

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

PHẠM VĂN DƯƠNG

**NGHIÊN CỨU MỘT SỐ THUẬT TOÁN LỌC
THƯ RÁC VÀ ỨNG DỤNG TRONG LỌC EMAIL NỘI BỘ
CỦA VIỆN THÔNG TIN BẮC KẠN**

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

Thái Nguyên - 2017

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

PHẠM VĂN DƯƠNG

**NGHIÊN CỨU MỘT SỐ THUẬT TOÁN LỌC
THƯ RÁC VÀ ỨNG DỤNG TRONG LỌC EMAIL
NỘI BỘ CỦA VIỆN THÔNG TIN BẮC KẠN**

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: TS. NGUYỄN HẢI MINH

Thái Nguyên - 2017

LỜI CAM ĐOAN

Tôi xin cam đoan luận văn là công trình nghiên cứu của riêng cá nhân tôi, không sao chép của ai do tôi tự nghiên cứu, đọc, dịch tài liệu, tổng hợp và thực hiện. Nội dung lý thuyết trong luận văn tôi có sử dụng một số tài liệu tham khảo như đã trình bày trong phần tài liệu tham khảo. Các số liệu, chương trình phần mềm và những kết quả trong luận văn là trung thực và chưa được công bố trong bất kỳ một công trình nào khác.

Thái Nguyên, tháng 4 năm 2017

Học viên thực hiện

Phạm Văn Dương

LỜI CẢM ƠN

Lời đầu tiên, em xin gửi lời biết ơn sâu sắc đến TS. Nguyễn Hải Minh người đã tận tình hướng dẫn, chỉ bảo, giúp đỡ em trong suốt quá trình làm luận văn.

Em cũng xin gửi lời cảm ơn đến các thầy cô giáo trường Đại học Công nghệ thông tin và Truyền thông, các thầy cô Viện Công nghệ thông tin đã truyền đạt những kiến thức và giúp đỡ em trong suốt quá trình học của mình.

Và cuối cùng tôi xin gửi lời cảm ơn tới các đồng nghiệp, gia đình và bạn bè những người đã ủng hộ, động viên tạo mọi điều kiện giúp đỡ để tôi có được kết quả như ngày hôm nay.

Thái Nguyên, tháng 4 năm 2017

Học viên

Phạm Văn Dương

MỤC LỤC

LỜI CAM ĐOAN	1
LỜI CẢM ƠN	2
MỤC LỤC.....	3
CÁC HÌNH VẼ, BẢNG BIỂU TRONG LUẬN VĂN.....	5
MỞ ĐẦU.....	6
Chương 1. THƯ RÁC VÀ CÁC PHƯƠNG PHÁP LỌC THƯ RÁC	8
1.1. Một số khái niệm cơ bản	8
1.1.1. Định nghĩa thư rác.....	8
1.1.2. Phân loại thư rác.....	9
1.2. Các phương pháp lọc thư rác	10
1.2.1. Lọc thư rác thông qua việc đưa ra luật lệ nhằm hạn chế, ngăn chặn việc gửi thư rác.	10
1.2.2. Lọc thư rác dựa trên địa chỉ IP	11
1.2.3. Lọc dựa trên chuỗi hỏi/ đáp.....	13
1.2.4. Phương pháp lọc dựa trên mạng xã hội.....	13
1.2.5. Phương pháp lọc nội dung.....	14
Chương 2. TỔNG QUAN CÁC THUẬT TOÁN NSA, PSA, PNSA TRONG LỌC THƯ RÁC	19
2.1. Cơ sở lý thuyết về hệ miễn dịch nhân tạo.....	19
2.1.1. Khái niệm về hệ miễn dịch nhân tạo.....	19
2.1.2. Phạm vi ứng dụng của hệ miễn dịch nhân tạo	19
2.1.3. Cấu trúc cơ bản của hệ miễn dịch nhân tạo	20
2.2. Cơ sở lý thuyết về thuật toán chọn lọc tiêu cực (Negative Selection Algorithms - NSA).....	24
2.3. Cơ sở lý thuyết về thuật toán chọn lọc tích cực (Positive Selection Algorithms – PSA).....	26

2.4. Cơ sở lý thuyết thuật toán cải tiến chọn lọc thư rác (Positive and Negative Selection Algorithms – PNSA)	27
2.4.1. Một số định nghĩa.....	27
2.4.2. Thuật toán sinh tập bộ dò r-chunk.....	30
2.4.3. Thuật toán sinh tập bộ dò dạng r – contiguous	33
2.5. Các nghiên cứu gần đây	36
Chương 3. KẾT QUẢ CÀI ĐẶT CÁC THUẬT TOÁN.	38
3.1. Tổng quan ứng dụng CNTT tại tỉnh Bắc Kạn.....	38
3.2. Mô hình tổng quát	39
3.3. Mô hình thực tế ứng dụng lọc email Spam tại hệ thống email nội bộ của Viễn thông tỉnh Bắc Kạn	40
3.4. Ứng dụng hệ miễn dịch nhân tạo trong lọc thư rác.....	40
3.4.1. Phát biểu bài toán	41
3.4.2. Cơ sở dữ liệu TREC'07	42
3.4.3. Phương pháp.....	42
3.4.4. Phân tích thuật toán.....	43
3.4.5. Đánh giá	45
3.5. So sánh với các thuật toán trên WEKA	46
3.5.1. Phát biểu bài toán	46
3.5.2. Cơ sở dữ liệu SpamBase	46
3.5.3. Phần mềm WEKA.....	49
3.5.4. Thiết kế phần mềm.....	52
3.5.5 Phân tích thuật toán kết hợp chọn lọc tích cực và chọn lọc tiêu cực PNSA	53
3.5.6 Giao diện chương trình và kết quả	56
3.5.7. Đánh giá	59
KẾT LUẬN	61
TÀI LIỆU THAM KHẢO.....	63

CÁC HÌNH VẼ, BẢNG BIỂU TRONG LUẬN VĂN

Hình 1.1: Tất cả các thư điện tử	9
Hình 1.2 : Mô tả tổng quan quá trình hoạt động của honeyd	15
Hình 2.1: Cấu trúc phân tầng của HMD nhân tạo.....	20
Hình 2.2: Kháng thể nhận diện kháng nguyên dựa vào phân bù	22
Hình 2.3 Sơ đồ khối thuật toán chọn lọc tiêu cực	25
Hình 2.4 Sơ đồ khối thuật toán chọn lọc tích cực	27
Hình 3.1. Mô hình tổng quát của quá trình gửi và nhận thư điện tử.....	39
Hình 3.2. Mô hình mạng nội bộ của Viễn Thông Tỉnh Bắc Kạn.....	40
Hình 3.3 Giao diện phần mềm Weka	50
Hình 3.4 Giao diện Weka Explorer.....	51
Hình 3.5 Giao diện Weka Explorer sau khi chọn CSDL Spambase.....	51
Hình 3.6 Phân loại dữ liệu.....	52
Hình 3.7 Giao diện chương trình	56
Bảng 2.1. Kết quả bảng băm A.....	33
Bảng 2.2. Các thông số bảng băm A.....	34
Bảng 3.1. Kết quả khi chạy chương trình với 9 bộ test	45
Bảng 3.2. So sánh kết quả	45
Bảng 3.3. Kết quả thử nghiệm trên WEKA và PNSA	57
Bảng 3.4. So sánh PNSA với một số phương pháp cho kết quả tốt hơn ...	58
Bảng 3.5. So sánh PNSA với một số phương pháp cho kết quả thấp hơn.	58
Bảng 3.6. Kết quả so khớp với giá trị tham số r thay đổi	59

MỞ ĐẦU

Mạng Internet ra đời đã mang lại cho con người những tiện ích hết sức to lớn và quan trọng, một trong những tiện ích đó là dịch vụ thư điện tử. Vì, đó là phương tiện giao tiếp đơn giản, tiện lợi, rẻ và hiệu quả giúp mọi người gắn kết và liên lạc với nhau thường xuyên hơn. Tuy nhiên, lợi dụng tính mở của công nghệ và cơ chế trao đổi thư mà hàng ngày người dùng nhận được một số thư ngoài mong đợi đó là thư rác (Spam). Thư rác thường được gửi với số lượng rất lớn thường vì mục đích quảng cáo, thậm trí là đính kèm mã độc dưới dạng Virus gây phiền toái cho người dùng, làm giảm tốc độ xử lý của máy chủ mail server.

Thư rác (spam) là thư điện tử được gửi hàng loạt với nội dung mà người nhận không mong đợi, không muốn xem, hay chứa những nội dung không liên quan đến người nhận và thường được sử dụng để gửi thông tin quảng cáo. Do có giá thành tương đối thấp so với các phương pháp quảng cáo khác, thư rác hiện chiếm một tỷ lệ lớn và ngày càng tăng trong tổng số thư điện tử được gửi qua Internet. Sự xuất hiện và gia tăng thư rác không những gây khó chịu và làm mất thời gian của người nhận mà còn ảnh hưởng tới đường truyền Internet và làm chậm tốc độ xử lý của máy chủ thư điện tử, gây thiệt hại lớn về kinh tế.

Xuất phát từ lý do đó, đề tài đặt vấn đề nghiên cứu một số thuật toán LỌC THƯ RÁC, một trong những thuật toán mới được công bố gần đây để đề xuất một mô hình thực nghiệm trên một dịch vụ email thực tế. Qua đó hướng tới xây dựng ứng dụng bằng cách tích hợp thêm một số Module trong hỗ trợ sử dụng dịch vụ sử dụng email.

Nội dung luận văn gồm có 3 chương:

Dự kiến nội dung báo cáo của luận văn gồm: Phần mở đầu, 3 chương chính, phần kết luận, tài liệu tham khảo, phụ lục. Bộ cục được trình bày như sau:

Phần mở đầu: Nêu lý do chọn đề tài và hướng nghiên cứu chính

Chương 1: THƯ RÁC VÀ CÁC PHƯƠNG PHÁP LỌC THƯ RÁC

Chương 2: TỔNG QUAN CÁC THUẬT TOÁN NSA, PSA, PNSA TRONG LỘC THƯ RÁC.

Chương 3: KẾT QUẢ CÀI ĐẶT CÁC THUẬT TOÁN.

Phần kết luận: Tóm tắt các kết quả đã đạt được và hướng phát triển tiếp theo của đề tài.

Chương 1.

THƯ RÁC VÀ CÁC PHƯƠNG PHÁP LỌC THƯ RÁC

Một trong những dịch vụ Internet mang lại đó là dịch vụ thư điện tử, đây là phương pháp giao tiếp rất đơn giản, tiện lợi, rẻ và hiệu quả giữa mọi người. Tuy nhiên, chính vì những lợi ích của dịch vụ thư điện tử mang lại mà số lượng thư trao đổi trên Internet ngày càng tăng và hầu hết trong số những thư đó là thư rác (Email spam). Thư rác thường được gửi với số lượng lớn, người dùng không mong đợi với nhiều mục đích khác nhau như: Quảng cáo, đính kèm virus, gây phiền toái khó chịu cho người dùng, làm giảm tốc độ Internet và tốc độ xử lý của server, gây thiệt hại lớn về kinh tế. Chương này khái quát các vấn đề về thư rác, ảnh hưởng của thư rác trong cuộc sống và các phương pháp ngăn chặn thư rác. Các khái niệm trong chương này được tham khảo trong [1], [2], [3], [4].

1.1. Một số khái niệm cơ bản

1.1.1. Định nghĩa thư rác

Có nhiều tranh cãi về việc đâu là định nghĩa chính xác của thư rác (spam email), bởi vì thư rác mang tính cá nhân hóa nên khó mà nói lên được hết ý nghĩa của thư rác. Nhiều ý kiến cho rằng thư rác là những “thư điện tử (email) không mong muốn”. Định nghĩa này cũng không thực sự chính xác, như một nhân viên nhận những thư điện tử về công việc từ sếp của họ, đây là những thư điện tử người nhân viên không mong muốn nhưng chúng không phải là thư rác. Lại có ý kiến khác cho rằng thư rác là những “thư điện tử thương mại không được yêu cầu từ phía người nhận” - những thư này bao gồm các thư điện tử quảng cáo về các sản phẩm và thư điện tử lừa gạt. Nhưng định nghĩa này cũng không thực sự chính xác, nó làm mọi người nghĩ rằng thư rác giống như là thư đáng bỏ đi (junk mail). Sau đây sẽ đưa ra một định nghĩa thông dụng nhất về thư rác và giải thích các đặc điểm của nó để phân biệt thư rác với thư thông thường [1,2]: