

ỨNG DỤNG PHƯƠNG PHÁP TÍNH ENTROPY THÔNG TIN TRONG VIỆC PHÒNG CHỐNG RANSOMWARE

Hoàng Văn Quyết*, Đoàn Văn Minh

***Tóm tắt:** Bài báo trình bày phương pháp tính Entropy thông tin, tiến hành thực nghiệm xác định giá trị Entropy của các tập tin dữ liệu nguyên bản và khi bị mã hóa bởi mã độc tống tiền trong hệ thống mạng máy tính. Đề xuất giải pháp và xây dựng phần mềm bảo vệ dữ liệu máy tính trước sự tấn công của Ransomware dựa trên phương pháp tính Entropy thông tin.*

Từ khóa: Entropy thông tin, Entropy liên tục, Phần mềm độc hại, Mã hóa dữ liệu.

1. MỞ ĐẦU

Ngày nay, sự thuận tiện của các phương thức thanh toán điện tử, cũng như công nghệ mã hóa hiện đại là cơ sở để mã độc tống tiền (Ransomware) bùng phát và có xu hướng ngày càng gia tăng, gây hậu quả nghiêm trọng cho người sử dụng máy tính. Chính vì vậy, bảo vệ dữ liệu điện tử khỏi Ransomware là một vấn đề rất cấp thiết.

Bài báo này không đề cập đến phương pháp tiêu diệt Ransomware (một phương pháp mà các hãng bảo mật lớn trên thế giới thường làm), hay phương pháp khắc phục hậu quả khi dữ liệu đã bị mã hóa bởi Ransomware do công nghệ mã hóa ngày càng hiện đại. Tác giả chỉ tập trung vào phương pháp bảo vệ dữ liệu máy tính khi bị nhiễm Ransomware, một phương pháp hiện nay chưa có hãng bảo mật nào nghiên cứu.

Đặc điểm đặc trưng của dữ liệu mã hóa bởi Ransomware đó là mức độ ngẫu nhiên lớn và để đo mức độ ngẫu nhiên đó có thể sử dụng đại lượng Entropy thông tin của C. E. Shannon [3]. Từ công thức tính Entropy của C. E. Shannon, tác giả xây dựng công thức tính Entropy cho tập tin dữ liệu được lưu trên bộ nhớ, làm cơ sở tiến hành tính thực nghiệm lượng Entropy của các tập tin dữ liệu cụ thể, từ đó rút ra ngưỡng Entropy phân biệt giữa tập tin mã hóa và chưa mã hóa. Trên những cơ sở lý thuyết và thực nghiệm đó, tác giả đã phát triển chương trình bảo vệ dữ liệu máy tính khi bị nhiễm Ransomware.

2. PHƯƠNG PHÁP TÍNH ENTROPY THÔNG TIN

2.1. Khái quát về Entropy và Ransomware

Năm 1927, Von Neumann đã xây dựng công thức thống kê trong nhiệt động lực học và cơ học có chứa giá trị Entropy, tuy nhiên, đến năm 1948, khái niệm Entropy thông tin mới được C. E. Shannon đưa ra trong bài báo "A Mathematical

Theory of Communication".

Entropy thông tin là một khái niệm mở rộng của Entropy trong nhiệt động lực học và cơ học thống kê, được áp dụng sang lý thuyết thông tin, mô tả mức độ hỗn loạn trong tín hiệu lấy từ một sự kiện ngẫu nhiên, cho phép chỉ ra số lượng thông tin (các phần không hỗn loạn ngẫu nhiên) có trong tín hiệu. Xem xét trường hợp cụ thể, một câu có ý nghĩa được viết bằng tiếng Việt và được thể hiện bởi các ký tự (chữ cái, khoảng cách và dấu câu), sẽ không hiện ra một cách hoàn toàn hỗn loạn ngẫu nhiên. Cụ thể, tần số xuất hiện của ký tự "x" sẽ không giống với tần số xuất hiện của ký tự phổ biến hơn là "t", đồng thời, nếu dòng chữ vẫn đang được viết hay đang truyền tải, sẽ khó đoán trước được ký tự tiếp theo. Việc xuất hiện ký tự tiếp theo có mức độ ngẫu nhiên nhất định, giá trị này được xác định bởi Entropy thông tin.

Ransomware là một loại phần mềm độc hại sử dụng hệ thống mật mã để mã hóa dữ liệu, ngăn chặn người dùng truy cập và sử dụng máy tính, yêu cầu nạn nhân phải nộp một khoản tiền chuộc nếu muốn lấy lại dữ liệu [1].

Ransomware mã hóa dữ liệu bằng cách đọc dữ liệu lên bộ nhớ RAM, mã hóa dữ liệu, sau đó ghi ngược xuống ổ đĩa. Quá trình mã hóa thành công khi dữ liệu mã hóa trên bộ nhớ RAM được ghi ngược hoàn toàn xuống ổ đĩa. Dữ liệu mã hóa bởi Ransomware sẽ có mức độ ngẫu nhiên (mức độ hỗn loạn) của các byte lớn hơn dữ liệu chưa bị mã hóa. Mức độ phức tạp của việc giải mã các dữ liệu tỉ lệ thuận với mức độ ngẫu nhiên [2].

2.2. Cơ sở nền tảng tính Entropy thông tin

C. E. Shannon [3] đã xây dựng định nghĩa về Entropy thông tin để thỏa mãn các giả định sau:

- Entropy phải tỷ lệ thuận liên tục với các xác suất xuất hiện của các phần tử ngẫu nhiên trong tín hiệu. Thay đổi nhỏ trong xác suất xuất hiện sẽ dẫn đến thay đổi nhỏ trong Entropy;
- Nếu các phần tử ngẫu nhiên đều có xác suất xuất hiện bằng nhau, việc tăng số lượng phần tử ngẫu nhiên phải làm tăng Entropy;
- Các chuỗi tín hiệu có thể tạo ra theo nhiều bước và giá trị Entropy phải bằng tổng trọng số của các Entropy của từng bước.

C. E. Shannon chỉ ra rằng, các định nghĩa Entropy cho một tín hiệu có thể nhận các giá trị rời rạc, thỏa mãn các giả định trên, đều được tính theo công thức:

$$-K \sum_{i=1}^n p(i) \log_2 p(i) \quad (1)$$

với K là một hằng số, chỉ phụ thuộc vào đơn vị đo; n là tổng số các giá trị có thể

nhận của tín hiệu; i là giá trị rời rạc thứ i ; $p(i)$ là xác suất xuất hiện của giá trị i .

2.3. Xác định Entropy trong trường hợp ngẫu nhiên rời rạc

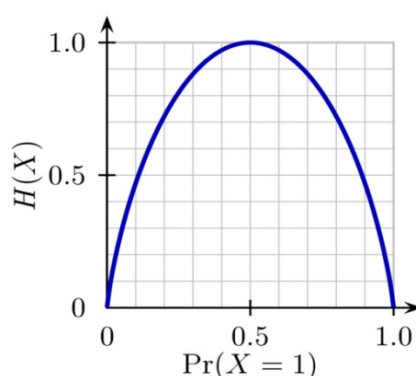
Nếu một sự kiện ngẫu nhiên rời rạc x , có thể nhận các giá trị là $\overline{1..n}$, thì Entropy $H(x)$ của nó được tính theo công thức (2):

$$H(x) = \sum_{i=1}^n p(i) \log_2\left(\frac{1}{p(i)}\right) = - \sum_{i=1}^n p(i) \log_2 p(i) \quad (2)$$

Như vậy, Entropy của x cũng là giá trị kì vọng của các độ ngẫu nhiên của các giá trị mà x có thể nhận. Entropy thông tin trong trường hợp phân tử tín hiệu ngẫu nhiên rời rạc còn được gọi là Entropy Shannon.

Thực nghiệm kết quả tính Entropy của phép thử Bernoulli: $X=1$ với xác suất p và $X=0$ với xác suất $(1-p)$, khi đó:

$$H(x) = H(p) = -p \log_2 p - (1-p) \log_2(1-p) \quad (3)$$



Hình 1. Biểu đồ giá trị Entropy rời rạc của X .

2.4. Xác định giá trị Entropy trong trường hợp ngẫu nhiên liên tục

Nếu x là số thực ngẫu nhiên liên tục, thì Entropy liên tục được tính theo công thức:

$$H(f) = \int_{-\infty}^{+\infty} f(x)(\log f(x))dx \quad (4)$$

với f là hàm mật độ xác suất.

Entropy liên tục thường được gọi là Entropy vi phân hay Entropy Boltzmann [4]. Entropy Boltzmann không phải là giới hạn của Entropy Shannon khi $n \rightarrow \infty$, do đó không phải là độ đo mức độ hỗn loạn của thông tin.

Trường hợp một dòng chữ luôn chỉ có các ký tự "a" sẽ có Entropy bằng 0, vì ký tự tiếp theo sẽ luôn là "a". Một dòng chữ chỉ có hai ký tự 0 và 1 ngẫu nhiên hoàn toàn sẽ có Entropy là 1 bit cho mỗi ký tự.

Một dòng chữ tiếng Anh thông thường có Entropy khoảng 1,1 đến 1,6 bit cho mỗi ký tự. Thuật toán nén PPM [4] có thể tạo ra tỷ lệ nén 1,5 bit cho mỗi ký tự. Trên thực tế, tỷ lệ nén của các thuật toán nén thông dụng có thể được dùng làm ước lượng cho Entropy của dữ liệu.

Entropy của dòng văn bản S thông thường được định nghĩa dựa trên mô hình Markov. Nếu các ký tự tiếp theo hoàn toàn độc lập với các ký tự trước đó, Entropy nhị phân $H[S]$ sẽ là:

$$H[S] = \sum p(i) \log_2 p(i) \quad (5)$$

2.5. Công cụ phân tích Entropy nhị phân

Bintropy là công cụ phân tích mẫu, ước tính khả năng một tệp tin có chứa các thông tin nén hoặc mã hóa. Bintropy có hai chế độ hoạt động:

- Chế độ thứ nhất, công cụ sẽ phân tích Entropy của mỗi đoạn thực thi có định dạng PE, được xác định trong phần đầu của tệp thực thi. Điều này giúp người phân tích xác định đoạn mã thực thi có thể bị mã hóa và nén. Một bộ biên dịch chuẩn tạo ra PE thực thi có các phần theo định dạng chuẩn (.text, .data, .reloc, .rsrc). Tuy nhiên, nhiều công cụ đóng gói biến đổi định dạng của tệp thực thi gốc, nén các đoạn mã, dữ liệu và dồn chúng vào một hay hai đoạn mới. Trong chế độ này, Bintropy tính giá trị Entropy cho mỗi đoạn nó cần. Tuy nhiên, không tính Entropy cho phần đầu tệp tin bởi vì phần này không chứa các byte dữ liệu nén hay mã hóa;

- Chế độ thứ hai bỏ qua định dạng tệp, thay vào đó Bintropy phân tích Entropy của toàn bộ tệp, từ byte đầu tiên cho đến byte cuối cùng. Với tệp định dạng PE, người dùng có thể phân tích Entropy của đoạn mã và dữ liệu ẩn tại cuối tệp hoặc ở giữa các đoạn định dạng PE.

Entropy của một khối dữ liệu là một phép thống kê lượng thông tin chứa bên trong. Trong bài báo “Sử dụng phân tích Entropy để tìm ra mã độc nén và mã hóa”, hai tác giả Hamrock và Lyda đưa ra một quan sát đáng chú ý là các dữ liệu nén và mã hóa trong mẫu mã dữ liệu độc hại đóng gói có mức Entropy cao. Mã chương trình và dữ liệu bình thường có mức Entropy thấp hơn nhiều. Mã độc hại sử dụng kỹ thuật đóng gói được xác định bởi mức Entropy cao trong nội dung của nó.

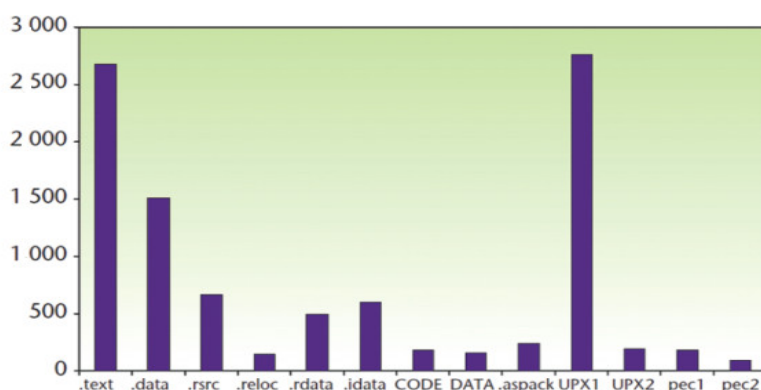
Để đánh giá khả năng công cụ Bintropy dựa trên phân tích Entropy, Lynda và Hamrock đã tiến hành đánh giá thử nghiệm trên bốn tập dữ liệu với các phân loại tệp khác nhau: plain text, thực thi thông thường, thực thi nén và thực thi mã hóa. Mỗi tập dữ liệu gồm 100 tệp khác nhau, mỗi tệp được tính Entropy dựa trên các khối dữ liệu có độ dài 256 byte. Công cụ Bintropy tính Entropy mức trung bình của các khối và khối có mức Entropy cao nhất. Mục đích thử nghiệm này là xác

định mức Entropy tối ưu để phân loại tệp thực thi thông thường và tệp thực thi đã biến đổi sử dụng kỹ thuật mã hóa hoặc kỹ thuật nén. Sau khi sử dụng tập dữ liệu training, Bintropy có khả năng phát hiện các tệp thực thi bị nén hoặc mã hóa khi đặc tính Entropy vượt qua một mức định trước.

Bảng 1. Độ chính xác thống kê Entropy dựa trên tập dữ liệu.

Tập dữ liệu	Mức Entropy trung bình	Entropy trong khoảng	Mức Entropy cao nhất
Plain text	4.347	4.066–4.629	4.715
Thực thi thông thường	5.099	4.941–5.258	6.227
Thực thi nén	6.801	6.677–6.926	7.233
Thực thi mã hóa	7.175	7.174–7.177	7.303

Bảng kết quả cho thấy, với độ chính xác đạt 99% và mức Entropy trong khoảng 6.677 đến 7.177, công cụ Bintropy sẽ phát hiện tệp nén hay mã hóa.



Hình 2. Phân bố số lượng tệp mã độc theo đoạn (session) bị mã hóa hoặc nén.

Lyda và Hamrock cũng đã thực hiện xác định xu hướng Entropy bởi công cụ Bintropy và để tạo độ tin cậy của đánh giá đã áp dụng trên một tập 21.567 mã độc Win32 - với thực thi định dạng PE từ bộ thu thập của các hãng phần mềm chống virus nổi tiếng trên thế giới trong khoảng thời gian từ tháng 01/2000 đến tháng 12/2005. Dựa trên khảo sát bởi sử dụng công cụ Bintropy để phân tích, kết quả chỉ ra rằng, UPX1 là phần được kẻ viết mã độc sử dụng kỹ thuật đóng gói phổ biến nhất, sau đó là phần text (hình 2).

2.6. Xây dựng công thức tính Entropy cho dữ liệu lưu trên RAM

Dữ liệu lưu trên bộ nhớ máy tính được mô tả dưới dạng các byte nhớ, một byte

có giá trị từ 1 đến 255, do đó đơn vị rời rạc i sẽ có giá trị từ 1 đến 255.

Từ thực nghiệm và áp dụng công thức tính Entropy nhị phân (5), tác giả đã xây dựng công thức tính lượng Entropy cho một tập tin dữ liệu được lưu trên RAM, cụ thể như sau:

$$Entropy = -K \sum_{i=1}^{255} p(i) \log_2 p(i) \quad (6)$$

với hằng số $K = 1.4426950408$; $p(i)$ là xác suất xuất hiện của giá trị i (từ 1 đến 255) trên độ dài của tập tin dữ liệu được lưu trên RAM.

Dựa trên công thức tính Entropy thông tin, đã tiến hành kiểm tra thực nghiệm đối với các tập tin dữ liệu, trong trạng thái nguyên bản và trạng thái bị mã hóa bởi Ransomware và được lưu trữ ở các định dạng khác nhau như *.doc*, *.pdf*, *.htm*, *.xls*. Kết quả thực nghiệm được thể hiện trong bảng 2.

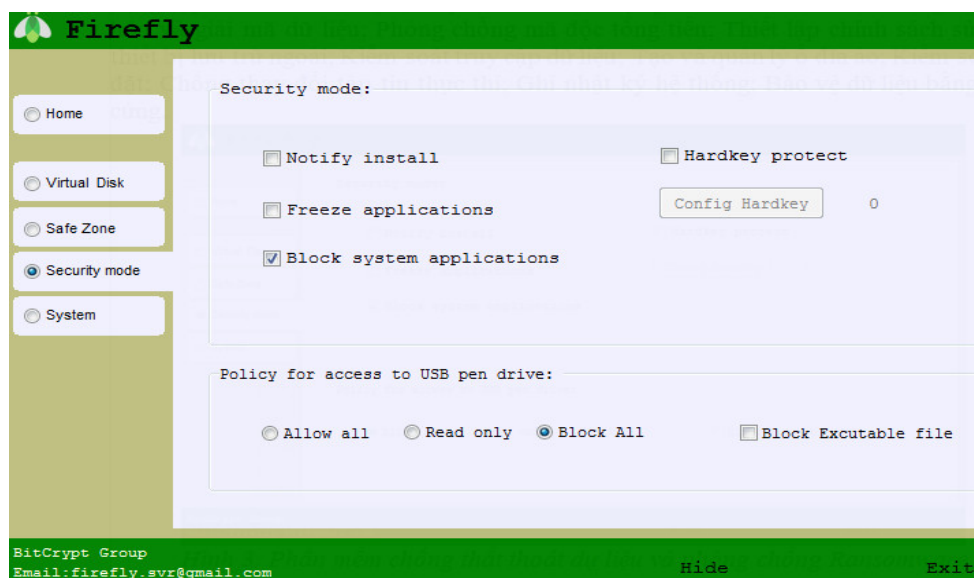
Bảng 2. Kết quả thực nghiệm tính Entropy cho các tập tin có định dạng khác nhau.

Tập tin đã bị mã hóa	Tập tin định dạng .pdf	Tập tin định dạng .doc	Tập tin định dạng .htm	Tập tin định dạng .xls
9.230	1.631	0.101	0.104	0.097
9.408	1.178	0.039	0.104	0.107
10.000	1.458	0.039	0.109	0.097
10.000	2.928	0.068	0.098	0.109
10.000	1.055	0.067	0.097	0.109
9.454	1.648	0.071	0.100	0.109
10.000	2.108	0.119	0.082	0.076
10.000	1.468	0.070	0.105	0.102
10.000	1.155	0.039	0.104	0.098

Kết quả thực nghiệm cho thấy sự khác biệt về lượng Entropy của tập tin dữ liệu bị mã hóa bởi Ransomware so với của các tập tin dữ liệu chưa bị mã hóa. Kết quả này cho phép xác định được ngưỡng Entropy để làm cơ sở thiết lập quyền ghi dữ liệu từ Ram xuống ổ đĩa. Nếu lượng Entropy lớn hơn ngưỡng cho phép, việc ghi xuống ổ đĩa bị cấm, ngược lại việc ghi xuống ổ đĩa diễn ra như bình thường.

Trên nền tảng tính toán Entropy thông tin, tác giả đã xây dựng phần mềm bảo vệ dữ liệu có tên là Firefly [6], bảo vệ người dùng trước các cuộc tấn công của Ransomware, hành động có chủ đích của tin tặc bằng cách xác định các ngưỡng Entropy của các định dạng tập tin khác nhau (hình 3). Phần mềm được phát triển có các chức năng sau: Mã hóa và giải mã dữ liệu; Phòng chống mã độc tổng tiền;

Thiết lập chính sách sử dụng thiết bị lưu trữ ngoài; Kiểm soát truy cập dữ liệu; Tạo và quản lý ổ đĩa ảo; Kiểm soát cài đặt; Chống thay đổi tập tin thực thi; Ghi nhật ký hệ thống; Bảo vệ dữ liệu bằng khóa cứng.



Hình 3. Phần mềm chống thất thoát dữ liệu và phòng chống Ransomware.

3. KẾT LUẬN

Trong nhiệm vụ bảo mật, bảo đảm an ninh, an toàn thông tin, việc xây dựng các hệ thống, ứng dụng, dịch vụ ngăn cản các hành vi của mã độc là rất quan trọng. Tác giả đã giới thiệu phương pháp tính Entropy thông tin cho tín hiệu, làm cơ sở cho phương pháp tính Entropy thông tin cho dữ liệu máy tính. Kết quả nghiên cứu, thực nghiệm tính toán giá trị Entropy dữ liệu máy tính, đã xác định được ngưỡng Entropy để thiết lập quyền ghi xuống ổ đĩa, từ đó, cho phép xây dựng các phần mềm, dịch vụ nhằm ngăn chặn Ransomware mã hóa dữ liệu trên máy tính người dùng.

TÀI LIỆU THAM KHẢO

- [1]. N. M. Abramson, “*Information Theory and Coding*”, McGraw-Hill, New York, 1963.
- [2]. M. Sikorski, “*Practical Malware Analysis*”, Physical model and science Journal, **vol. 10**, No.1 (1995), pp. 79-93.
- [3]. C. E. Shannon, “*Prediction and Entropy of Printed English*”, Bell System Technical Journal, **vol.30**, No. 1 (1951), pp. 50-64.
- [4]. C. Tom, “*An introduction to information theory and Entropy*”, Physical

technology (2014), pp. 67-80.

[5]. C. E. Shannon, “*A Mathematical Theory of Communication*”, Bell System Technical Journal, **vol. 27**, No.3 (1953), pp. 379-423.

[6]. Firefly, Data Threat Prevention, <http://firefly.com.vn>.

ABSTRACT

THE USAGE OF INFORMATION ENTROPY CALCULATION IN RANSOMWARE PREVENTION

This article shows how to calculate information Entropy changes for the original files and the files that were encrypted by Ransomware in a computer system. It is also proposing a method using information Entropy calculation in order to develop software that protects data from Ransomware.

Keywords: Information Entropy, Boltzmann’s Entropy, Malware, Data encryption.

Nhận bài ngày 15 tháng 2 năm 2017

Hoàn thiện ngày 15 tháng 3 năm 2017

Chấp nhận đăng ngày 01 tháng 5 năm 2017

Địa chỉ: Cục Công nghệ thông tin, BTTM.

*Email: hoangquyetik@gmail.com.