

네트워크 플로우의 연관성 모델을 이용한 트래픽 분류 방법

(A Traffic-Classification Method Using
the Correlation of the Network Flow)

구영훈[†] 심규석^{**} 이성호^{***} Baraka D. Sija^{***} 김명섭^{****}
(YoungHoon Goo) (Sungho Lee) (Kyuseok Shim) (Baraka D. Sija) (MyungSup Kim)

요약 오늘날의 네트워크는 고속화와 유비쿼터스 환경으로 인해 다양한 응용이 급속도로 생성되고 있으며 네트워크 트래픽도 매우 복잡해지고 있다. 이에 효율적인 네트워크 운용 및 관리를 위한 구체적인 단위의 트래픽 분류가 필수적이다. 다양한 트래픽 분류 방법이 연구되고 있는 가운데 아직 트래픽을 완벽하게 분류해내는 방법론은 개발되지 않은 실정이다. 이에 본 논문에서는 네트워크 플로우의 연관성 모델을 정의하고 이를 기반으로 트래픽을 분류하는 방법을 제안한다. 트래픽 분류를 위한 네트워크 플로우의 연관성 모델은 크게 유사성 모델과 연결성 모델로 이루어진다. 제안하는 방법론을 효과적으로 적용하기 위한 방안을 제시하며 실험을 통해 본 분류 방법론이 높은 정확도와 분석력의 방법론이라는 것을 증명한다.

키워드: 트래픽 분류, 네트워크 플로우의 연관성 모델, 유사성 모델, 연결성 모델, 가이드라인

Abstract Presently, the ubiquitous emergence of high-speed-network environments has led to a rapid increase of various applications, leading to constantly complicated network traffic. To manage networks efficiently, the traffic classification of specific units is essential. While various traffic-classification methods have been studied, a methods for the complete classification of network traffic has not yet been developed. In this paper, a correlation model of the network flow is defined, and a traffic-classification method for which this model is used is proposed. The proposed network-correlation model for traffic classification consists of a similarity model and a connectivity model. Suggestion for the effectiveness of the proposed method is demonstrated in terms of accuracy and completeness through experiments.

Keywords: traffic classification, correlation model of the network flow, similarity model, connectivity model, guideline

· 본 연구는 이 논문은 2015년도 정부(교육부)의 재원으로 한국연구재단 기초연구사업(No.2015R1D1A3A01018057) 및 2016년도 정부(미래창조과학부)의 재원으로 한국연구재단의 지원을 받아 수행된 과학기술인문융합연구사업임(No.NRF-2016M3C1B6929228)

· 이 논문은 제43회 통계학술발표회에서 '네트워크 플로우의 연관성 모델을 이용한 트래픽 분류 방법'의 제목으로 발표된 논문을 확장한 것임

논문접수 : 2017년 1월 18일

(Received 18 January 2017)

논문수정 : 2017년 2월 22일

(Revised 22 February 2017)

심사완료 : 2017년 2월 22일

(Accepted 22 February 2017)

[†] 비회원 : 고려대학교 컴퓨터정보학과
gyh0808@korea.ac.kr

^{**} 학생회원 : 고려대학교 컴퓨터정보학과
kusuk007@korea.ac.kr

^{***} 비회원 : 고려대학교 컴퓨터정보학과
gaek5@korea.ac.kr
sijabarakajia25@korea.ac.kr

^{****} 종신회원 : 고려대학교 컴퓨터정보학과 교수(Korea Univ.)
tmskim@korea.ac.kr
(Corresponding author임)

Copyright©2017 한국정보과학회 : 개인 목적이거나 교육 목적인 경우, 이 저작물의 전체 또는 일부에 대한 복사본 혹은 디지털 사본의 제작을 허가합니다. 이 때, 사본은 상업적 수단으로 사용할 수 없으며 첫 페이지에 본 문구와 출처를 반드시 명시해야 합니다. 이 외의 목적으로 복제, 배포, 출판, 전송 등 모든 유형의 사용행위를 하는 경우에 대하여는 사전에 허가를 얻고 비용을 지불해야 합니다.
정보과학회논문지 제44권 제4호(2017. 4)

1. 서론

오늘날의 네트워크는 고속화와 유비쿼터스 환경으로 인하여 다양한 응용이 급속도로 생성되고 있으며 이에 따라 네트워크 트래픽도 매우 복잡해지고 있다. 이러한 상황 속에서 네트워크 트래픽 모니터링과 분석은 효율적인 네트워크 운영과 안정적 서비스 제공에 있어 필수적이다. 인터넷 트래픽의 정확한 응용별 분류는 다양한 트래픽 분석 요구에 대처하기 위해 반드시 선행되어야만 한다.

다양한 트래픽 분류 방법이 연구되고 있는 가운데 아직 트래픽을 완벽하게 분류해내는 방법론은 개발되지 않은 실정이다. 먼저, 헤더 시그니처 기반 분류 방법[1]은 다양한 프로토콜을 사용하는 응용, 둘 이상의 포트 또는 임의의 포트를 설정할 수 있는 기능을 제공하는 응용 등 현대의 복잡한 구조를 갖는 응용에 대해서는 신뢰성을 갖기 힘들다. 페이로드 시그니처 기반 분류 방법[2]은 분석률과 정확도 측면에서 가장 높은 분석 성능을 보이지만 시그니처 추출 작업이 수작업으로 이루어져 응용 프로그램의 변화에 신속하게 대처할 수 없으며, 시그니처 생성 시 전문성을 가진 인력과 많은 시간이 필요하고 시그니처가 생성자의 능력에 큰 차이가 있을 수 있다. 또한 암호화된 트래픽은 분류할 수 없는 큰 단점을 가지고 있다. 이를 극복하기 위한 통계 시그니처 기반 분류 방법[3]은 암호화된 트래픽을 분류할 수 있지만, 통계 정보를 사용할 경우, 특정 응용 프로그램에서 사용한 엔진 또는 응용 프로토콜에 의존적인 시그니처가 생성될 가능성이 높다. 앞서 언급한 세 가지 분류 방법은 모두 시그니처가 매칭되는 플로우만을 분석하므로, 미탐되는 플로우가 다수 존재하는 것을 방지하기 위해서는 다수의 시그니처를 사용하여야 하며 이는 처리 시간의 증가를 야기한다. 또한, 기계학습 기반 트래픽 분류 방법[4,5]은 많은 양의 샘플 데이터 확보가 정확도에

큰 영향을 미치며 동일한 응용 계층 프로토콜을 이용하는 응용에 대해서는 분류가 어려운 문제가 발생된다.

본 논문에서는 네트워크 플로우의 연관성 모델을 정의하고 이를 기반으로 트래픽을 분류하는 방법을 제안한다. 제안하는 방법론은 입력 트래픽에 대해 연관성을 자동으로 계산하여 그룹화하므로 분류 속도가 빠르고 암호화된 트래픽 분석이 가능하며 연속적인 그룹화로 분석률을 최대화할 수 있다.

본 논문은 본 장의 서론에 이어 2장에서 제안하는 방법론에 대해 기술하고, 3장에서는 본 방법론의 효과적인 적용 방안에 대해 기술한다. 4장에서는 본 방법론의 타당성을 실험을 통해 증명하고, 마지막으로 5장에서는 결론 및 향후 연구를 기술한다.

2. 네트워크 플로우의 연관성 모델을 이용한 트래픽 분류 방법

본 장에서는 제안하는 방법론을 기술한다. 네트워크 플로우의 연관성 모델은 크게 유사성 모델과 연관성 모델로 구성된다. 그림 1은 제안하는 방법론의 분류 체계도이다. 먼저 타겟 트래픽과 타겟 트래픽에서 분류하고자 하는 탐지 대상(응용 또는 악성행위)의 플로우 중 일부를 Seed Group으로 선정하여 유사성 모델에 입력한다. Seed Group은 악성행위의 경우 IDS 또는 IPS의 경보, 응용의 경우 다양한 시그니처를 통해 얻은 정보를 통해 선정할 수 있다. 유사성 모델에서는 Seed Group과 타겟 트래픽과의 유사성을 계산하여 일정 임계치 이상인 트래픽을 그룹화하여 Group1을 출력한다. 다음으로 Group1을 연결성 모델의 입력으로 사용하면 Group1과 그룹화되지 않은 트래픽과의 연결성을 계산하여 일정 임계치 이상인 트래픽을 그룹화하여 Group2를 출력한다. 더 이상 그룹화되지 않을 때까지 해당 출력인

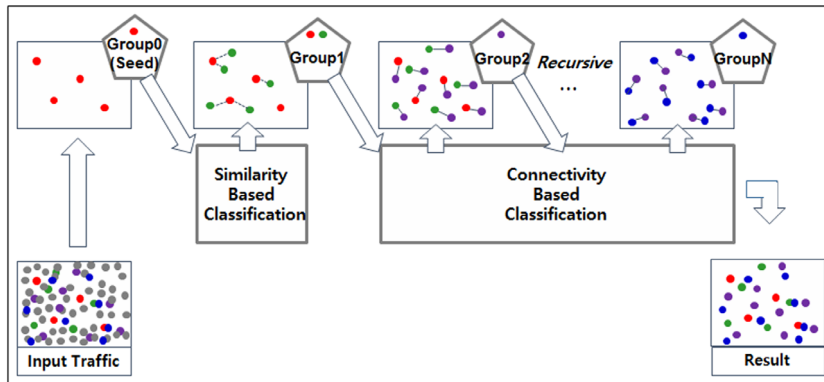


그림 1 플로우의 연관성을 이용한 트래픽 분류 체계

Fig. 1 Process of traffic classification using the flow correlation

Group을 반복적으로 연결성 모델에 입력한다. 최종적으로 Seed Group과 Group1부터 GroupN까지의 트래픽을 분류 결과로 출력한다.

플로우간 유사성은 플로우의 통계 정보들을 벡터로 구성하여 유클리드 거리 값으로 계산한다. 플로우간 유사성에 사용되는 통계 정보들은 다른 응용을 중복 탐지하지 않고 특정 응용만을 유일하게 탐지할 수 있도록 선정하는 것이 중요하다. 이를 위하여 유사성에 사용되는 각 통계 정보뿐만 아니라 특정 응용만을 유일하게 구분할 수 있도록 사용되는 통계 정보의 개수도 경우에 따라서 증가시킬 필요도 있다. 본 연구진은 이를 위하여 다양한 플로우의 통계 정보를 유사성 모델로써 적용시켜 보았으며 최종적으로 21가지 플로우의 통계 정보를 유사성 모델의 속성으로 선정하였다. 사용되는 통계 정보들은 플로우를 구성하는 모든 패킷의 inter-arrival time, 처음 16개 패킷간 inter-arrival time, 모든 패킷 페이로드 크기, 첫 5개 패킷의 페이로드 크기, 이 4가지 항목의 최대값, 최소값, 평균, 표준편차 총 16개의 속성들과 첫 5개 패킷의 방향성을 고려한 페이로드 크기를 사용하여 총 21가지의 속성을 사용한다. 방향성을 고려한 페이로드 크기란 TCP의 경우, 클라이언트에서 서버로 전송되는 패킷에는 페이로드 크기에 +1을 곱한 값을 의미하며 반대의 경우 -1을 곱한 값을 의미하고, UDP의 경우, 두 호스트 중 첫 패킷의 방향을 + 방향이라고 규정한다. 식 (1)은 플로우 f_x 와 f_y 의 유사성을 계산하는 식이다.

$$\begin{aligned} Sim(f_x, f_y) &= 1 - \frac{Euclidean\ Distance}{\sqrt{21}} \\ &= 1 - \frac{\sqrt{\sum_{i=1}^{21} (a_i(f_x) - a_i(f_y))^2}}{\sqrt{21}} \end{aligned}$$

식 (1) 플로우간 유사성 계산

Expr 1. Calculation of the similarity

a_i 는 각 속성을 의미하며 각 속성들의 가중치를 동일하게 하기 위하여 Min-Max 정규화를 각 속성에 적용하여 0~1 사이의 값으로 척도를 맞춘다. 계산되는 유클리드 거리를 다시 Min-Max 정규화를 통해 0~1 사이의 값으로 나타내고 비유사도를 유사도로 전환하기 위해 1에서 해당 값을 뺀다. 최종적으로 플로우간 유사성은 0~1 사이의 값으로 표현되고 두 플로우가 유사할수록 1에 가까운 값을 가진다. 유사성 모델에서는 Seed Group과 타겟 트래픽의 유사성을 계산하여 특정 임계치 이상의 유사성을 가진 플로우들을 그룹화한다.

플로우간 연결성은 플로우의 발생시간, 출발지 IP 주소와 도착지 IP 주소쌍, 출발지 포트와 도착지 포트쌍,

L4 프로토콜 총 4개의 속성 값을 이용하여 계산한다. 식 (2)는 플로우 f_x 와 f_y 의 연결성을 계산하는 식이다.

$$\begin{aligned} Conn(f_x, f_y) &= \sum_{i=1}^4 (w_i \times a_i(f_x, f_y)), \\ \text{where } 0 \leq a_i(f_x, f_y) &\leq 1, \sum_{i=1}^4 w_i = 1 \end{aligned}$$

식 (2) 플로우간 연결성 계산

Expr 2. Calculation of the connectivity

w_i 는 속성별 가중치를 나타낸다. 가중치의 합은 1이며 탐지 대상의 종류에 따라 발생하는 트래픽이 상이할 수 있음을 염두하여 속성들의 가중치를 조절할 수 있다. 대상별에 맞는 가중치 조절을 통해 분류 성능을 최대화할 수 있다. a_i 는 연결성의 각 속성을 의미한다.

$$a_{ST}(f_x, f_y) = 1 - \sqrt{\frac{dist(f_x, f_y)}{maxdist(F)}}$$

식 (3) 시작시간 연결성 계산

Expr 3. Calculation of the connectivity of the start time

식 (3)은 플로우간 시작시간의 연결성을 계산하는 식이다. $dist(f_x, f_y)$ 는 플로우의 시작시간 차를 의미하며 $maxdist(F)$ 는 타겟 트래픽의 전체 플로우들의 시작시간 차 중 최대값을 의미한다. 본 수식의 의미는 모든 플로우중 최대값과 최소값인 0, 그리고 f_x, f_y 의 시작 시간 차이를 가지고 min-max 정규화를 통해 0~1 사이의 값으로 표현한 값이다. 또한 같은 탐지 대상의 플로우이지만 시작 시간 차이가 너무 클 경우, 다른 응용으로 판별하게 되는 가능성이 커지는 것을 완화시키기 위하여 제곱근을 통해 최종적으로 0~1 사이의 값으로 표현한다. 본 수식은 0에 가까울수록 두 플로우가 유사하지 않다는 것을 의미하며 1에 가까울수록 유사하다는 것을 의미한다.

$$a_{IP}(f_x, f_y) = \left(\frac{prefixlenIP_{src}(f_x, f_y)}{32} \right)^2 + \left(\frac{prefixlenIP_{dst}(f_x, f_y)}{32} \right)^2$$

식 (4) 출발지, 도착지 IP 주소 쌍 연결성 계산

Expr 4. Calculation of the connectivity of the source IP and the destination IP

식 (4)는 출발지, 도착지 IP주소 쌍 연결성을 계산하는 수식이다. 식 (4)의 $prefixlenIP_{src}(f_x, f_y)$ 는 f_x 의 출발지 주소와 f_y 의 출발지 주소 32비트 중 앞부분이 같은 비트의 개수와 f_x 의 출발지 주소와 f_y 의 도착지 주

소 중 큰 값을 의미하며, $prefixlenIP_{dst}(f_x, f_y)$ 는 f_x 의 도착지 주소와 f_y 의 도착지 주소 32비트 중 앞부분이 같은 비트의 개수와 f_x 의 도착지 주소와 f_y 의 출발지 주소 중 큰 값을 의미한다. 수식의 결과로 0~1 사이의 값이 나오며, 0에 가까울수록 두 플로우가 유사하지 않으며, 1에 가까울수록 두 플로우가 유사함을 의미한다.

$$a_{Port}(f_x, f_y) = \left(\frac{prefixlenPort_{src}(f_x, f_y)}{16} \right)^2 + \left(\frac{prefixlenPort_{dst}(f_x, f_y)}{16} \right)^2$$

식 (5) 출발지, 도착지 포트 쌍 연결성 계산

Expr 5. Calculation of the connectivity of the source port and the destination port

식 (5)도 식 (4)의 계산식과 마찬가지로 출발지, 도착지 포트 쌍의 16비트의 앞부분이 얼마나 같은지를 계산하는 식이며, $prefixlenPort_{src}(f_x, f_y)$ 는 f_x 의 출발지 포트와 f_y 의 출발지 포트 16비트 중 앞부분이 같은 비트의 개수와 f_x 의 출발지 포트와 f_y 의 도착지 포트 중 큰 값을 의미하며, $prefixlenPort_{dst}(f_x, f_y)$ 는 f_x 의 도착지 포트와 f_y 의 도착지 포트 16비트 중 앞부분이 같은 비트의 개수와 f_x 의 도착지 포트와 f_y 의 출발지 포트 중 큰 값을 의미한다. 수식의 결과로 0~1 사이의 값이 나오며, 0에 가까울수록 두 플로우가 유사하지 않으며, 1에 가까울수록 두 플로우가 유사함을 의미한다. 이는 한 응용이 실행되었을 때 이에 의해서 발생하는 통신 플로우들의 포트는 증가하는 경향이 있으며 서로 유사한 범위의 포트를 사용한다는 점에 착안하였다.

$$a_{PROT}(f_x, f_y) = \begin{cases} 0: f_x.PROT \neq f_y.PROT \\ 1: f_x.PROT = f_y.PROT \end{cases}$$

식 (6) 프로토콜 연결성 계산

Expr 6. Calculation of the connectivity of the protocol

식 (6)은 프로토콜 연결성을 계산하는 식으로 프로토콜이 같으면 1 아니면 0으로 계산된다. 최종적으로 플로우간 연결성은 위 4가지 속성의 연결성에 가중치를 곱한 뒤 합하여 0~1 사이의 값으로 표현되고 두 플로우간 유사할수록 1에 가까운 값을 가진다. 연결성 모델에서는 유사성 모델에서 그룹화된 Group1과 타겟 트래픽의 연결성을 계산하여 특정 임계치 이상의 연결성을 가진 플로우들을 그룹화하고, 이 그룹을 다시 연결성 모델에 반복적으로 입력하여 탐지 대상의 트래픽들을 계속하여 그룹화할 수 있다.

3. 네트워크 플로우의 연관성 모델을 이용한 트래픽 분류 방법의 효과적인 적용 방안

유사성 모델은 통계 정보를 이용하고 연결성 모델은 헤더정보 및 플로우 발생 시간을 이용한다는 점에서 통계 시그니처 기반 분류 방법과 헤더 시그니처 기반 분류 방법과 비슷하다. 하지만, 제안하는 방법론은 유사성 모델을 먼저 1차적으로 사용하고 이에 분류하지 못한 나머지 트래픽을 연결성 모델을 연쇄적으로 사용함으로써 분석률을 최대화할 수 있으며 암호화된 트래픽 또한 분류가 가능하다. 연결성 모델에서는 헤더 시그니처처럼 특정 IP 리스트와 포트, 프로토콜이 같은 플로우만 분류하는 것이 아닌 연결성의 계산을 통하여 기준 플로우의 헤더정보와 유사한 플로우를 그룹화하며 1차적으로 유사성 모델의 출력을 입력으로 사용하기 때문에 포트 기반 분류 방법의 한계를 극복할 수 있다. 또한, 시그니처 기반 분류 방법이 아니므로 시그니처를 생성할 필요가 없으며 단지, Seed Group으로 사용할 플로우를 입력으로 설정하면 이를 탐지 대상 트래픽 분류의 첫 단서로 하여 대상 트래픽 그룹핑이 이루어진다.

본 방법론을 효과적으로 적용하기 위해서는 탐지 대상별로 유사성 모델의 가이드라인과 연결성 모델의 가이드라인을 구축하고 이를 참고하여 분류하여야 한다. 또한 탐지 대상별 가이드라인을 구축할 시 네트워크 관리자처럼 숙련된 전문가가 아닌 사용자도 쉽게 본 방법론을 통해 트래픽을 분류하는 것이 가능하다. 유사성 모델의 탐지 대상별 가이드라인이란 모든 정답지 플로우 중 어떤 플로우를 Seed Group으로 사용해도 정확도가 100%가 되는 정교한 임계치를 의미하며, 연결성 모델의 탐지 대상별 가이드라인이란 분석률을 최대화할 수 있도록 탐지 대상의 특성을 잘 반영한 연결성의 속성별 가중치와 임계치를 의미한다. 이와 같이 가이드라인이 필요한 이유는 다음과 같다.

먼저, 유사성 모델에서는 그룹화된 플로우들을 연결성 모델의 입력으로 사용하며 이를 연쇄적으로 사용하기 때문에 유사성 모델에서 False Positive가 조금이라도 있을 경우 연결성 모델에서는 반복되는 횟수에 따라 False Positive가 급증하게 된다. 따라서, 유사성 모델에서는 반드시 정확도가 100%가 되어야 한다. 때문에, 탐지 대상별로 어떤 플로우를 Seed로 사용하여도 False Positive가 0이 나오는 정교한 임계치가 설정되어야 한다. 두번째로 연결성 모델에서는 응용 또는 공격의 종류에 따라 발생하는 트래픽의 특성이 상이할 수 있음을 염두하여 각 대상별로 가중치를 조절할 수 있게함으로써 보다 더 정확도와 분석률을 향상시킬 수 있다.

4. 실험 및 결과

본 장에서는 본 방법론이 높은 정확도와 분석률의 방법론이며 이를 위한 연관성 모델의 가이드라인이 분명히 존재하고 이를 찾는 것이 가능하다는 것을 실험을 통해 증명한다.

실험 환경은 다음과 같다. 먼저 여러 호스트에서 탐지하고자 하는 대상의 정답지 트래픽과 노이즈 트래픽을 수집한다. 본 실험에서는 P2P 응용, 파일 공유 응용, SNS 응용의 대표로 각각 Torrent, Dropbox, Facebook을 선정하였다. 그리고 4개의 호스트에서 각 응용의 정답지 트래픽과 이에 대한 노이즈 트래픽을 수집하였다. 다음으로 정답지 플로우 중 하나를 Seed로 선택하고 유사성 모델의 임계치를 0~1까지 0.0001씩 증가시키며 분류를 한다. 결과에서 False Positive가 0이 되는 임계치를 확인한다. 모든 정답지 플로우를 Seed로 선정하여 위의 과정을 반복한 후 모든 Seed에 대하여 False Positive가 0이 되는 임계치의 최대값을 확인한다. 임계치가 높을수록 분석률은 감소하지만 정확도는 향상되기 때문에 이 임계치는 최대값이 된다. 이 유사성 모델의 임계치 중 최대값을 사용하여 유사성 모델로 1차적으로 분류를 한 후 연결성 모델에 입력으로 넣고 연결성 모델 속성의 가중치와 연결성 모델의 임계치를 0.1씩 조정하는 모든 경우의 수로 그룹화된 결과를 확인하고 분석률이 최대이며 정확도는 100%인 조합을 확인한다. 본 실험을 각 응용에 대하여 다른 정답지 트래이스들을 사용하여 4회 반복함으로써 각 응용의 강건한 유사성 및 연결성 가이드라인이 존재함을 확인한다.

표 1은 실험에 사용된 트래픽의 정보이다.

첫 행부터 4행씩 각각 Torrent, Dropbox, Facebook 응용의 4회분 실험에 사용된 트래픽 정보이다.

표 2는 실험 결과다. 1열은 실험번호. 2열은 모든 플로우를 한 번씩 SeedGroup으로 설정하여 유사성 모델에 적용하였을 시 False Positive가 0가 되는 임계치들 중 최대값이며, 3열~6열은 연결성 모델 임계치와 속성별 가중치를 나타내고 8열은 True Positive, 9열은 False Positive를 의미한다.

첫 행부터 4행씩 각각 Torrent, Dropbox, Facebook에 대한 실험 결과이다.

각 응용별로 4번의 실험 결과 Torrent의 유사성 모델의 정확도 100%의 임계치는 0.996임을 확인되었다. 이로써 Torrent에 대한 유사성 모델의 가이드라인은 0.996이라 할 수 있다. 4번의 실험 결과 정확도는 100%이며 분석률을 최대로 하는 연결성 모델 임계치와 속성별 가중치 조합이 모두 0.6, 0.2, 0.4, 0.1, 0.3로 일정하므로 보아 Torrent의 연결성 모델 가이드라인이 존재함을

표 1 실험에 사용된 트래픽의 정보

Table 1 Traffic information of the experiment

Experiment	Flow	Packet	Byte
1	Torrent	582	237205
	Noise	3189	85747
2	Torrent	924	173018
	Noise	2573	197871
3	Torrent	732	182642
	Noise	2669	257903
4	Torrent	579	230648
	Noise	2486	285159
1	Dropbox	39	190243
	Noise	3624	72800
2	Dropbox	45	196262
	Noise	4277	253646
3	Dropbox	44	320014
	Noise	4200	187214
4	Dropbox	39	288179
	Noise	4791	126271
1	Facebook	136	72026
	Noise	1554	57878
2	Facebook	132	83672
	Noise	776	27763
3	Facebook	158	61714
	Noise	2788	96882
4	Facebook	127	50786
	Noise	3017	190482

표 2 실험 결과

Table 2 Experiment result

Ex	Sim thres	Con thres	w ip	w port	w prot	w time	TP	FP
1	0.996	0.6	0.2	0.4	0.1	0.3	93.0%	0%
2	0.9759	0.6	0.2	0.4	0.1	0.3	99.2%	0%
3	0.9949	0.6	0.2	0.4	0.1	0.3	99.8%	0%
4	0.9870	0.6	0.2	0.4	0.1	0.3	95.4%	0%
1	0.9759	0.9	0.1	0.1	0.4	0.4	100%	0%
2	0.9696	0.9	0.1	0.1	0.4	0.4	84.2%	0%
3	0.973	0.9	0.1	0.1	0.4	0.4	100%	0%
4	0.958	0.9	0.1	0.1	0.4	0.4	100%	0%
1	0.9964	0.8	0.1	0.2	0.5	0.2	89.3%	0%
2	0.9913	0.8	0.1	0.2	0.5	0.2	83.2%	0%
3	0.9947	0.8	0.1	0.2	0.5	0.2	96.0%	0%
4	0.9971	0.8	0.1	0.2	0.5	0.2	98.0%	0%

확인하였다. Dropbox의 유사성 모델의 가이드라인은 0.9759, 연결성 모델의 가이드라인은 0.9, 0.1, 0.1, 0.4, 0.4로 확인하였으며, Facebook의 유사성 모델의 가이드라인은 0.9971, 연결성 모델의 가이드라인은 0.8, 0.1, 0.2, 0.5, 0.2로 확인하였다. 본 실험에서는 연결성 모델을 1회만 사용하였으며 Seed Group으로 임의의 플로우

하나만을 사용하였음에도 불구하고 100%의 정확도를 유지하면서 높은 분석률을 보였다. Seed Group으로 여러 개의 플로우를 사용하거나 연결성 모델을 반복적으로 사용할 경우 더욱 높은 분석률을 보일 것으로 예상된다.

본 실험에서는 분류하고자 하는 탐지 대상마다의 특정한 유사성 모델의 가이드라인과 연결성 모델의 가이드라인이 분명히 존재하고 이를 실험을 통해 찾을 수 있으며 거듭되는 실험으로 더욱더 강건하게 만들 수 있음을 증명하는 것에 의의를 둔다.

5. 결론 및 향후 연구

본 논문에서는 네트워크 플로우의 연관성 모델을 이용한 트래픽 분류 방법을 제안하고 이를 효과적으로 사용하기 위한 방안을 제시하였다. 그리고 본 방법론의 타당성을 실험을 통해 증명하였다. 향후 연구로는 다양한 응용에 대한 가이드라인을 구축하고 연쇄적인 연결성 기반을 통해 분석률을 100%로 만드는 연쇄적 가이드라인에 관한 연구와 본 방법론을 악성 탐지에 적용하는 연구를 진행할 계획이다.

References

- [1] S. H. Yoon, and M. S. Kim, "Research on Header Signature Maintenance Method for Internet Application Traffic Identification," *Proc KICS ICC 2011*, pp. 1200-1201, Jeju Island, Korea, Jun. 2011.
- [2] J. S. Park, S. H. Yoon, and M. S. Kim, "Performance Improvement of the Payload Signature based Traffic Classification System using Application Traffic Temporal Locality," *Proc APNOMS 2013*, pp. 1-6, Hiroshima, Japan, Sep. 2013.
- [3] H. M. An, J. H. Ham, and M. S. Kim, "Performance Improvement of the Statistical Information Based Traffic Identification System," *KTCCS*, Vol. 2, No. 8, pp. 335-342, Aug. 2013.
- [4] T. Nguyen, and G. A., "A survey of techniques for internet traffic classification using machine learning," *IEEE Communications Surveys & Tutorial*, Vol. 10, No. 4, pp. 56-76, Jan. 2009.
- [5] S. H. Lee, and M. S. Kim, "Application Traffic Classification using TensorFlow Machine Learning Tool," *Proc KICS 2016*, pp. 224-225, ChungAng Univ, Korea, Nov. 2009.



구 영 훈

2016년 고려대학교 컴퓨터정보학과 졸업(학사). 2016년~현재 고려대학교 컴퓨터정보학과 석사과정. 관심분야는 네트워크 관리 및 보안, 트래픽 모니터링 및 분석



심 규 석

2014년 고려대학교 컴퓨터정보학과 졸업(학사). 2016년 고려대학교 컴퓨터정보학과 졸업(석사). 2016년~현재 고려대학교 컴퓨터정보학과 박사과정. 관심분야는 네트워크 관리 및 보안, 트래픽 모니터링 및 분석



이 성 호

2016년 고려대학교 컴퓨터정보학과 졸업(학사). 2016년~현재 고려대학교 컴퓨터정보학과 석사과정. 관심분야는 네트워크 관리 및 보안, 트래픽 모니터링 및 분석, 트래픽 분류



Baraka D. Sija

2015년 세명대학교 정보통신학부 졸업(학사). 2016년~현재 고려대학교 컴퓨터정보학과 석사과정. 관심분야는 네트워크 관리 및 보안, 트래픽 모니터링 및 분석, 트래픽 분류



김 명 섭

1998년 포항공과대학교 전자계산학과 졸업(학사). 2000년 포항공과대학교 전자계산학과 졸업(석사). 2004년 포항공과대학교 전자계산학과 졸업(박사). 2006년 Dept. of ECS, Univ. of Toronto Canada. 2006년~현재 고려대학교 컴퓨터정보학과 교수. 관심분야는 네트워크 관리 및 보안, 트래픽 모니터링 및 분석, 멀티미디어 네트워크