

**Phương pháp biểu diễn ngữ nghĩa lân cận siêu liên kết trong máy
tìm kiếm VietSeek**

Đặng Tiểu Hùng

Người hướng dẫn: TS. Hà Quang Thụy

MỤC LỤC

PHẦN MỞ ĐẦU.....	5
CHƯƠNG 1. TỔNG QUAN VỀ TÌM KIẾM THÔNG TIN TRÊN WEB.....	7
1.1 Giới thiệu về tìm kiếm thông tin	7
1.2 Bài toán tìm kiếm thông tin	7
1.2.1 Giai đoạn 1: Thu thập và phân tích thông tin	Error! Bookmark not defined.
1.2.2 Giai đoạn 2: Xử lý câu hỏi và trả lời	Error! Bookmark not defined.
1.3 Mô hình biểu diễn thông tin của văn bản	Error! Bookmark not defined.
1.3.1 Mô hình biểu diễn thông tin theo từ khoá .	Error! Bookmark not defined.
1.3.2 Mô hình biểu diễn thông tin theo nội dung	Error! Bookmark not defined.
1.4 Phân tích cú pháp và ngữ nghĩa.....	Error! Bookmark not defined.
1.5 Phân lớp văn bản.....	Error! Bookmark not defined.
1.6 Phân cụm văn bản.....	Error! Bookmark not defined.
1.7 Khai thác thông tin cấu trúc web.....	Error! Bookmark not defined.
1.8 Khai thác thông tin sử dụng web.....	Error! Bookmark not defined.
CHƯƠNG 2. PHƯƠNG PHÁP BIỂU DIỄN TRANG WEB THEO NGỆ NGHĨA LÂN CẬN SIÊU LIÊN KẾT.....	ERROR! BOOKMARK NOT DEFINED.
2.1 Giới thiệu.....	Error! Bookmark not defined.
2.2 Phương pháp đánh giá chất lượng độ đo tự động	Error! Bookmark not defined.
2.2.1 Chọn phương pháp đánh giá	Error! Bookmark not defined.
2.2.2 Xác định thứ tự nền trong ODP..	Error! Bookmark not defined.

- 2.2.3 So sánh sự t- ơng quan giữa các tập thứ tự. **Error! Bookmark not defined.**
- 2.2.4 Miền của tập thứ tự..... **Error! Bookmark not defined.**
- 2.3 Định nghĩa mô hình vector biểu diễn thông tin văn bản **Error! Bookmark not defined.**
 - 2.3.1 Vector biểu diễn thông tin văn bản **Error! Bookmark not defined.**
 - 2.3.2 Lựa chọn từ khoá biểu diễn **Error! Bookmark not defined.**
 - 2.3.3 L- ọc bớt từ khoá..... **Error! Bookmark not defined.**
 - 2.3.4 Xác định trọng số của từ khoá.... **Error! Bookmark not defined.**
- 2.4 Định nghĩa độ đo t- ơng tự..... **Error! Bookmark not defined.**
- 2.5 Đánh giá chất l- ợng xếp hạng đối với mỗi ph- ơng pháp xây dựng vector **Error! Bookmark not defined.**
 - 2.5.1 Đánh giá chất l- ợng đối với cách chọn từ khoá **Error! Bookmark not defined.**
 - 2.5.2 Đánh giá chất l- ợng đối với cách chuẩn hoá trọng số từ khoá **Error! Bookmark not defined.**
 - 2.5.3 Đánh giá chất l- ợng đối với ph- ơng pháp l- ọc bớt từ khoá **Error! Bookmark not defined.**
- 2.6 Thiết kế các thuật toán tìm kiếm theo mô hình vector **Error! Bookmark not defined.**

CHƯƠNG 3. MÁY TÌM KIẾM VIETSEEK VÀ THỰC NGHIỆM THUẬT TOÁN TÌM KIẾM THEO NGHỆ NGHĨA LÂN CẬN SIÊU LIÊN KẾT. **Error! Bookmark not defined.**

- 3.1 Máy tìm kiếm VietSeek..... **Error! Bookmark not defined.**
 - 3.1.1 Các đặc điểm cơ bản của VietSeek **Error! Bookmark not defined.**
 - 3.1.2 Cơ sở dữ liệu của VietSeek..... **Error! Bookmark not defined.**

3.2 Đề xuất thuật toán tìm kiếm mới cho máy tìm kiếm VietSeek ..**Error!**

Bookmark not defined.

3.2.1 Những cơ sở để đề xuất thuật toán**Error! Bookmark not defined.**

3.2.2 Xây dựng các thuật toán áp dụng cho máy tìm kiếm VietSeek**Error!**

Bookmark not defined.

3.2.3 Kết quả thực hiện..... **Error! Bookmark not defined.**

PHỤ LỤC LUẬN..... **ERROR! BOOKMARK NOT DEFINED.**

TÀI LIỆU THAM KHẢO..... 10

PHỤ LỤC **ERROR! BOOKMARK NOT DEFINED.**

LỜI CẢM ƠN

Tôi xin bày tỏ lòng kính trọng và biết ơn tới các thầy giáo, cô giáo khoa Công nghệ tr- ờng Đại học Quốc gia Hà Nội đã dìu dắt tôi trong suốt quá trình học tập và nghiên cứu, cũng nh- ớng góp những ý kiến quý báu cho luận văn.

Đặc biệt tôi xin bày tỏ lòng kính trọng và biết ơn sâu sắc Thầy giáo Tiến sĩ Hà Quang Thụy cùng gia đình đã tận tình, dành nhiều thời gian h- ớng dẫn, động viên, khích lệ cho tôi hoàn thành luận văn này.

Tôi xin bày tỏ lòng biết ơn sâu sắc tới gia đình, bạn bè và đồng nghiệp đã tạo điều kiện thuận lợi giúp đỡ cũng nh- ớng có nhiều ý kiến đóng góp bổ ích cho luận văn.

Tôi xin kính chúc các thầy giáo, cô giáo cùng gia đình mạnh khoẻ, hạnh phúc; Tiếp tục sự nghiệp đào tạo cho các thế hệ học sinh, sinh viên đạt đ- ợc nhiều thành công hơn nữa trên con đ- ờng học tập và nghiên cứu khoa học.

Tôi xin chúc các bạn bè, đồng nghiệp mạnh khoẻ, thành công; áp dụng hiệu quả và sáng tạo các kiến thức đ- ợc học vào thực tiễn.

Xin trân trọng cảm ơn.

Hà Nội ngày 25/03/2004

Học viên

Đặng Tiểu Hùng

PHẦN MỞ ĐẦU

Cùng với sự phát triển mạnh mẽ của Internet là một số lượng khổng lồ dữ liệu được phát sinh, tuy nhiên (theo thông tin từ tập đoàn Oracle) thì khoảng 90% dữ liệu ở dạng phi cấu trúc hoặc nửa cấu trúc. Trong khi nhu cầu khai thác, tìm kiếm thông tin một cách chính xác trên internet đã ngày càng trở nên bức thiết hơn, do đó xuất hiện các hệ tìm kiếm theo từ khoá (cụm từ khoá) như Yahoo, Google ... Tuy nhiên việc tìm kiếm theo từ khoá vẫn chưa đủ để giúp người sử dụng nhanh chóng tìm được trang Web cần thiết vì số lượng kết quả trả lại rất lớn và nhiều khi chỉ là các trang Web ít có liên quan. Vì vậy các hệ thống tìm kiếm ngày càng được cải tiến để ngày càng thông minh hơn. Xuất hiện những hệ thống tối mục tiêu cụ thể như tra cứu thông tin về các chủ đề y tế, giáo dục, luật pháp, âm nhạc ... Tuy vậy, việc nghiên cứu các giải pháp để tìm được một trang thông tin theo một nội dung nào đó sát với yêu cầu người sử dụng thì vẫn còn nhiều hạn chế. Đã có nhiều mô hình tìm kiếm được đề xuất, song những mô hình lý thuyết về mặt lý thuyết thì lại chưa có tính khả thi khi cài đặt. Do đó, trong các hệ tìm kiếm, người ta tìm cách cải tiến các phương pháp đơn giản có sẵn để có áp dụng trong thực tế. Luận văn này hướng tới việc nghiên cứu, phân tích, đánh giá kết quả của một số thuật toán tìm kiếm theo nội dung, từ đó đề xuất một phương án cải tiến để nâng cao hiệu quả về tính chính xác của nội dung cũng như về tốc độ.

Từ việc tìm hiểu, đánh giá và phân tích ưu, nhược điểm của các phương pháp tiếp cận khác nhau, dựa theo mục tiêu trên ý tưởng nâng cao hiệu quả tìm kiếm, luận văn đề xuất giải pháp thực hiện “*Phương pháp biểu diễn ngữ nghĩa lân cận siêu liên kết cho máy tìm kiếm VietSeek*”.

Nội dung của luận văn được định hướng vào các vấn đề sau:

1. Mô hình toán học biểu diễn trang văn bản Web.
2. Khái quát các phương pháp tiếp cận trong tìm kiếm trang Web có nội dung tự động. Đánh giá ưu điểm và nhược điểm của mỗi phương pháp được khảo sát.

3. Đề xuất phương pháp kết hợp để đạtnâng cao hiệu quả cao hơn trong tìm kiếm trang Web có nội dung tương tự.

Luận văn bao gồm Phần mở đầu, ba chương nội dung và Phần kết luận vớimà nội dung các chương được trình bày như dưới đây.

Chương 1 với tiêu đề là *Tổng quan về các phương pháp biểu diễn và tìm kiếm thông tin trên web* giới thiệu khái quát về các phương pháp biểu diễn và tìm kiếm trên web.

Tiêu đề của chương 2 là *Phương pháp biểu diễn trang web theo ngữ nghĩa lân cận siêu liên kết*. Chương này sẽ trình bày cơ sở, nội dung của phương pháp đọc đề xuất cũng nhưtrình bày đánh giá phương pháp được đề xuất với các phương pháp khác. Luận văn cũng trình bày chi tiếtđánh các lựa chọn được đề xuất trong mỗi bước của phương pháp, từ đó chọn ra giải pháp tốt nhất.

Chương 3 *Máy tìm kiếm VietSeek và thử nghiệm Thuật toán tìm kiếm theo ngữ nghĩa lân cận siêu liên kết* giới thiệu kiến trúc logic của máy tìm kiếm VietSeek, thiết kế logic về dữ liệu theo biểu diễn vector và thuật toán tìm kiếm theo nội dung trên cơ sở biểu diễn trang web do luận văn đề xuất. Trongch chương này cũng đề xuất những cải tiến khi áp dụng vào thực tế để nâng cao hiệu suất thực hiện của phương pháp biểu diễn.

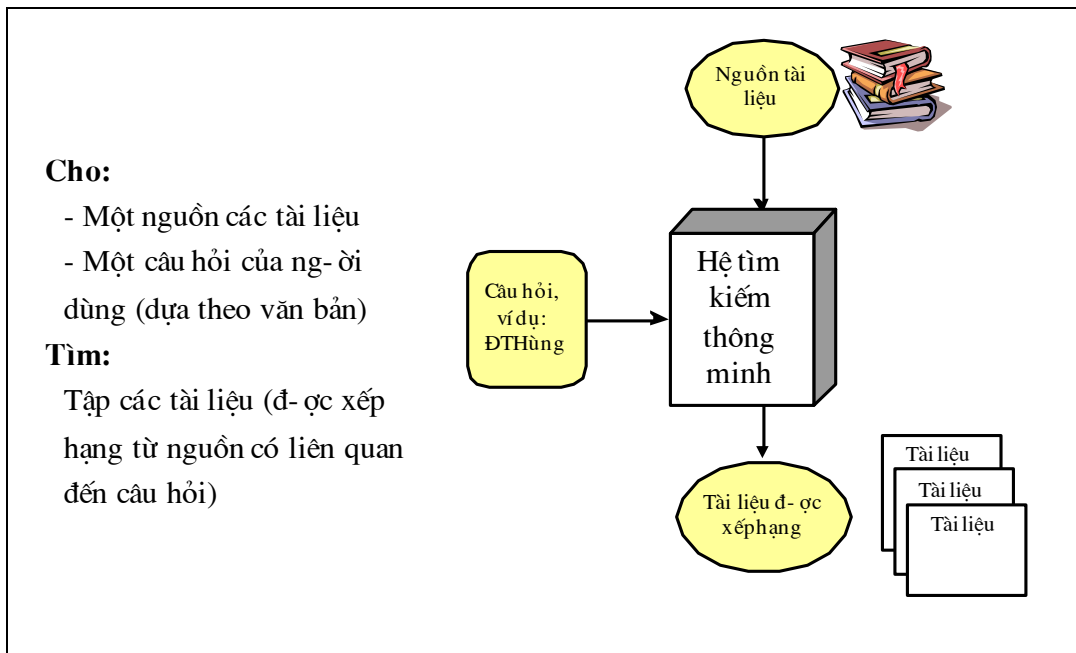
Phần kết luận tổng hợp những kết quả nghiên cứu chính của luận văn, và chỉ ra một số hạn chế của luận văn. Đồng thời luận văn cũng đề xuất một số hướng nghiên cứu cụ thể tiếp theo của luận văn.

Phần phụ lục bổ sung một số thông tin về chi tiết về việc áp dụng thuật toán cho máy tìm kiếm VietSeek như sơ đồ khối một số module cần bổ sung chức năng, những lệnh bổ sung vào cơ sở dữ liệu của VietSeek.

CHƯƠNG 1. TỔNG QUAN VỀ TÌM KIẾM THÔNG TIN TRÊN WEB

1.1 Giới thiệu về tìm kiếm thông tin

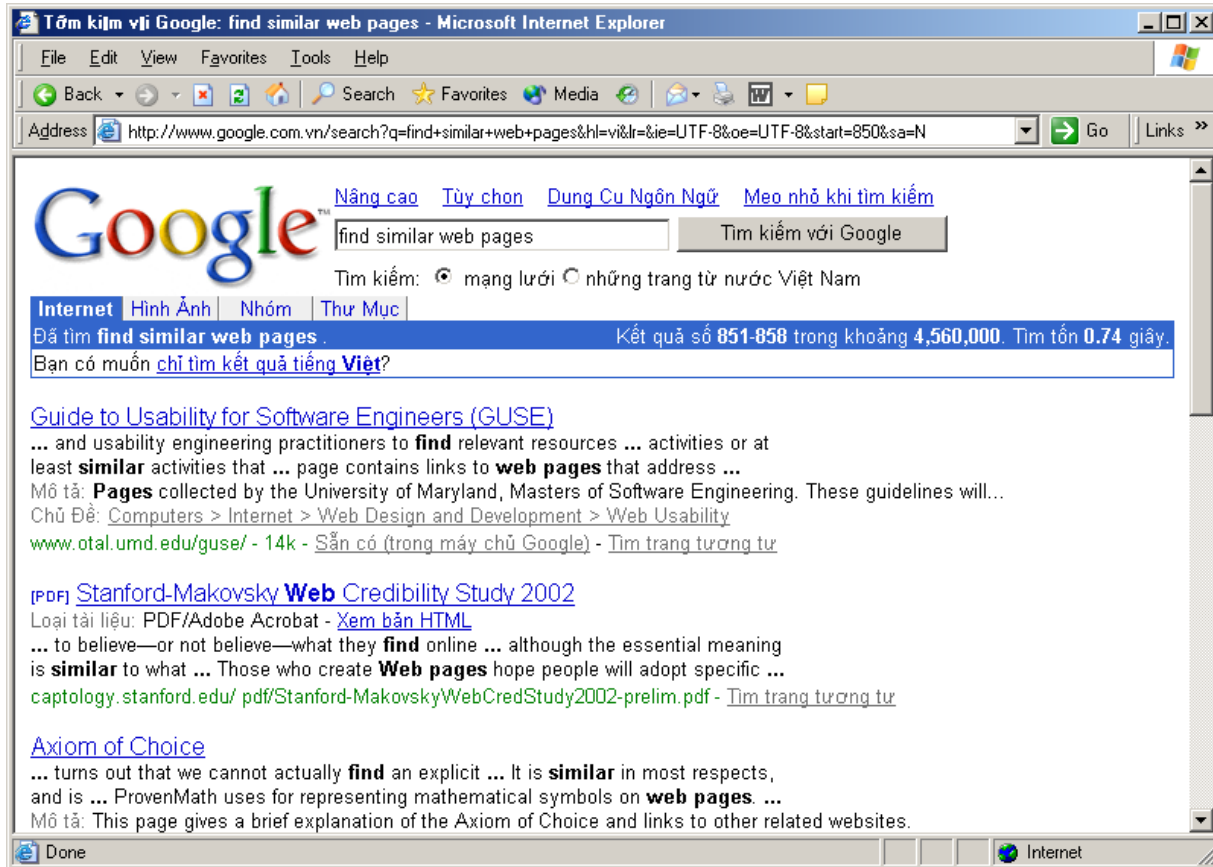
Khai phá dữ liệu thông tin trên web (web mining) là quá trình khảo sát và phân tích dữ liệu web một cách tự động hoặc bán tự động để phát hiện ra thông tin. Từ thông tin đã được khai phá, và tìm kiếm thông tin (Information Retrieval) trên web là phương pháp để truy cập một cách hiệu quả nhất đến thông tin mà người dùng quan tâm, đó có thể là kỳ vọng cung cấp một tập hợp nhỏ các văn bản gần nhất đến lĩnh vực hoặc chủ đề mà người dùng mong muốn tiếp cận.



Hình 1. Tìm kiếm thông tin

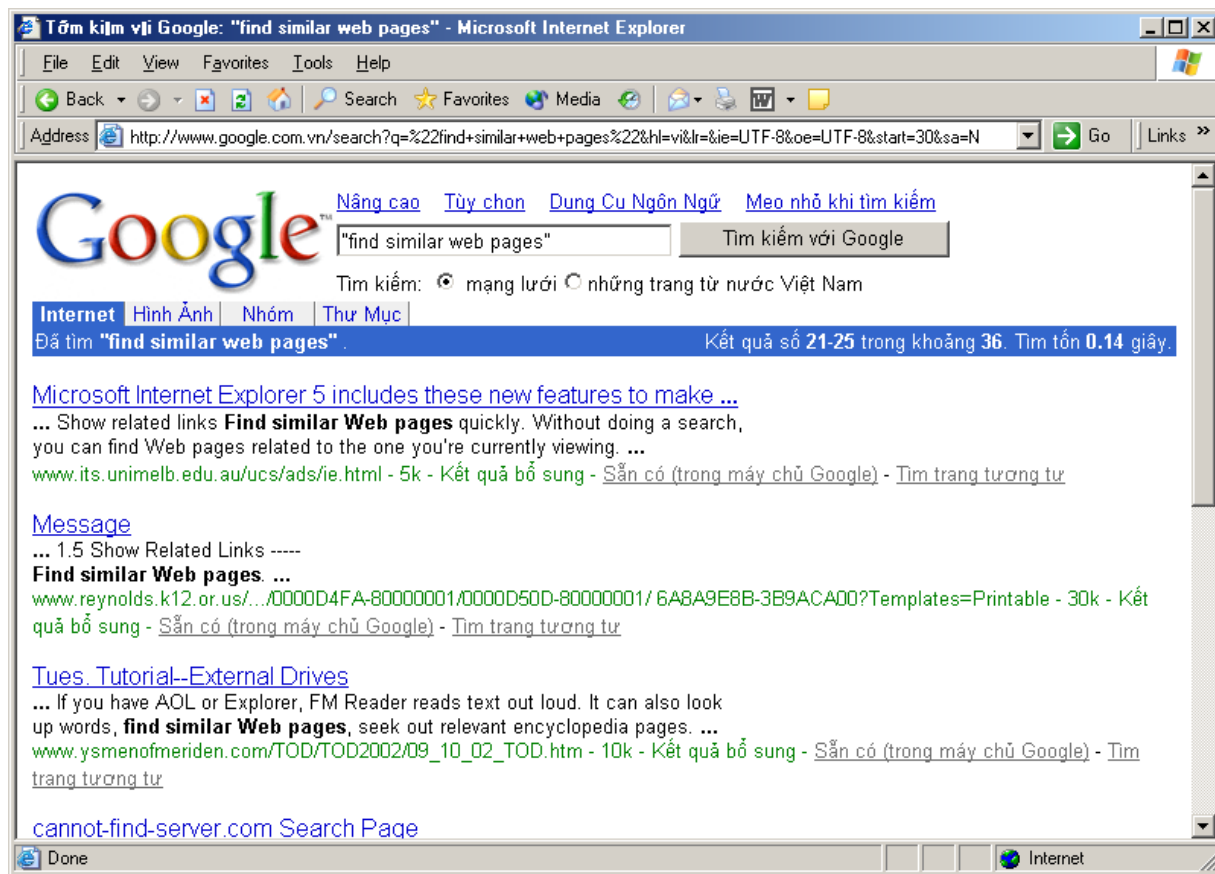
1.2 Bài toán tìm kiếm thông tin

Có 2 bài toán cơ bản trong tìm kiếm thông tin là tìm kiếm theo từ khóa và tìm kiếm theo nội dung. Bài toán tìm kiếm theo từ khóa là bài toán tìm kiếm thông tin theo các từ khóa do người dùng cung cấp [1]. Hệ tìm kiếm sẽ trả về cho người dùng các trang web có chứa những từ khóa trong câu hỏi. Tuy vậy, với số lượng khổng lồ các trang web trên internet hiện nay thì số lượng kết quả tìm đ-ợc theo từ khóa là quá lớn. Ví dụ nếu tìm các trang web có từ khóa *find similar web page* thì cho kết quả 858 trang web.



Hình 2. Tìm kiếm thông tin theo từ khoá

Bằng cách tìm kiếm theo cụm từ khoá thì số lượng kết quả trả về chính xác hơn, số kết quả trả về là 25 trang web.



Hình 3. Tìm kiếm thông tin theo cụm từ khoá

TÀI LIỆU THAM KHẢO

Tiếng Việt

- [1]. Phạm Thanh Nam (2003). *Một số giải pháp cho bài toán tìm kiếm trong cơ sở dữ liệu Hypertext*. Luận văn thạc sĩ Công nghệ thông tin - Đại học Quốc gia Hà Nội.
- [2]. Phạm Thanh Nam, Bùi Quang Minh, Hà Quang Thuy (2004). *Giải pháp tìm kiếm trang Web t-ong tự trong máy tìm kiếm VietSeek*. Tạp chí Tin học và Điều khiển học (nhận đăng 1-2004).
- [3]. Đoàn Sơn (2002). *Các ph-ong pháp biểu diễn và ứng dụng trong khai phá dữ liệu văn bản*. Luận văn thạc sĩ Công nghệ thông tin - Đại học Quốc gia Hà Nội.

Tiếng Anh

- [4]. J. Dean and M. Henzinger (1999). *Finding Related Pages in the World Wide Web*. Proceedings of WWW8, 1999.
- [5]. L. A. Goodman and W. H. Kruskal (1954). *Measures of association for cross classifications*. J. of Amer. Stat. Assoc, 1954.
- [6]. T.H. Haveliwala, A. Gionis, and P. Indyk (2000). *Scalable Techniques for Clustering the Web*. Informal Proceedings of the International Workshop on the Web and Databases, WebDB, 2000.
- [7]. J. Hirai, S. Raghavan, H. Garcia-Molina, and A. Paepcke (2000). *WebBase: A Repository of Web Pages*. Proceedings of WWW9, 2000.
- [8]. A.K. Jain, M. Narasimha Murty, and P.J. Flynn (1999). *Data clustering: A review* ACM Computing Surveys, 31(3), 1999.
- [9]. H. P. Luhn. *The Automatic Creation of Literature Abstracts*. IBM Journal of Research and Development, 2:159-165, 1958.
- [10]. Nguyen Ngoc Minh, Nguyen Tri Thanh, Ha Quang Thuy, Luong Song Van, Nguyen Thi Van (2001). *A Knowledge Discovery Model in Full-text*

Databases. Proceedings of the First Workshop of International Joint Research: "Parallel Computing, Data Mining and Optical Networks". March 7, 2001, Japan Advanced Institute of Science and Technology (JAIST), Tatsunokuchi, Japan, 59-68.

- [11]. M. Porter (1980). *An Algorithm for Suffix Stripping*. Program: Automated Library and Information Systems, 14(3):130-137, 1980.
- [12]. G. Salton and M.J. McGill (1983). *Introduction to Modern Information Retrieval*. McGraw-Hill, 1983.
- [13]. Sen Slattery (2002). *Hypertext Classification*. Doctoral dissertation (CMU-CS-02-142). School of Computer Science. Carnegie Mellon University.
- [14]. S. Siegel and N. J. Castellan (1988). *Nonparametric Statistics for the Behavioral Sciences*. McGraw-Hill, 1988.
- [15]. M. Steinbach, G. Karypis, and V. Kumar (2000). *A comparison of document clustering techniques*. TextMining Workshop, KDD, 2000.
- [16]. Taher H. Haveliwala, Aristides Gionis, Dan Klein, Piotr Indyk (2002). *Evaluating Strategies for Similarity Search on the Web*. WWW2002 - USA.
- [17]. BBC. <http://www.bbc.com>.
- [18]. CNN <http://www.cnn.com>.
- [19]. Open Directory Project (ODP). <http://www.dmoz.com>.
- [20]. Web page www.InfoWorld.com (Theo công bố ngày 17/02/2004 thì trong kho dữ liệu của Google đã có 4,28 tỷ trang web, 880 triệu hình ảnh và 845 triệu thông điệp Internet. Mạng thông tin đang tăng nhanh gần đây là các trang web liên quan đến sách, bao gồm các ch-ong đầu, phần phê bình, tham khảo. Hệ thống thông tin này đ-ợc Google truy xuất qua dịch vụ Google Print đang đ-ợc vận hành thử nghiệm. Số liệu thống kê gần đây của Google là 3,3 tỷ trang web đ-ợc kết nối vào tháng 8-2003, là 400 triệu hình ảnh vào tháng 11/2002).
- [21]. Yahoo! <http://www.yahoo.com/>.