

ĐẠI HỌC BÁCH KHOA HÀ NỘI

ĐỒ ÁN TỐT NGHIỆP

Phân loại lưu lượng mạng mã hóa sử dụng Class
Incremental Learning

LÊ XUÂN NAM

nam.lx162814@sis.hust.edu.vn

Ngành: Kỹ thuật máy tính

Giảng viên hướng dẫn: TS. Tống Văn Vạn

Chữ kí GVHD

Khoa: Kỹ thuật máy tính

Trường: Công nghệ Thông tin và Truyền thông

HÀ NỘI, 08/2024

LỜI CẢM ƠN

Trong suốt những năm tháng học tập tại Đại học Bách khoa Hà Nội, em luôn nhận được sự giúp đỡ và động viên từ gia đình, thầy cô, bạn bè. Là sinh viên của Trường, em đã học được rất nhiều điều từ đạo đức, lối sống, phong cách làm việc đến kiến thức chuyên môn từ những người Thầy, người Cô luôn hết mình giúp đỡ sinh viên, em được tiếp xúc và học hỏi từ rất nhiều sinh viên tài năng của trường. Em xin gửi lời cảm ơn tới bố mẹ, vợ, thầy cô bạn bè đã luôn giúp đỡ em trong quá trình học tập và hoàn thiện đồ án tốt nghiệp. Đặc biệt em xin bày tỏ lòng biết ơn sâu sắc đến thầy Tổng Văn Vạn, người đã trực tiếp giúp đỡ, hướng dẫn và luôn sát sao đến em trong quá trình em làm đồ án tốt nghiệp. Em xin chân thành cảm ơn.

TÓM TẮT NỘI DUNG ĐỒ ÁN

Internet ngày càng được sử dụng nhiều trên toàn cầu khiến nó trở thành một phần không thể thiếu trong hoạt động hàng ngày của con người cả trong hoạt động cá nhân lẫn doanh nghiệp. Khi Internet ngày càng phổ biến đồng nghĩa với việc xuất hiện càng nhiều các ứng dụng. Đối với các nhà cung cấp dịch vụ mạng (Vietel, VNPT, v.v), phân loại lưu lượng mạng đóng một vai trò rất quan trọng trong việc quản lí mạng, vì nhà mạng có thể triển khai nhiều giải pháp hướng ứng dụng (truyền video, định tuyến, v.v) hay các giải pháp phát hiện bất thường trong hệ thống mạng. Với sự xuất hiện của các thuật toán học sâu (Deep Learning), rất nhiều nghiên cứu đã tập trung vào các giải pháp phân loại lưu lượng mạng sử dụng học sâu để phát hiện các kiểu ứng dụng đã biết.

Tuy nhiên, mỗi khi ứng dụng mới xuất hiện, nếu chỉ đào tạo lại mô hình với bộ dữ liệu mới thì mô hình không thể hoạt động tốt trên bộ dữ liệu cũ đã được đào tạo, còn việc đào tạo lại mô hình phân loại lưu lượng mạng với toàn bộ dữ liệu cũ và dữ liệu mới là không khả thi vì tốn kém về mặt thời gian và tiền bạc. Trong đề tài này, em đề xuất giải pháp phân loại lưu lượng mạng mã hóa sử dụng Class Incremental Learning nhằm giải quyết các hạn chế trên. Đề tài này gồm hai đóng góp chính. Thứ nhất, đề tài mô hình phân loại lưu lượng mạng mã hóa sử dụng thuật toán học sâu. Thứ hai, đề tài triển khai giải thuật Class Incremental Learning, iCaRL và đánh giá hiệu năng của các giải thuật này đối với bài toán phân loại lưu lượng mạng.

Sinh viên thực hiện
(Ký và ghi rõ họ tên)

MỤC LỤC

CHƯƠNG 1. GIỚI THIỆU ĐỀ TÀI.....	1
1.1 Đặt vấn đề.....	1
1.2 Các giải pháp hiện tại và hạn chế	1
1.3 Mục tiêu và định hướng giải pháp	2
1.3.1 Mục tiêu.....	2
1.3.2 Định hướng giải pháp	3
1.4 Đóng góp của đề án	3
1.5 Bố cục đề án	3
CHƯƠNG 2. NỀN TẢNG LÝ THUYẾT	5
2.1 Ngữ cảnh của bài toán.....	5
2.2 Các kết quả nghiên cứu tương tự	6
2.3 Mạng nơ ron tích chập.....	7
2.3.1 Tổng quan mạng nơ ron tích chập	7
2.3.2 Mô hình mạng nơ ron tích chập	7
2.4 Thuật toán iCaRL	11
2.4.1 Tổng quan về học gia tăng.....	11
2.4.2 Thách thức của học gia tăng	11
2.4.3 Chi tiết về thuật toán iCaRL.....	11
CHƯƠNG 3. PHƯƠNG PHÁP ĐỀ XUẤT.....	17
3.1 Tổng quan giải pháp.....	17
3.2 Thu thập và tiền xử lý dữ liệu.....	18
3.3 Mô hình phân loại.....	18
3.3.1 Mô hình CNN.....	18
3.3.2 Giải pháp phân loại lưu lượng mạng cho các lớp ứng dụng mới	20

CHƯƠNG 4. ĐÁNH GIÁ THỰC NGHIỆM.....	23
4.1 Các tham số đánh giá	23
4.2 Bộ dữ liệu	25
4.3 Phương pháp thí nghiệm.....	25
4.3.1 Môi trường lập trình	25
4.4 Kết quả thí nghiệm mô hình phân loại CNN và ResNet32.....	26
4.5 Kết quả thí nghiệm thuật toán iCaRL và học chuyển giao CNN.....	27
4.5.1 Kịch bản iCaRL	27
CHƯƠNG 5. KẾT LUẬN	33
5.1 Kết luận	33
5.2 Định hướng phát triển trong tương lai	33
TÀI LIỆU THAM KHẢO.....	35

DANH MỤC HÌNH VẼ

Hình 2.1	Ứng dụng mạng CNN trong phát hiện đối tượng	7
Hình 2.2	Kiến trúc mạng nơ ron tích chập cơ bản	8
Hình 2.3	Phép tích chập	8
Hình 2.4	Lớp gộp maxpooling	9
Hình 2.5	Lớp gộp avgpooling	10
Hình 2.6	Lớp kết nối đầy đủ	10
Hình 2.7	Thuật toán phân loại trong iCaRL	12
Hình 2.8	Thuật toán cập nhật lớp mới trong iCaRL	13
Hình 2.9	Thuật toán huấn luyện lớp mới trong iCaRL	14
Hình 2.10	Thuật toán chọn các ví dụ đại diện trong iCaRL	15
Hình 2.11	Thuật toán giảm kích thước bộ nhớ trong iCaRL	16
Hình 3.1	Luồng hoạt động của bài toán	17
Hình 3.2	Mô hình CNN được sử dụng trong bài toán	20
Hình 3.3	Mô hình transfer learning được sử dụng trong bài toán	21
Hình 3.4	Khối dư	22
Hình 4.1	Ví dụ về confusion matrix	24
Hình 4.2	Kết quả thu được khi đào tạo mô hình CNN	26
Hình 4.3	Confusion Matrix thu được khi đào tạo mô hình phân loại với ResNet32	27
Hình 4.4	Confusion matrix thu được khi bổ sung thêm lớp ứng dụng FileTransfer	28
Hình 4.5	Confusion matrix thu được khi bổ sung thêm lớp ứng dụng Music	29
Hình 4.6	Confusion matrix thu được khi bổ sung thêm lớp ứng dụng VoIP	30
Hình 4.7	Confusion matrix thu được khi bổ sung thêm lớp ứng dụng Photo	31

DANH MỤC BẢNG BIỂU

Bảng 4.1	Số lượng gói tin theo lớp ứng dụng	25
Bảng 4.2	Precision, Recall, F1-Score khi đào tạo với mô hình ResNet32	27
Bảng 4.3	Precision, Recall, F1-Score khi thêm lớp ứng dụng FileTransfer	28
Bảng 4.4	Precision, Recall, F1-Score khi thêm lớp ứng dụng Music . . .	29
Bảng 4.5	Precision, Recall, F1-Score khi thêm lớp ứng dụng VoIP . . .	30
Bảng 4.6	Precision, Recall, F1-Score khi thêm lớp ứng dụng Photo . . .	31

DANH MỤC THUẬT NGỮ VÀ TỪ VIẾT TẮT

Thuật ngữ	Ý nghĩa
CIL	Học tăng dần theo lớp (class incremental learning)
CNN	Mạng nơ ron tích chập
iCaRL	Phân loại tăng dần và học biểu diễn (Incremental Classifier and Representation Learning)

CHƯƠNG 1. GIỚI THIỆU ĐỀ TÀI

1.1 Đặt vấn đề

Internet ngày càng được sử dụng nhiều trên toàn cầu khiến nó trở thành một phần không thể thiếu trong hoạt động hàng ngày của con người cả trong hoạt động cá nhân lẫn doanh nghiệp. Theo báo cáo của Liên minh Viễn thông Quốc tế (ITU), người dùng internet trên toàn cầu đạt 4,1 tỷ vào năm 2019, tăng hơn 53% so với năm 2005. Internet phát triển đồng nghĩa với việc giám sát quản lý và bảo mật hệ thống đang gặp phải nhiều thách thức lớn. Chính vì vậy phân loại lưu lượng mạng rất cần thiết để các nhà cung cấp dịch vụ internet có thể lọc nội dung và xác định các hành vi hoặc ứng dụng có hại đối với người dùng hoặc tổ chức.

Khi internet ngày càng tăng đồng nghĩa với việc xuất hiện càng nhiều các ứng dụng. Mỗi khi ứng dụng mới xuất hiện, nếu chỉ đào tạo lại mô hình với bộ dữ liệu mới thì mô hình không thể hoạt động tốt trên bộ dữ liệu cũ đã được đào tạo, còn việc đào tạo lại mô hình phân loại lưu lượng mạng với toàn bộ dữ liệu cũ và dữ liệu mới là không khả thi vì tốn kém về mặt thời gian và tiền bạc. Chính vì vậy Class Incremental Learning được kỳ vọng giải quyết vấn đề gia tăng dữ liệu nhanh chóng.

Phân tích các gói tin mã hóa là một nhiệm vụ quan trọng đã được các nhà khoa học nghiên cứu từ lâu. Tuy nhiên các phương pháp truyền thống dựa vào các quy tắc và đặc điểm của gói tin gặp khó khăn trong việc phân loại và phát hiện các đe dọa từ các gói tin do sự hạn chế trong việc phát hiện các đặc điểm lưu lượng mạng được mã hóa mới.

Ngày nay, học sâu là một nhánh của trí tuệ nhân tạo đã cho thấy khả năng vượt trội trong việc phân tích và trích xuất các đặc điểm phức tạp của các bộ dữ liệu mà không cần sự can thiệp thủ công từ con người. Học sâu đã mang lại nhiều tiến bộ trong xử lý ngôn ngữ tự nhiên, thị giác máy tính, nhận dạng giọng nói, ... Việc áp dụng học sâu vào phân tích các gói tin mã hóa đã đạt được những tiến bộ nhất định mặc dù còn nhiều thách thức cần được giải quyết.

1.2 Các giải pháp hiện tại và hạn chế

Các phương pháp phân loại lưu lượng mạng truyền thống, chẳng hạn như dựa trên cổng, kiểm tra gói sâu và dựa trên thống kê đều bị hạn chế trong việc xác định các đặc điểm lưu lượng được mã hóa mới.

Giải pháp sử dụng cổng: Sử dụng cổng (port) để phân loại lưu lượng mạng mã hóa là phương pháp dựa trên giả định rằng các dịch vụ mạng cụ thể thường sử

dùng các cổng tiêu chuẩn như HTTPS sử dụng cổng 443, nó thường được sử dụng để truyền tải dữ liệu nhạy cảm trên web, SSL/TLS: các giao thức bảo mật được sử dụng để mã hóa lưu lượng truy cập mạng thường được sử dụng trên các cổng như 443(HTTPS), 993(IMAPS), ... Tuy nhiên các dịch vụ có thể được cấu hình để sử dụng các cổng không tiêu chuẩn, một cổng có thể được sử dụng cho nhiều loại dịch vụ khác nhau dẫn đến việc phân loại dựa trên cổng gặp khó khăn. Ngoài ra nếu chỉ dựa trên cổng để phân loại lưu lượng mạng sẽ không biết được thông tin chi tiết về nội dung của lưu lượng đó.

Kiểm tra gói sâu: Kiểm tra gói sâu (Deep Packet Inspection) là kỹ thuật phân tích tiêu đề và nội dung bên trong gói tin, nó có khả năng kiểm tra nội dung gói tin trong thời gian thực, giúp xác định chính xác loại ứng dụng hoặc giao thức mà khi dùng cổng không thể phát hiện. DPI thường được sử dụng trong các hệ thống tường lửa, hệ thống phát hiện xâm nhập hoặc trong các giải pháp quản lý băng thông hoặc giám sát an ninh mạng. So với phương pháp sử dụng cổng, DPI có độ chính xác cao hơn ngoài ra nó có thể giúp phát hiện các cuộc tấn công mạng hoặc phần mềm độc hại bằng cách phân tích sâu vào nội dung tin. Tuy nhiên ngoài sử dụng chi phí cao do yêu cầu tài nguyên tính toán mạnh mẽ, DPI không thể giải mã nội dung mã hóa mà không có khóa mã hóa.

Dựa trên thống kê: Phân loại lưu lượng mạng dựa vào thống kê là phương pháp phân tích các đặc trưng thống kê của lưu lượng mạng, chẳng hạn như kích thước gói tin, tần suất, độ trễ và các đặc trưng khác sau đó sử dụng các mô hình học máy để xác định loại ứng dụng hoặc giao thức đang được sử dụng. Mặc dù có nhiều ưu điểm hơn phương pháp sử dụng cổng và phương pháp kiểm tra gói sâu tuy nhiên phương pháp dựa trên thống kê vẫn tồn tại các hạn chế: cần dữ liệu huấn luyện để xây dựng mô hình học máy chính xác, khó khăn khi lưu lượng biến đổi và yêu cầu tài nguyên tính toán lớn.

Các phương pháp phân loại dựa trên học sâu xem xét các tính năng dựa trên gói được khám phá để giải quyết các hạn chế của các phương pháp truyền thống khi phân loại lưu lượng mạng mã hóa.

1.3 Mục tiêu và định hướng giải pháp

1.3.1 Mục tiêu

Với những vấn đề và bối cảnh hiện tại được nêu ra ở các mục trên, mục tiêu đặt ra của đề án là đề xuất, triển khai và đánh giá được mô hình phân loại lưu lượng mạng mã hóa. Thứ nhất mô hình đề xuất phải phân loại chính xác được các lưu lượng mạng thuộc các lớp đã được đào tạo. Thứ hai, với những lớp lưu lượng mạng mới chưa gặp bao giờ, mô hình đề xuất phải có khả năng phân loại nó mà không