

ĐẠI HỌC THÁI NGUYÊN
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

NGUYỄN DANH HÙNG

**NGHIÊN CỨU XÂY DỰNG HỆ THỐNG TỔNG HỢP,
PHÂN LOẠI THÔNG TIN TỰ ĐỘNG TRÊN WEB**

Chuyên ngành: Khoa học máy tính

Mã số : 60.48.0101

LUẬN VĂN THẠC SĨ CÔNG NGHỆ THÔNG TIN

Người hướng dẫn khoa học: PGS.TS ĐOÀN VĂN BAN

Thái nguyên – Năm 2014

MỤC LỤC

MỤC LỤC	i
DANH MỤC CÁC KÝ HIỆU, CÁC CHỮ VIẾT TẮT	iv
DANH MỤC CÁC BẢNG	v
DANH MỤC CÁC HÌNH	vi
MỞ ĐẦU	1
CHƯƠNG 1: KHAI PHÁ DỮ LIỆU.....	4
1.1. Khai phá dữ liệu	4
1.1.1. Giới thiệu khai phá dữ liệu	4
1.1.2. Quá trình khai phá dữ liệu	6
1.1.3. Các bài toán thông dụng trong khai phá dữ liệu.....	7
1.1.4. Ứng dụng của khai phá dữ liệu.....	7
1.2. Khai phá Web	8
1.2.1. Giới thiệu về khai phá Web.....	8
1.2.2. Khó khăn và thuận lợi	9
1.2.3. Quá trình khai phá Web.....	12
1.2.4. Các lĩnh vực của khai phá dữ liệu web.....	15
1.2.5. Các kiểu dữ liệu Web	16
1.3. Phân cụm tài liệu web.....	17
1.4. Phân lớp văn bản	19
1.4.1. Bài toán phân lớp văn bản	19
1.4.2. Dữ liệu văn bản.....	21
1.4.3. Biểu diễn văn bản	21
1.4.4. Một số vấn đề trong xử lý dữ liệu văn bản.....	23
1.5. Tổng kết chương 1	29
CHƯƠNG 2: MÔ HÌNH HỆ THỐNG TỔNG HỢP, PHÂN LOẠI THÔNG TIN TỰ ĐỘNG.....	30
2.1. Các phương pháp tách từ tiếng Việt.....	30
2.1.1. Phương pháp Maximum Matching: forward/backward	30

2.1.2. Phương pháp giải thuật học cải biến (Transformation-based Learning)	31
2.1.3. Mô hình tách từ bằng WFST và mạng Neural	32
2.1.4. Phương pháp quy hoạch động (Dynamic Programming).....	34
2.1.5. Phương pháp tách từ tiếng việt dựa trên thống kê từ Internet và thuật toán di truyền IGATEC	35
2.2. Các phương pháp phân loại văn bản	37
2.2.1. Phương pháp phân lớp Bayes (Naïve Bayes).....	37
2.2.2. Phương pháp k-người láng giềng gần nhất (K-Nearest Neighbor)	39
2.2.3. Phương pháp máy hỗ trợ vector (Support vector Machine).....	40
2.2.4. Phương pháp mạng nơron (Neural Network).....	42
2.2.5. Phương pháp Linear Least Square Fit.....	43
2.2.6. Phương pháp Centroid-based vector	44
2.3. Phân tích và xác định yêu cầu	46
2.3.1. Đặt vấn đề.....	46
2.3.2. Xác định yêu cầu của hệ thống.....	46
2.4. Mô hình hệ thống.....	47
2.4.1 Kiến trúc chung	47
2.4.2. Thành phần Web Crawler.....	48
2.4.3. Thành phần Extractor	49
2.4.4. Xử lý tài liệu	50
2.4.5. Phân loại văn bản tiếng Việt.....	52
2.5. Tổng kết chương 2.....	56
CHƯƠNG 3: XÂY DỰNG HỆ THỐNG TỔNG HỢP, PHÂN LOẠI THÔNG TIN	
VIỆC LÀM TỰ ĐỘNG.....	57
3.1. Mô tả chức năng hệ thống	57
3.1.1. Chức năng thu thập và xử lý tin tức	57
3.1.2. Chức năng người dùng	57
3.1.3. Chức năng quản trị	57
3.2. Giải pháp và công nghệ sử dụng	58

3.2.1. Công cụ rút trích dữ liệu HtmlAgility Pack	58
3.2.2. Ngôn ngữ truy vấn Xpath	60
3.3. Thiết kế cơ sở dữ liệu	64
3.4. Phát triển chương trình	65
3.4.1. Xây dựng phân hệ Crawler	65
3.4.2. Xây dựng phân hệ Extractor	66
3.4.3. Xây dựng phân hệ xử lý dữ liệu	69
3.4.4. Xây dựng cổng thông tin tổng hợp	69
3.5. Kết quả thử nghiệm hệ thống	69
3.6. Tổng kết chương 3	73
KẾT LUẬN	74
TÀI LIỆU THAM KHẢO	74

DANH MỤC CÁC KÝ HIỆU, CÁC CHỮ VIẾT TẮT

KDD	Knowledge Discovery in Database
KPDL	Khai phá dữ liệu
IGATEC	Internet and Genetics Algorithm-based Text Categorization for Documents in Vietnamese
kNN	K-Nearest Neighbor
LLSF	Linear Least Square Fit
NB	Naïve Bayes
NNet	Neural Network
LLSF	Linear Lest Square Fit
DF	Tần suất tài liệu (Document Frequency)
TBL	Phương pháp giải thuật học cải biên (Transformation – based Learning)
IDF	Tần suất tài liệu ngược (Inverse document frequency)
TF	Tần suất từ (Term frequency)

DANH MỤC CÁC BẢNG

Bảng 1.1: Thống kê các từ tần số xuất hiện cao (thống kê của B. Croft, UMass)	24
Bảng 3.1. Một số cú pháp của XPath	62
Bảng 3.2. Bảng tin tức	64
Bảng 3.3. Bảng chuyên mục tin.....	65
Bảng 3.4. Kênh tin.....	65
Bảng 3.5. Cấu hình và yêu cầu của máy thử nghiệm	69

DANH MỤC CÁC HÌNH

Hình 1.1. Các bước trong khám phá tri thức	5
Hình 1.2. Quá trình khai phá dữ liệu	6
Hình 1.3. Quá trình khai phá văn bản Web	12
Hình 1.4. Nội dung khai phá dữ liệu Web	16
Hình 1.5. Phân loại dữ liệu Web	17
Hình 1.6. Phân lớp văn bản	20
Hình 1.7. Biểu diễn văn bản	22
Hình 1.8. Lược đồ thống kê tần số của từ theo Định luật Zipf.....	25
Hình 2.1. Sơ đồ hệ thống WFST	32
Hình 2.2. Hệ thống IGATEC.....	35
Hình 2.3. Siêu mặt phẳng h phân chia dữ liệu huấn luyện thành 2 lớp + và – với khoảng cách biên lớn nhất.	41
Hình 2.4. Kiến trúc mô đun (Modular Architecture)	43
Hình 2.5. Mô hình kiến trúc hệ thống thu thập tin	48
Hình 3.1. Giải thuật hoạt động phân hệ Crawler.....	66
Hình 3.2. Ví dụ sơ đồ cây DOM.....	67
Hình 3.2. Giải thuật hoạt động của phân hệ Extractor	69
Hình 3.3. Giao diện trang chủ	70
Hình 3.4. Quản lý kênh tin	71
Hình 3.5. Quản lý cập nhập tin.....	71
Hình 3.6. Quản lý chuyên mục tin.....	72
Hình 3.7. Quản lý tin tức	72

MỞ ĐẦU

1. Lý do chọn đề tài

Trong những năm gần đây cùng với sự phát triển nhanh chóng của khoa học kỹ thuật là sự bùng nổ về tri thức. Kho dữ liệu, nguồn tri thức của nhân loại cũng trở nên đồ sộ, vô tận làm cho vấn đề khai thác các nguồn tri thức đó ngày càng trở nên nóng bỏng và đặt ra thách thức lớn cho nền công nghệ thông tin thế giới.

Cùng với những tiến bộ vượt bậc của công nghệ thông tin là sự phát triển mạnh mẽ của mạng thông tin toàn cầu, nguồn dữ liệu Web trở thành kho dữ liệu khổng lồ. Nhu cầu khai thác và xử lý thông tin phục vụ cho công tác quản lý, hoạt động sản xuất, kinh doanh, học tập... đã trở nên cấp thiết trong xã hội hiện đại. Do đó số lượng văn bản xuất hiện trên mạng Internet cũng tăng theo một tốc độ chóng mặt. Với lượng thông tin đồ sộ như vậy, một yêu cầu lớn đặt ra là làm sao tổ chức, tìm kiếm và có được thông tin nhanh chóng, hiệu quả nhất.

Để giải quyết vấn đề này, có một hướng giải quyết là nghiên cứu và áp dụng kỹ thuật khai phá dữ liệu trong môi trường Web. Vì vậy tôi chọn đề tài “nghiên cứu xây dựng hệ thống tổng hợp, phân loại thông tin tự động trên web” nhằm tìm hiểu phương pháp tổng hợp tin từ nhiều website và tự động phân loại các tin được lấy về.

2. Đối tượng và phạm vi nghiên cứu

Đối tượng nghiên cứu:

Tìm hiểu về khai phá dữ liệu web, các thuật toán phân loại tài liệu và ứng dụng trong truy xuất thông tin tự động. Trên cơ sở đó, xây dựng hệ thống tổng hợp, phân loại thông tin tự động trên web.

Phạm vi nghiên cứu:

- ✓ Khai phá dữ liệu web.
- ✓ Các giải thuật phân cụm tài liệu.

- ✓ Các kỹ thuật và công nghệ hỗ trợ trích xuất thông tin tự động.
- ✓ Kết hợp các yếu tố trên để xây dựng hệ thống tổng hợp, phân loại thông tin trực tuyến.

3. Hướng nghiên cứu của đề tài

Về lý thuyết: Nghiên cứu các giải pháp kỹ thuật trong việc thu thập thông tin tự động trên internet, ứng dụng kỹ thuật khai phá dữ liệu cho việc phân tích thông tin thu thập được theo các lĩnh vực khác nhau nhằm giúp người dùng theo dõi, tìm kiếm thông tin dễ dàng, thuận tiện.

Về thực tiễn: Ứng dụng hệ thống này trong việc xây dựng hệ thống tổng hợp, phân loại thông tin việc làm tự động.

4. Những nội dung chính

Luận văn được trình bày trong 3 chương, có phần mở đầu, phần kết luận, phần mục lục, phần tài liệu tham khảo. Các nội dung cơ bản của luận văn được trình bày như sau:

Chương 1: Trình bày những nội dung tổng quan về khai phá dữ liệu, khai phá web, phân loại văn bản.

Chương 2: Trình bày một số phương pháp tách, phân loại từ tiếng Việt và mô hình hệ thống tổng hợp, phân loại tin tức.

Chương 3: Trình bày giải pháp xây dựng thử nghiệm hệ thống tổng hợp, phân loại thông tin việc làm tự động.

5. Phương pháp nghiên cứu

Nghiên cứu lý thuyết:

- Tìm hiểu lý thuyết về khai phá dữ liệu và khai phá dữ liệu web.
- Tìm hiểu các thuật toán phân cụm tài liệu.
- Tìm hiểu cơ chế hoạt động của các hệ thống tìm kiếm thu thập thông tin.

Nghiên cứu thực nghiệm:

- Dựa trên lý thuyết đã nghiên cứu, tiến hành xây dựng hệ thống thu thập và phân loại thông tin từ các kênh tin được cấu hình trước.
- Thử nghiệm trên máy đơn qua localhost có kết nối internet.

6. Ý nghĩa khoa học

Về mặt lý thuyết: Giới thiệu tổng quan, ứng dụng của khai phá dữ liệu web, các thuật toán phân loại tài liệu và cơ chế của hệ thống thu thập tin.

Về mặt thực tiễn: Xây dựng hệ thống tổng hợp, phân loại thông tin tự động trên web. Cho phép người dùng cập nhật các thông tin mới nhất từ các website khác, lưu trữ, tìm kiếm thông tin theo các chuyên mục.