

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

HÀ MẠNH KIÊN

**KỸ THUẬT PHÂN LỚP DỮ LIỆU VÀ ỨNG DỤNG
TRONG PHÁT HIỆN MÃ ĐỘC**

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.Lương Thế Dũng

THÁI NGUYÊN - 2015

LỜI CAM ĐOAN

Tôi cam đoan đây là công trình của riêng tôi.

Các số liệu, kết quả nêu trong luận văn là trung thực và chưa từng được ai công bố trong bất kỳ công trình nào khác.

Qua đây tôi xin chân thành cảm ơn toàn thể các thầy cô trong khoa đào tạo sau đại học Trường Đại học Công nghệ Thông tin và Truyền thông – Đại học Thái Nguyên, những người đã trực tiếp giảng dạy, truyền đạt cho tôi kiến thức chuyên môn và phương pháp làm việc khoa học.

Đặc biệt, tôi xin chân thành cảm ơn **TS. Lương Thế Dũng** ,đã tận tình hướng dẫn để tôi có thể hoàn thành luận văn này.

Tôi cũng xin gửi lời cảm ơn tới gia đình, bạn bè, đồng nghiệp đã giúp đỡ, động viên và tạo điều kiện cho tôi trong quá trình làm luận văn.

Tác giả luận văn

Hà Mạnh Kiên

MỤC LỤC

ĐẶT VẤN ĐỀ	1
CHƯƠNG 1: TỔNG QUAN VỀ MÃ ĐỘC HẠI	2
1.1. Các loại mã độc.....	2
1.1.1. Virus.....	2
1.1.2. Worm.....	3
1.1.3. Trojan Horse.....	3
1.1.4. Malicious Mobile Code.....	5
1.1.5. Tracking Cookie.....	6
1.1.6. Phần mềm gián điệp (Spyware)	6
1.1.7. Attacker Tool	7
1.1.8. Phishing.....	9
1.2. Phương pháp phát hiện mã độc hại	9
1.2.1. Phần mềm phát hiện mã độc	9
1.2.2. Kỹ thuật phát hiện phần mềm mã độc	10
1.2.3. Kỹ thuật phát hiện dựa mẫu nhận dạng	10
1.2.4. Phát hiện dựa trên đặc điểm	12
1.2.5. Phát hiện dựa trên hành vi.....	12
1.2.6. Kỹ thuật gây nhiễu	13
1.2.7. Phân tích sự tương tự	14
1.2.8. Chuẩn hóa mã độc	15
CHƯƠNG 2: MỘT SỐ KỸ THUẬT PHÂN LỚP	16
2.1. Tổng quan về khai phá dữ liệu	16
2.1.1. Khái niệm về khai phá dữ liệu.....	16
2.1.2. Ứng dụng trong khai phá dữ liệu.....	16
2.1.3. Các bài toán chính trong khai phá dữ liệu.....	17
2.1.4. Tiến trình khai phá dữ liệu.	20

2.2. Một số kỹ thuật phân lớp dữ liệu.....	22
2.2.1. Khái niệm phân lớp.	22
2.2.2. Mục đích của phân lớp.....	24
2.2.3. Các tiêu chí để đánh giá thuật toán phân lớp.	24
2.2.4. Các phương pháp đánh giá độ chính xác của mô hình phân lớp phương pháp holdout.....	25
2.3. Phân lớp dựa trên phương pháp học Naïve bayes.....	26
2.3.1 Giới thiệu	26
2.3.2. Bộ phân lớp Naïve Bayes.	28
2.4. Phân lớp dựa trên cây quyết định (Decision Tree).....	29
2.4.1. Khái niệm cây quyết định:	29
2.4.2. Các vấn đề cần xem xét khi phân lớp dựa cây quyết định.....	42
2.5. Kỹ thuật phân loại máy vector hỗ trợ.....	44
2.5.1. Giới thiệu	44
2.5.2. SVM với tuyến tính.	46
CHƯƠNG 3: ỨNG DỤNG KỸ THUẬT PHÂN LỚP TRONG PHÁT HIỆN MÃ ĐỘC	52
3.1. Mô hình bài toán.....	52
3.1.1. Thu thập dữ liệu	52
3.1.2 Tiền xử lý dữ liệu	53
3.1.3 Lựa chọn thuộc tính	54
3.1.4. Xây dựng bộ phân lớp	58
3.2. Tiến hành thực nghiệm	59
3.2.1. Phân lớp cây quyết định	59
3.2.2. Phân lớp SVM.....	60
3.3 Phân tích và bình luận.....	61
KẾT LUẬN	63
TÀI LIỆU THAM KHẢO	64

DANH MỤC BẢNG

Bảng 3.1.	Bảng kết quả độ chính xác cây quyết định bộ phân lớp đa lớp.....	60
Bảng 3.2	Bảng kết quả độ chính xác cây quyết định bộ phân lớp nhị phân ...	60
Bảng 3.3.	Bảng kết quả xây dựng bộ phân lớp SVM:	61

DANH MỤC HÌNH

Hình 1.1.	Mô tả về Phishing	9
Hình 1.2.	Kiểu phân mềm mã độc cơ bản	10
Hình 1.3.	Mã độc đa hình	11
Hình 1.4.	Phần mềm độc hại siêu đa hình	11
Hình 1.5.	Bộ phát hiện mã độc dựa trên hành vi	13
Hình 1.6.	Kỹ thuật gây nhiễu	14
Hình 2.1.	Quy trình phát hiện tri thức	20
Hình 2.2.	Ước lượng độ chính xác của mô hình phân lớp với phương pháp holdout.	25
Hình 3.1.	Các bước xây dựng mô hình phát hiện mã độc	52
Hình 3.2	Quá trình trích rút các hàm API	56
Hình 3.3	Chi tiết quá trình xây dựng mô hình phát hiện mã độc	58
Hình3.4	Biểu đồ so sánh độ chính xác (%) của hai thuật toán.....	62

ĐẶT VẤN ĐỀ

Khi nhu cầu về việc sử dụng internet của con người ngày càng tăng thì cũng là lúc mối đe dọa xuất hiện ngày càng nhiều, nổi bật là đe dọa của mã độc hại. Mã độc là một loại phần mềm hệ thống do các tin tặc hay các kẻ nghịch ngợm tạo ra nhằm gây hại cho máy tính. Tùy theo cách thức mà tin tặc dung, sự nguy hại của các loại phần mềm khác nhau từ chỗ chỉ hiển thị các cửa sổ thông báo cho đến việc tấn công chiếm máy và lây lan sang máy khác như virus. Xuất hiện bất kỳ đâu trên môi trường của các thiết bị điện tử như các đĩa mềm, usb, đến môi trường Internet trong các website, trong các tin nhắn, trong hòm thư điện tử của người dùng, trong các phần mềm tiện ích.....Khi mã độc hại đã nhiễm vào một máy tính nào đó thì nó sẽ lây lan sang máy tính khác là khá nhanh và khó lường trước được.

Công nghệ thông tin liên tục phát triển và thay đổi, nhiều phần mềm mới ra đời mang đến cho con người nhiều tiện ích hơn. Do vậy để chống lại các loại mã độc hại người ta thường sử dụng các chương trình phát hiện và loại bỏ mã độc hại. Tuy nhiên việc phát hiện mã độc hại của các chương trình hiện nay thường dựa trên các thuật toán đối sánh mẫu và quan trọng là một cơ sở dữ liệu đầy đủ và cập nhật thường xuyên những mẫu mới. Để có một cơ sở dữ liệu như đã nêu cần một chương trình quản lý một cách hiệu quả và tốt rất nhiều công sức để tạo ra các mẫu mã độc hại. Một phương pháp mới hiện nay là dựa trên các mô hình toán học để phát hiện ra các mã độc hại mới mà không sử dụng các cơ sở dữ liệu mẫu, trong đó khai phá dữ liệu là một phương pháp quan trọng và đang được nhiều người quan tâm. Chính vì vậy luận văn này tiến hành nghiên cứu, tìm hiểu các kỹ thuật phân lớp dữ liệu và ứng dụng trong phát hiện mã độc. Nhằm xây dựng ra các mô hình, thuật toán để phát hiện và đánh giá các mô hình đó.

CHƯƠNG 1

TỔNG QUAN VỀ MÃ ĐỘC HẠI

1.1. Các loại mã độc

1.1.1. Virus

Virus là một loại mã độc hại (Malicious code) có khả năng tự nhân bản và lây nhiễm chính nó vào các file, chương trình hoặc máy tính. Như vậy virus máy tính phải luôn luôn bám vào một vật chủ (đó là file dữ liệu hoặc file ứng dụng) để lây lan. Các chương trình diệt virus dựa vào đặc tính này để thực thi việc phòng chống và diệt virus, để quét các file trên thiết bị lưu, quét các file trước khi lưu xuống ổ cứng... Điều này cũng giải thích vì sao đôi khi các phần mềm diệt virus tại PC đưa ra thông báo “phát hiện ra virus nhưng không diệt được” khi thấy có dấu hiệu hoạt động của virus trên PC, bởi vì “vật mang virus” lại nằm ở máy khác nên không thể thực thi việc xóa đoạn mã độc hại đó.

Compiled Virus là virus mà mã thực thi của nó đã được dịch hoàn chỉnh bởi một trình biên dịch để nó có thể thực thi trực tiếp từ hệ điều hành. Các loại boot virus như (Michelangelo và Stoned), file virus (như Jerusalem) rất phổ biến trong những năm 80 là virus thuộc nhóm này, compiled virus cũng có thể là pha trộn bởi cả boot virus và file virus trong cùng một phiên bản.

Interpreted Virus là một tổ hợp của mã nguồn mã chỉ thực thi được dưới sự hỗ trợ của một ứng dụng cụ thể hoặc một dịch vụ cụ thể trong hệ thống. Một cách đơn giản, virus kiểu này chỉ là một tập lệnh, cho đến khi ứng dụng gọi thì nó mới được thực thi. Macro virus, scripting virus là các virus nằm trong dạng này. Macro virus rất phổ biến trong các ứng dụng Microsoft Office khi tận dụng khả năng kiểm soát việc tạo và mở file để thực thi và lây nhiễm. Sự khác nhau giữa macro virus và scripting virus là: Macro virus là

tập lệnh thực thi bởi một ứng dụng cụ thể, còn scripting virus là tập lệnh chạy bằng một service của hệ điều hành. Melissa là một ví dụ xuất sắc về Macro virus, Love Stages là ví dụ cho scripting virus.

1.1.2. Worm

Worm cũng là một chương trình có khả năng tự nhân bản và tự lây nhiễm trong hệ thống tuy nhiên nó có khả năng “tự đóng gói”, điều đó có nghĩa là Worm không cần phải có “file chủ” để mang nó khi nhiễm vào hệ thống. Như vậy, có thể thấy rằng chỉ dùng các chương trình quét file sẽ không diệt được Worm trong hệ thống vì Worm không “bám” vào một file hoặc một vùng nào đó trên đĩa cứng. Mục tiêu của Worm bao gồm cả làm lãng phí nguồn lực băng thông của mạng và phá hoại hệ thống như xoá file, tạo backdoor, thả keylogger,... Tấn công của Worm có đặc trưng là lan rộng cực kỳ nhanh chóng do không cần tác động của con người (như khởi động máy, copy file hay đóng/mở file). Worm có thể chia làm 2 loại:

Network Service Worm lan truyền bằng cách lợi dụng các lỗ hổng bảo mật của mạng, của hệ điều hành hoặc của ứng dụng. Sasser là ví dụ cho loại sâu này.

Mass Mailing Worm là một dạng tấn công qua dịch vụ mail, tuy nhiên nó tự đóng gói để tấn công và lây nhiễm chứ không bám vào vật chủ là email. Khi sâu này lây nhiễm vào hệ thống, nó thường cố gắng tìm kiếm số địa chỉ và tự gửi bản thân nó đến các địa chỉ thu nhận được. Việc gửi đồng thời cho toàn bộ các địa chỉ thường gây quá tải cho mạng hoặc cho máy chủ mail. Netsky, Mydoom là ví dụ cho thể loại này.

1.1.3. Trojan Horse

Trojan Horse là loại mã độc hại được đặt theo sự tích “Ngựa thành Troy”. Trojan horse không có khả năng tự nhân bản tuy nhiên nó lây vào hệ thống với biểu hiện rất bình thường nhưng thực chất bên trong có ẩn chứa các đoạn mã với mục đích gây hại. Trojan có thể gây hại theo ba cách sau:

Tiếp tục thực thi các chức năng của chương trình mà nó bám vào, bên cạnh đó thực thi các hoạt động gây hại một cách riêng biệt (ví dụ như gửi một trò chơi dụ cho người dùng sử dụng, bên cạnh đó là một chương trình đánh cắp password).

Tiếp tục thực thi các chức năng của chương trình mà nó bám vào, nhưng sửa đổi một số chức năng để gây tổn hại (ví dụ như một trojan giả lập một cửa sổ login để lấy password) hoặc che dấu các hành động phá hoại khác (ví dụ như trojan che dấu cho các tiến trình độc hại khác bằng cách tắt các hiển thị của hệ thống).

Thực thi luôn một chương trình gây hại bằng cách núp dưới danh một chương trình không có hại (ví dụ như một trojan được giới thiệu như là một trò chơi hoặc một tool trên mạng, người dùng chỉ cần kích hoạt file này là lập tức dữ liệu trên PC sẽ bị xoá hết).

Có 7 loại trojan chính:

Trojan truy cập từ xa: Được thiết kế để cho kẻ tấn công có khả năng từ xa chiếm quyền điều khiển của máy bị hại. Các trojan này thường dấu vào các trò chơi và các chương trình nhỏ làm cho người dùng mất cảnh giác.

Trojan gửi dữ liệu: Nó thực hiện việc lấy và gửi dữ liệu nhạy cảm như mật khẩu, thông tin thẻ tín dụng, các tệp nhật ký, địa chỉ email... cho kẻ tấn công. Trojan này có thể tìm kiếm cụ thể thông tin hoặc cài phần mềm đọc trộm bàn phím và gửi toàn bộ các phím bấm về cho kẻ tấn công.

Trojan hủy hoại: Thực hiện việc xóa các tệp tin. Loại trojan này giống với virus và thường có thể bị phát hiện bởi các chương trình diệt virus.

Trojan kiểu proxy: Sử dụng máy tính bị hại làm proxy, qua đó có thể sử dụng máy bị hại để thực hiện các hành vi lừa gạt hay đánh phá các máy tính khác.

Trojan FTP: Được thiết kế để mở cổng 21 và cho phép tin tặc kết nối vào máy bị hại sử dụng FTP.