



Críston Pereira de Souza

**Políticas Eficientes para Revisão de
Páginas Web**

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. Eduardo Sany Laber

Rio de Janeiro
Maio de 2010



Críston Pereira de Souza

Políticas Eficientes para Revisão de Páginas Web

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. Eduardo Sany Laber

Orientador
Departamento de Informática — PUC-Rio

Prof. Artur Alves Pessoa

Departamento de Engenharia de Produção — UFF

Prof. Marcus Felipe Fontoura

Yahoo! Research

Prof. Marcus Vinicius Soledade Poggi de Aragão

Departamento de Informática — PUC-Rio

Prof. Raúl Pierre Rentería

Departamento de Informática — PUC-Rio

Prof. José Eugenio Leal

Coordenador Setorial do Centro Técnico Científico — PUC-Rio

Rio de Janeiro, 25 de Maio de 2010

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

Críston Pereira de Souza

Bacharel em Ciência da Computação pela Universidade Federal da Bahia e Mestre em Informática pela Pontifícia Universidade Católica do Rio de Janeiro. Tem experiência profissional como analista de sistemas, e atuou em projetos de pesquisa nas áreas de Otimização Combinatória, Aprendizado de Máquina e Sistemas Distribuídos.

Ficha Catalográfica

Souza, Críston Pereira de

Políticas eficientes para revisitação de páginas Web / Críston Pereira de Souza; orientador: Eduardo Sany Laber. — Rio de Janeiro : PUC, Departamento de Informática, 2010.

v., 84 f: il. ; 29,7 cm

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

Inclui referências bibliográficas.

1. Informática – Teses. 2. Máquinas de Busca. 3. Coletores Web. 4. Internet. 5. Algoritmos Aproximativos. 6. Simulação. I. Laber, Eduardo Sany. II. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática. III. Título.

Agradecimentos

Ao meu orientador Professor Eduardo Laber pelo tempo dedicado durante todos estes anos, ajudando com muitas idéias e sua forma objetiva de trabalhar.

Ao CNPq, à FAPERJ e à PUC–Rio, pelos auxílios concedidos que permitiram realizar este trabalho.

À minha mãe Silene, que não mediou esforços para a educação dos filhos.

À Andréia pelo companheirismo e tantas ajudas importantes, sem os quais tudo seria mais difícil.

Aos colegas do LEARN/PUC–Rio pelos momentos de descontração, em especial os alunos Eduardo e Caio que colaboraram com esta pesquisa.

Aos funcionários do departamento de informática pela orientação e auxílio nas questões relacionadas ao curso.

Aos Professores Macêdo, Aline, Flávio e George do LaSiD/UFBA, por despertar em seus alunos o interesse pela pesquisa científica.

Resumo

Souza, Críston Pereira de; Laber, Eduardo Sany. **Políticas Eficientes para Revisão de Páginas Web.** Rio de Janeiro, 2010. 84p. Tese de Doutorado — Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Uma máquina de busca precisa constantemente revisitar páginas *Web* para manter seu repositório local atualizado. Uma política de revisitação deve ser empregada para construir um escalonamento de revisitações que mantenha o repositório o mais atualizado possível utilizando os recursos disponíveis. Para evitar sobrecarga de servidores *Web*, a política de revisitação deve respeitar um tempo mínimo entre requisições consecutivas a um mesmo servidor. Esta regra é chamada restrição de *politeness*. Devido ao porte do problema, consideramos que uma política de revisitação é eficiente se o tempo médio para escalar uma revisitação é sublinear no número de páginas do repositório. Neste sentido, quando a restrição de *politeness* é considerada, não conhecemos política eficiente com garantia teórica de qualidade. Nesta pesquisa investigamos três políticas eficientes que respeitam a restrição de *politeness*, chamadas MERGE, RANDOM e DELAYED. Fornecemos fatores de aproximação para o nível de atualização do repositório quando empregamos as políticas MERGE ou RANDOM. Demonstramos que 0,77 é um limite inferior para este fator de aproximação quando empregamos a política RANDOM, e apresentamos uma conjectura de que 0,927 é um limite inferior para este fator de aproximação quando empregamos a política MERGE. As políticas também são avaliadas através da simulação da execução destas políticas para manter o nível de atualização de um repositório contendo 14,5 milhões de páginas *Web*. Um repositório contendo artigos da Wikipedia também é utilizado nos experimentos, onde podemos observar que a política MERGE apresenta melhores resultados que uma estratégia gulosa natural para este repositório. A principal conclusão desta pesquisa é que existem políticas simples e eficientes para o problema de revisitação de páginas *Web*, que perdem pouco em termos do nível de atualização do repositório mesmo quando consideramos a restrição de *politeness*.

Palavras-chave

Máquinas de Busca; Coletores Web; Internet; Algoritmos Aproximativos; Simulação.

Abstract

Souza, Críston Pereira de; Laber, Eduardo Sany (Advisor).

Efficient Web Page Refresh Policies. Rio de Janeiro, 2010.

84p. DSc Thesis — Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

A search engine needs to continuously revisit web pages in order to keep its local repository up-to-date. A page revisiting schedule must be defined to keep the repository up-to-date using the available resources. In order to avoid web server overload, the revisiting policy must respect a minimum amount of time between consecutive requests to the same server. This rule is called politeness constraint. Due to the large number of web pages, we consider that a revisiting policy is efficient when the mean time to schedule a revisit is sublinear on the number of pages in the repository. Therefore, when the politeness constraint is considered, there are no existing efficient policies with theoretical quality guarantees. We investigate three efficient policies that respect the politeness constraint, called MERGE, RANDOM and DELAYED. We provide approximation factors for the repository's up-to-date level for the MERGE and RANDOM policies. Based on these approximation factors, we devise a 0.77 lower bound for the approximation factor provided by the RANDOM policy and we present a conjecture that 0.927 is a lower bound for the approximation factor provided by the MERGE policy. We evaluate these policies through simulation experiments which try to keep a repository with 14.5 million web pages up-to-date. Additional experiments based on a repository with Wikipedia's articles concluded that the MERGE policy provides better results than a natural greedy strategy. The main conclusion of this research is that there are simple and efficient policies that can be applied to this problem, even when the politeness constraint must be respected, resulting in a small loss of repository's up-to-date level.

Keywords

Search Engines; Web Crawlers; Internet; Approximation Algorithms; Simulation.

Sumário

| | | |
|-----|---|-----------|
| 1 | Introdução | 11 |
| 1.1 | Objetivo e Justificativa | 16 |
| 1.2 | Metodologia | 17 |
| 1.3 | Contribuições | 17 |
| 1.4 | Organização da Tese | 20 |
| 2 | Fundamentação Teórica e Trabalhos Relacionados | 21 |
| 2.1 | Medidas para o Nível de Atualização do Repositório | 21 |
| 2.2 | Restrição de Politeness | 27 |
| 2.3 | Políticas de Revisão | 28 |
| 2.4 | Modelo Adotado nesta Tese | 34 |
| 3 | Caracterização do Repositório Experimental | 35 |
| 3.1 | Construção e Monitoramento | 36 |
| 3.2 | Número de Páginas por Servidor | 38 |
| 3.3 | Modificações segundo um Processo de Poisson | 39 |
| 3.4 | Distribuição das Taxas de Modificação | 40 |
| 4 | Política de Tempo Igualmente Espaçada por Página | 41 |
| 4.1 | Política de Tempo Igualmente Espaçada por Página | 42 |
| 4.2 | Política de Tempo DELAYED | 44 |
| 4.3 | Limites Superiores para o Freshness do Repositório | 45 |
| 4.4 | Alocação de Recursos considerando a Restrição de Politeness | 47 |
| 4.5 | Resultados Experimentais | 50 |
| 5 | Política de Tempo Igualmente Espaçada por Servidor | 52 |
| 5.1 | Política de Tempo Igualmente Espaçada por Servidor | 54 |
| 5.2 | Política de Seleção de Páginas RANDOM | 55 |
| 5.3 | Política de Seleção de Páginas MERGE | 58 |
| 5.4 | Uso Efetivo do Canal de Comunicação | 63 |
| 5.5 | Resultados Experimentais | 66 |
| 6 | Experimentos com um Repositório de Artigos da Wikipedia | 68 |
| 6.1 | Repositório WIKIPEDIA | 69 |
| 6.2 | Políticas de Revisão | 72 |
| 6.3 | Experimentos | 75 |
| 7 | Conclusões | 78 |
| 7.1 | Trabalhos Futuros | 79 |
| | Referências Bibliográficas | 81 |

Lista de figuras

| | |
|---|---|
| <p>1.1 Principais elementos do problema de revisitação de páginas Web.</p> <p>2.1 <i>Freshness</i> e <i>age</i> de uma página. O <i>freshness</i> vale 1 quando a página está atualizada, ou 0 caso contrário. O <i>age</i> vale 0 quando a página está atualizada, ou vale o tempo desde a última modificação quando está desatualizada.</p> <p>3.1 Concentração das páginas nos servidores com mais páginas (repositório WEBBASE).</p> <p>3.2 Concentração de páginas por servidor. Para cada quantidade x de páginas por servidor, o gráfico fornece a proporção das páginas hospedadas em servidores com até x páginas (repositório WEBBASE).</p> <p>3.3 Histograma das durações dos intervalos entre modificações das páginas que se modificam em média 1 vez a cada 10 dias (repositório WEBBASE).</p> <p>3.4 Distribuição das taxas de modificação das páginas - lei de potências (repositório WEBBASE).</p> <p>4.1 Duas páginas de um mesmo servidor com atualizações igualmente espaçadas. P é o tempo mínimo permitido entre requisições a um servidor. Pode ocorrer violação da restrição de <i>politeness</i>.</p> <p>4.2 Política de tempo <i>DELAYED</i>: a $(i+1)$-ésima requisição ao servidor é escalonada para o instante t_{i+1} e viola a restrição de <i>politeness</i>, sendo portanto atrasada para o instante $t_i + P$, onde P é o tempo mínimo permitido entre requisições a um servidor.</p> <p>4.3 Pseudo-código da política de tempo <i>DELAYED</i>.</p> <p>4.4 Algoritmo para a alocação de recursos <i>OPT_POLITE</i>.</p> <p>4.5 <i>Freshness</i> do repositório WEBBASE fornecido pela política <i>DELAYED</i> durante 4 anos de operação do <i>crawler</i>, para a frequência total C de revisitação igual a 1%, 10% e 90% da frequência máxima permitida pela restrição de <i>politeness</i>.</p> <p>5.1 Pseudo-código da política de tempo igualmente espaçada por servidor.</p> <p>5.2 Fator de aproximação para o <i>freshness</i> de uma página i modificada com taxa λ_i, e revisitada com frequência f_i pela política <i>RANDOM</i>.</p> <p>5.3 Exemplo de aplicação da política de seleção de páginas <i>MERGE</i>. Duas páginas de um mesmo servidor são revisitadas com frequências f_1 e f_2. A frequência de requisições ao servidor vale $f = f_1 + f_2$. Os deslocamentos d_1, d_2 e d são uniformemente distribuídos nos intervalos $[0, f_1^{-1})$, $[0, f_2^{-1})$ e $[0, f^{-1})$, respectivamente.</p> <p>5.4 Pseudo-código da política de seleção de páginas <i>MERGE</i>.</p> <p>5.5 Revisitações igualmente espaçadas que ocorrem em um intervalo I_t com duração t.</p> | <p style="text-align: right;">12</p> <p style="text-align: right;">22</p> <p style="text-align: right;">37</p> <p style="text-align: right;">38</p> <p style="text-align: right;">39</p> <p style="text-align: right;">40</p> <p style="text-align: right;">43</p> <p style="text-align: right;">44</p> <p style="text-align: right;">45</p> <p style="text-align: right;">48</p> <p style="text-align: right;">51</p> <p style="text-align: right;">54</p> <p style="text-align: right;">56</p> <p style="text-align: right;">59</p> <p style="text-align: right;">60</p> <p style="text-align: right;">60</p> |
|---|---|

| | | |
|-----|--|----|
| 5.6 | Limite inferior fornecido pela Equação (5-9) para o fator de aproximação do <i>freshness</i> do repositório quando empregamos a política de seleção de páginas MERGE. | 63 |
| 5.7 | Limite superior para a probabilidade de ocorrer pelo menos x requisições acima da quantidade esperada em um intervalo de tempo arbitrário. As requisições são igualmente espaçadas para cada um dos n elementos. | 65 |
| 5.8 | <i>Freshness</i> do repositório WEBBASE fornecido pela política DELAYED durante 6 anos de operação do <i>crawler</i> , para a frequência total C de revisitação igual a 1%, 10% e 90% da frequência máxima permitida pela restrição de <i>politeness</i> . | 67 |
| 6.1 | Distribuição acumulada empírica do tempo entre modificações consecutivas (em minutos) de 10 artigos escolhido aleatoriamente dentre os artigos do repositório WIKIPEDIA. Os gráficos mostram também as distribuições Pareto, Gama e Weibull ajustadas com estimadores de máxima verossimilhança. | 72 |
| 6.2 | Quando um artigo i é avaliado no instante t , conhecemos o instante $u_i(t)$ da última revisitação deste artigo, e o instante $b_i(t)$ da sua última modificação antes de $u_i(t)$. O instante $A_i(t)$ da primeira modificação depois de t é uma variável aleatória. | 74 |
| 6.3 | <i>Freshness</i> do repositório WIKIPEDIA após a execução das políticas MERGE, MERGE2, RANDOM e GREEDY, variando a frequência de requisições ao servidor. | 76 |

As far as the laws of mathematics refer to reality, they are not certain; and as far as they are certain, they do not refer to reality.

Albert Einstein, *Sidelights on Relativity*.