



Francisco Benjamim Filho

**Classificação de páginas Web por
aprendizagem de múltiplas categorias latentes**

Tese de Doutorado

Tese apresentada ao Programa de Pós-graduação em Informática
do Departamento de Informática da PUC-Rio como requisito
parcial para obtenção do título de Doutor em Informática

Orientador: Prof. Ruy Luiz Milidiú

Rio de Janeiro
Março de 2012



Francisco Benjamim Filho

**Classificação de páginas Web por
aprendizagem de múltiplas categorias latentes**

Tese apresentada ao Programa de Pós-graduação em Informática do Departamento de Informática do Centro Técnico Científico da PUC-Rio como requisito parcial para obtenção do título de Doutor em Informática. Aprovada pela comissão examinadora abaixo assinada.

Prof. Ruy Luiz Milidiú

Orientador

Departamento de Informática — PUC-Rio

Prof. Marcus Vinicius Soledade Poggi de Aragão

Departamento de Informática — PUC-RIO

Prof. Raúl Pierre Rentería

Departamento de Informática — PUC-RIO

Prof. Bianca Zadrozny

UFF

Prof. Julio Cesar Duarte

IME

Prof. José Eugenio Leal

Coordenador do Centro Técnico Científico — PUC-Rio

Rio de Janeiro, 27 de Março de 2012

Todos os direitos reservados. Proibida a reprodução total ou parcial do trabalho sem autorização da universidade, do autor e do orientador.

Francisco Benjamim Filho

Graduou-se no ano de 2001 em Engenharia da Computação no Instituto Militar de Engenharia (Rio de Janeiro). Em 2005 obteve o título de Mestre em Informática pela Pontifícia Universidade Católica do Rio de Janeiro. Trabalha desde 2001 no Centro Tecnológico do Exército (Rio de Janeiro) integrando atualmente a equipe de desenvolvimento de Radio Definido por Software.

Ficha Catalográfica

Benjamim Filho, Francisco

Classificação de páginas Web por aprendizagem de múltiplas categorias latentes / Francisco Benjamim Filho; orientador: Ruy Luiz Milidiú. — Rio de Janeiro : PUC–Rio, Departamento de Informática, 2012.

v., 76 f: il. ; 29,7 cm

1. Tese (Doutorado em Informática) - Pontifícia Universidade Católica do Rio de Janeiro, Departamento de Informática.

Inclui referências bibliográficas.

1. Informática – Tese. 2. WWW. 3. Algoritmo. 4. Relevância. 5. Máquinas de busca. 6. Classificação baseada em hiperlink. I. Milidiú, Ruy Luiz. II. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática. III. Título.

CDD: 004

À minha esposa e ao meu filho.

Agradecimentos

À minha esposa Nathália por todo apoio, amor, carinho, dedicação e paciência, e pelo incentivo a seguir em frente, principalmente, nos momentos mais difíceis.

Ao meu filho Benício que me apresentou uma alegria ímpar e um amor incomensurável.

Aos demais familiares, pelo apoio.

Ao meu orientador, Professor Ruy Luiz Milidiú, pelas sugestões e correções a esse trabalho e pela atenção, incentivo e compreensão ao longo do curso.

Aos amigos que sempre estiveram presente, mesmo que de maneira virtual, me apoiando, incentivando.

À direção do CTEEx e à chefia da Divisão de Tecnologia da Informação, por permitirem a realização do curso de doutorado em tempo parcial.

À banca examinadora pelas críticas positivas ao presente trabalho.

Agradeço à PUC-Rio por me oferecer um ambiente tão propício à minha formação.

Aos meus professores pelo exemplo de seriedade, profissionalismo e por apontarem bons caminhos a serem trilhados.

Resumo

Benjamim Filho, Francisco; Milidiú, Ruy Luiz. **Classificação de páginas Web por aprendizagem de múltiplas categorias latentes**. Rio de Janeiro, 2012. 76p. Tese de Doutorado — Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

O crescimento explosivo e a acessibilidade generalizada da *World Wide Web* (*WWW*) levaram ao aumento da atividade de pesquisa na área da recuperação de informação para páginas *Web*. A *WWW* é um rico e imenso ambiente em que as páginas se assemelham a uma comunidade grande de elementos conectada através de hiperlinks em razão da semelhança entre o conteúdo das páginas, a popularidade da página, a autoridade sobre o assunto e assim por diante, sabendo-se que, em verdade, quando um autor de uma página a vincula à outra, está concebendo-a como importante para si. Por isso, a estrutura de hiperlink da *WWW* é conhecida por melhorar significativamente o desempenho das pesquisas para além do uso de estatísticas de distribuição simples de texto. Nesse sentido, a abordagem *Hyperlink Induced Topic Search* (*HITS*) introduz duas categorias básicas de páginas *Web*, *hubs* e autoridades, que revelam algumas informações semânticas ocultas a partir da estrutura de hiperlink. Em 2005, fizemos uma primeira extensão do *HITS*, denominada de *Extended Hyperlink Induced Topic Search* (*XHITS*), que inseriu duas novas categorias de páginas *Web*, quais sejam, novidades e portais. Na presente tese, revisamos o *XHITS*, transformando-o em uma generalização do *HITS*, ampliando o modelo de duas categorias para várias e apresentando um algoritmo eficiente de aprendizagem de máquina para calibrar o modelo proposto valendo-se de múltiplas categorias latentes. As descobertas aqui expostas indicam que a nova abordagem de aprendizagem fornece um modelo *XHITS* mais preciso. É importante registrar, por fim, que os experimentos realizados com a coleção *ClueWeb09* 25TB de páginas da *WWW*, baixadas em 2009, mostram que o *XHITS* pode melhorar significativamente a eficácia da pesquisa *Web* e produzir resultados comparáveis aos do *TREC 2009/2010 Web Track*, colocando-o na sexta posição, conforme os resultados publicados.

Palavras-chave

WWW ; Algoritmo ; Relevância ; Máquinas de busca ; Classificação baseada em hiperlink.

Abstract

Benjamim Filho, Francisco; Milidiú, Ruy Luiz (advisor). **Ranking of Web pages by learning multiple latent categories**. Rio de Janeiro, 2012. 76p. DSc. Thesis — Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

The rapid growth and generalized accessibility of the World Wide Web (WWW) have led to an increase in research in the field of the information retrieval for Web pages. The WWW is an immense and prodigious environment in which Web pages resemble a huge community of elements. These elements are connected via hyperlinks on the basis of similarity between the content of the pages, the popularity of a given page, the extent to which the information provided is authoritative in relation to a given field etc. In fact, when the author of a Web page links it to another, s/he is acknowledging the importance of the linked page to his/her information. As such the hyperlink structure of the WWW significantly improves research performance beyond the use of simple text distribution statistics. To this effect, the HITS approach introduces two basic categories of Web pages, hubs and authorities which uncover certain hidden semantic information using the hyperlink structure. In 2005, we made a first extension of HITS, called Extended Hyperlink Induced Topic Search (XHITS), which inserted two new categories of Web pages, which are novelties and portals. In this thesis, we revised the XHITS, transforming it into a generalization of HITS, broadening the model from two categories to various and presenting an efficient machine learning algorithm to calibrate the proposed model using multiple latent categories. The findings we set out here indicate that the new learning approach provides a more precise XHITS model. It is important to note, in closing, that experiments with the ClueWeb09 25TB collection of Web pages, downloaded in 2009, demonstrated that the XHITS is capable of significantly improving Web research efficiency and producing results comparable to those of the TREC 2009/2010 Web Track.

Keywords

WWW ; Algorithm ; Relevance ; search engines ; Classification based on hyperlink.

Sumário

1	Introdução	11
1.1	Objetivo	14
1.2	Contribuições	14
1.3	Organização do Trabalho	15
2	Trabalhos Relacionados	16
2.1	Recuperação da Informação	16
2.2	Recuperação da Informação na Web	18
2.2.1	Estrutura da Web	18
2.2.2	<i>TREC Web Tracks</i>	20
2.3	Algoritmos baseados em hiperlinks	23
2.3.1	<i>Hyperlink Induced Topic Search</i>	24
2.3.2	<i>PageRank</i>	27
2.3.3	<i>Stochastic Approach for Link-Structure Analysis</i>	30
2.4	Texto âncora e o <i>Grafo Web</i>	32
3	<i>Extended Hyperlink Induced Topic Search</i>	34
3.1	O <i>Hyperlink Induced Topic Search</i> Matricial	34
3.2	O Algoritmo <i>Extended Hyperlink Induced Topic Search</i>	36
3.3	Estrutura de Influência	37
3.4	Reforço Simétrico	38
3.5	Reforço Positivo	39
4	Aprendizagem de Máquina para o <i>Extended Hyperlink Induced Topic Search</i>	41
4.1	Aprendizagem com Múltiplas Categorias Latentes	41
4.2	Algoritmo de Aprendizagem com Múltiplas Categorias Latentes	47
4.3	Decomposição aplicada a Aprendizagem com Múltiplas Categorias Latentes	48
5	Experimentos	52
5.1	Base de Teste	52
5.2	Medidas de Avaliação dos Resultados	56
5.3	Metodologia	59
5.4	Experimentos e Resultados	63
5.4.1	Experimento com Conjunto de Treinamento <i>WT09</i>	64
5.4.2	Experimento com Conjunto de Treinamento <i>WT10</i>	66
6	Conclusão	70
6.1	Trabalhos Futuros	72

Lista de figuras

2.1	Representa a conectividade da <i>Web</i> segundo <i>Broder</i> . Extraída da URL http://www9.org/w9cdrom/160/160.html	20
2.2	Expansão de R_σ	25
2.3	<i>Hubs</i> e autoridades	27
2.4	Exemplo do cálculo dos pesos do <i>PageRank</i> sem a normalização	29
2.5	Exemplo do problema de <i>rank sink</i> no cálculo dos pesos do <i>PageRank</i>	29
2.6	Em (a) temos o grafo estendido para a consulta σ e em (b) o grafo bipartido.	31
4.1	Modelo básico de aprendizagem de máquina.	42
4.2	Modelo resumido do aprendizado de máquina baseado em correção de erro.	43
5.1	Gráfico da evolução da densidade de hiperlinks das coleções da <i>TREC</i>	54
5.2	Gráfico da evolução da erro com o aumento do conjunto de treinamento. O valor do erro está multiplicado por 100.	61
5.3	Gráfico da evolução da tempo médio para encontrar o erro mínimo por categoria.	65
5.4	Gráfico da evolução da métrica $ERR@20$ por categoria.	66
5.5	Gráfico da evolução da métrica $ERR@20$ por categoria.	68

Lista de tabelas

2.1	Resumo das informações das coleções da TREC	23
5.1	Evolução da densidade de hiperlinks das coleções da TREC	54
5.2	Resumo dos participantes da categoria A da <i>Web TREC Track 2010</i>	56
5.3	Tempo médio de execução do experimento por categoria	65
5.4	Os melhores resultados da busca <i>ad hoc</i> ordenado pela métrica <i>ERR@20</i>	67
5.5	Os melhores resultados da busca <i>ad hoc</i> ordenado pela métrica <i>P@20</i>	68
5.6	Os melhores resultados da busca <i>ad hoc</i> ordenado pela métrica <i>MAP</i>	68