

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

**HOÀNG MINH THỦY**

**NGHIÊN CỨU CÁC PHƯƠNG PHÁP TRÍCH CHỌN THÔNG TIN  
VÀ ỨNG DỤNG TRÍCH CHỌN THÔNG TIN DU LỊCH  
TRONG VĂN BẢN TIẾNG VIỆT**

**Chuyên ngành: KHOA HỌC MÁY TÍNH**

**Mã số: 60 48 01 01**

**LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH**

**NGƯỜI HƯỚNG DẪN KHOA HỌC**

**GS. VŨ ĐỨC THI**

**Thái Nguyên – 2015**

## LỜI CAM ĐOAN

Tác giả Hoàng Minh Thủy xin cam kết rằng nội dung của Luận văn này chưa được nộp cho bất kỳ một chương trình cấp bằng cao học nào cũng như bất kỳ một chương trình đào tạo cấp bằng nào khác.

Ngoài ra, tác giả cũng xin cam kết Luận văn thạc sĩ này là nỗ lực riêng của cá nhân tác giả. Các kết quả, phân tích, kết luận trong Luận văn thạc sĩ này (ngoài các phần được trích dẫn) đều là kết quả làm việc của cá nhân tác giả.

*Thái Nguyên, ngày 10 tháng 11 năm 2015*

**Tác Giả**

**Hoàng Minh Thủy**

## LỜI CẢM ƠN

Lời đầu tiên em xin gửi lời cảm ơn chân thành đến Các quý thầy cô giáo, Tổ chuyên môn Trường Đại học Công nghệ thông tin và Truyền thông - Đại học Thái Nguyên đã tận tình giảng dạy, truyền đạt những kiến thức, kinh nghiệm quý báu trong suốt thời gian em theo học tại trường. Các kiến thức, kinh nghiệm quý báu của các Quý thầy cô giáo không chỉ giúp cá nhân em hoàn thiện hệ thống kiến thức trong học tập mà còn giúp em ứng dụng các kiến thức đó trong công tác hiện tại tại đơn vị.

Đặc biệt, em xin chân thành cảm ơn thầy giáo GS. Vũ Đức Thi đã rất nhiệt tình và tâm huyết trong việc định hướng và giúp đỡ em hoàn thành luận văn này.

Ngoài ra, em cũng xin chân thành cảm ơn Ban lãnh đạo và cán bộ viên chức Trường Đại học Lâm nghiệp đã tạo điều kiện cung cấp những ý kiến quý báu và những kiến thức thực tiễn cho em thực hiện luận văn tốt nghiệp này.

Em cũng xin được bày tỏ tình cảm với gia đình, đồng nghiệp, bạn bè đã tạo điều kiện để cá nhân em có thể dành thời gian cho khóa học. Xin chân thành cảm ơn những người bạn lớp cao học CK13, trong 2 năm qua đã luôn luôn động viên, khích lệ và hỗ trợ em trong quá trình học tập.

Trong quá trình thực hiện Luận văn mặc dù đã cố gắng hết mình, song chắc chắn luận văn của em vẫn còn nhiều thiếu sót. Em rất mong nhận được sự chỉ bảo và đóng góp tận tình của các thầy cô để luận văn của em được hoàn thiện hơn.

*Thái Nguyên, ngày 10 tháng 11 năm 2015*

**Tác Giả**

**Hoàng Minh Thủy**

## MỤC LỤC

LỜI CAM ĐOAN .....	i
LỜI CẢM ƠN.....	iii
MỤC LỤC.....	iv
DANH MỤC CÁC BẢNG .....	vii
DANH MỤC CÁC HÌNH.....	viii
MỞ ĐẦU .....	1
1.1.Sự cần thiết lựa chọn đề tài.....	1
1.2.Mục tiêu đề tài.....	2
1.3.Đối tượng và phạm vi nghiên cứu.....	2
1.4.Phương pháp nghiên cứu.....	2
1.5.Cấu trúc của luận văn.....	2
Chương 1 .....	4
TỔNG QUAN VỀ TRÍCH CHỌN THÔNG TIN VÀ BÀI TOÁN TRÍCH CHỌN THÔNG TIN DU LỊCH.....	4
1.1.Tổng quan về trích chọn thông tin .....	4
1.1.1. Bài toán trích chọn thực thể .....	5
1.1.2. Bài toán trích chọn quan hệ.....	7
1.1.3. Bài toán trích chọn cụm từ khóa.....	8
1.2.Bài toán trích chọn thông tin du lịch.....	9
1.3.Ý nghĩa của bài toán trích chọn thông tin du lịch.....	10
1.3.1. Ý nghĩa khoa học.....	10
1.3.2. Ý nghĩa thực tế.....	10
1.4.Ứng dụng của bài toán trích chọn thông tin du lịch.....	10
1.4.1. Hệ thống tìm kiếm và tư vấn du lịch .....	10
1.4.2. Bài toán dự đoán xu hướng du lịch .....	11
1.5.Kết luận chương .....	11
Chương 2 .....	12
MỘT SỐ PHƯƠNG PHÁP TRÍCH CHỌN THÔNG TIN .....	12

2.1.Trích chọn thông tin dựa vào cây DOM .....	12
2.1.1. <i>Khái niệm cây DOM</i> .....	12
2.1.2. <i>Xây dựng cây DOM</i> .....	13
2.1.3. <i>Sử dụng cây DOM để trích chọn thông tin</i> .....	14
2.2.Trích chọn thông tin dựa trên tập luật.....	15
2.2.1. <i>Hình thức và biểu diễn của luật</i> .....	16
2.2.2. <i>Đặc trưng của từ tố (token)</i> .....	16
2.2.3. <i>Tập luật xác định thực thể đơn</i> .....	16
2.2.4. <i>Các luật đánh dấu biên của thực thể</i> .....	18
2.2.5. <i>Các luật xác định nhiều thực thể</i> .....	18
2.2.6. <i>Đánh giá phương pháp tiếp cận dựa trên luật</i> .....	19
2.3.Trích chọn thông tin dựa trên học máy.....	19
2.4.Phương pháp kết hợp giữa phân tích mã HTML và luật .....	20
2.5.Kết luận chương.....	21
Chương 3 .....	22
<b>BÀI TOÁN TRÍCH CHỌN TOUR DU LỊCH TRÊN MỘT SỐ TRANG</b>	
<b>THÔNG TIN ĐIỆN TỬ TIẾNG VIỆT.....</b>	<b>22</b>
3.1.Bài toán trích chọn thông tin du lịch trên một số trang thông tin điện tử	
tiếng Việt.....	22
3.1.1. <i>Phát biểu bài toán</i> .....	22
3.1.2. <i>Ý tưởng giải quyết</i> .....	23
3.2.Phương pháp giải quyết bài toán.....	23
3.2.1. <i>Bộ thu thập dữ liệu</i> .....	25
3.2.2. <i>Bộ lọc dữ liệu</i> .....	26
3.2.3. <i>Bộ trích chọn tour</i> .....	27
3.2.4. <i>Bộ trích chọn thuộc tính</i> .....	29
Chương 4 .....	38
<b>THỬ NGHIỆM VÀ ĐÁNH GIÁ KẾT QUẢ .....</b>	<b>38</b>
4.1.Bài toán thử nghiệm.....	38

4.2.Môi trường và các công cụ thử nghiệm .....	38
4.2.1. <i>Môi trường thử nghiệm</i> .....	38
4.2.2. <i>Công cụ phần mềm sử dụng để thử nghiệm</i> .....	39
4.3.Xây dựng cơ sở dữ liệu .....	39
4.4.Thử nghiệm quy trình trích chọn tour du lịch.....	41
4.4.1. <i>Thu thập dữ liệu (Web Crawler)</i> .....	41
4.4.2. <i>Lọc dữ liệu</i> .....	44
4.4.3. <i>Trích chọn các tour du lịch và các thuộc tính</i> .....	46
4.5.Phân tích lỗi.....	49
4.5.1. <i>Phân tích lỗi của bộ lọc dữ liệu</i> .....	49
4.5.2. <i>Phân tích lỗi của quá trình trích chọn</i> .....	51
4.6.Một số ứng dụng kết quả trích chọn tour du lịch.....	51
4.6.1. <i>Thống kê theo định danh</i> .....	52
4.6.2. <i>Thống kê theo giá tour</i> .....	54
4.6.3. <i>Thống kê theo thời gian</i> .....	55
4.7.Kết luận chương .....	57
KẾT LUẬN.....	58
TÀI LIỆU THAM KHẢO .....	59

**DANH MỤC CÁC BẢNG**

Bảng 1.1. Bảng phân loại thực thể.....	6
Bảng 4.1. Cấu hình hệ thống thử nghiệm .....	38
Bảng 4.2. Công cụ phần mềm có sẵn.....	39
Bảng 4.3. Kết quả lọc các bài viết chứa thông tin về các tour du lịch.....	45
Bảng 4.4. Kết quả trích chọn tour du lịch và trích chọn thuộc tính.....	47
Bảng 4.5. Bảng thống kê số tour theo địa danh du lịch .....	52
Bảng 4.6. Bảng thống kê số tour theo giá .....	54
Bảng 4.7. Bảng thống kê số tour theo thời gian du lịch.....	56

## DANH MỤC CÁC HÌNH

Hình 2.1. Mô hình biểu diễn cây DOM .....	12
Hình 2.2. Minh họa sử dụng visual cue .....	14
Hình 2.3. Minh họa cây DOM dùng trong mẫu trích chọn.....	15
Hình 3.1. Mô hình bài toán trích chọn .....	25
Hình 3.2. Mô hình làm việc của bộ thu thập dữ liệu .....	25
Hình 3.3. Mô hình làm việc của bộ lọc dữ liệu.....	26
Hình 3.4. Các thông tin chi tiết về tour của website Du lịch Dấu Chân.....	30
Hình 3.5. Các thông tin chi tiết về tour của website Du lịch Năm Châu.....	30
Hình 3.6. Các thông tin chi tiết về tour của website Du lịch Quốc tế Nét Việt.....	31
Hình 3.7. Các thông tin chi tiết về tour của website Du lịch AMI TOUR .....	31
Hình 3.8. Các thông tin chi tiết về tour của website Du lịch Giấc Mơ Việt...	32
Hình 3.9. Các thông tin chi tiết về tour của website Du lịch Việt .....	33
Hình 3.10. Các thông tin chi tiết về tour của website Du lịch Á Châu.....	34
Hình 3.11. Mô hình làm việc của bộ trích chọn thuộc tính .....	35
Hình 4.1. Thu thập dữ liệu từ trang <a href="http://www.dulichnamchau.vn">www.dulichnamchau.vn</a> . .....	43
Hình 4.2. Quá trình thu thập dữ liệu từ trang <a href="http://www.dulichnamchau.vn">www.dulichnamchau.vn</a> . .....	44
Hình 4.3. Kết quả lọc các bài viết chứa thông tin về các tour du lịch .....	46
Hình 4.4. Kết quả trích chọn các tour du lịch .....	48
Hình 4.5. Giao diện tra cứu tour du lịch .....	49
Hình 4.6. Lỗi lọc dữ liệu khi thông tin ở dạng lựa chọn.....	50
Hình 4.7. Lỗi lọc dữ liệu khi không có thông tin về tour du lịch .....	50
Hình 4.8. Biểu đồ thống kê số tour theo địa danh du lịch .....	53
Hình 4.9. Biểu đồ thống kê số tour theo giá tiền .....	55
Hình 4.10. Biểu đồ thống kê số tour theo thời gian.....	56



## MỞ ĐẦU

### 1.1. Sự cần thiết lựa chọn đề tài

Trích chọn thông tin (IE - Information Extraction) là một lĩnh vực nghiên cứu quan trọng trong khai phá dữ liệu văn bản [3, 4]. Trích chọn thông tin là quá trình thu thập thông tin từ các nguồn dữ liệu theo nhiều định dạng khác nhau, không đồng nhất, thậm chí không có định dạng cụ thể, sau đó chuyển thành một dạng đồng nhất. Dữ liệu sau khi trích chọn được lưu vào cơ sở dữ liệu để xử lý hay được sử dụng cho những hệ thống khai phá dữ liệu. Từ dữ liệu, thông tin được trích chọn ra có thể sử dụng các kỹ thuật phân tích, khai phá để khám phá ra các mẫu thông tin có ích, tiềm ẩn trong dữ liệu.

Ngày nay, cùng với sự phát triển của công nghệ thông tin, Tin học đã dần được ứng dụng rộng rãi trong nhiều lĩnh vực như kinh tế, du lịch, thương mại, y tế, ngân hàng và mang lại nhiều lợi ích to lớn. Nền kinh tế không ngừng phát triển, đời sống văn hoá - xã hội ngày càng được nâng cao thì du lịch đã trở thành một nhu cầu không thể thiếu trong cuộc sống của người dân, trên các trang web du lịch là hàng loạt thông tin về các tour du lịch trong nước và ngoài nước. Tuy nhiên lượng thông tin về các tour du lịch trên Internet là vô cùng lớn, gây khó khăn cho người có nhu cầu du lịch trong việc lựa chọn địa điểm tham quan, lựa chọn công ty cung cấp dịch vụ,... Do vậy, một bài toán đặt ra là cần phải xây dựng một hệ thống tìm kiếm và tư vấn du lịch, giúp người dùng có thể lựa chọn được những tour du lịch phù hợp nhất với yêu cầu đề ra. Để có một hệ thống tìm kiếm và tư vấn tốt thì trước tiên ta phải xây dựng được tập dữ liệu có độ chính xác cao. Cùng với nó là bài toán con trích chọn thông tin du lịch trong văn bản tiếng Việt.

Để có thể tiến đến tìm hiểu được những vấn đề trên, em lựa chọn đề tài ***“Nghiên cứu các phương pháp trích chọn thông tin và ứng dụng trích chọn thông tin du lịch trong văn bản Tiếng Việt”*** làm luận văn tốt nghiệp Thạc sĩ của mình.

## **1.2. Mục tiêu đề tài**

Tìm hiểu các phương pháp trích chọn thông tin và xây dựng mô hình giải quyết bài toán trích chọn thông tin về các tour du lịch từ các trang thông tin điện tử tiếng Việt trên Internet.

## **1.3. Đối tượng và phạm vi nghiên cứu**

Đối tượng nghiên cứu của đề tài là các phương pháp tiếp cận giải quyết bài toán trích chọn thông tin trong văn bản tiếng Việt và các trang thông tin điện tử tiếng Việt trên mạng Internet về lĩnh vực du lịch.

Phạm vi nghiên cứu của đề tài là bài toán trích chọn thông tin về các tour du lịch trên một số trang thông tin điện tử tiếng Việt (website) trên mạng Internet.

## **1.4. Phương pháp nghiên cứu**

Phương pháp nghiên cứu của đề tài là nghiên cứu lý thuyết và nghiên cứu thực nghiệm.

Về nghiên cứu lý thuyết, đề tài đã tổng hợp các kết quả nghiên cứu về các phương pháp trích chọn thông tin từ văn bản tiếng Việt phục vụ phân tích, thống kê, báo cáo, ra quyết định. Về nghiên cứu thực nghiệm, đề tài xây dựng và cài đặt, thử nghiệm mô hình trích chọn thông tin du lịch từ một số trang web về du lịch bằng tiếng Việt trên mạng Internet.

## **1.5. Cấu trúc của luận văn**

Cấu trúc luận văn gồm: mở đầu, bốn chương chính, kết luận và tài liệu tham khảo.

*Phần mở đầu:* Lý do chọn đề tài và bố cục luận văn

*Chương 1:* Giới thiệu tổng quan bài toán trích chọn thông tin và một số lĩnh vực nghiên cứu liên quan.