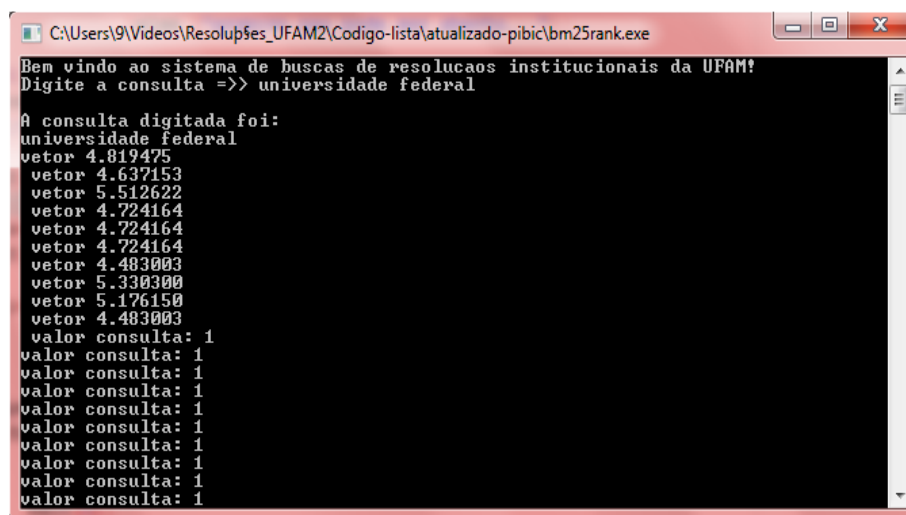


274 **Figura 1:** Estrutura do funcionamento da Máquina de busca, usando o Modelo BM25.
275



279

Figura 2: Consulta usando o Modelo BM25.



280

281 **Figura 3:** Aplicação dos graus (valores) de relevância na coleção de resoluções usando

282

o Modelo BM25.

283

284

285

286

287

288

289

290 Tabelas

291

Nº	Descrição	Ago 2014	Set	Out	No V	Dez	Jan 2015	Fev	Mar	Abr	Mai	Jun	Jul
1	Converter as resoluções levantadas de imagem para formato txt	X	X										
2	Implementar a máquina de busca utilizando o modelo bm25	X	X	X	X	X							
3	Aplicação de métodos de avaliação de desempenho para a verificação da funcionalidade da máquina de busca desenvolvida					X	X	X	X				
4	Elaboração do Resumo e Relatório Final									X	X	X	X
5	Preparação da Apresentação Final para o Congresso											X	X

292