

Mean FlowSpectrum을 활용한 암호화된 트래픽 분류 성능 개선 연구

김찬형¹, 이브라히모바-나일라, 김태윤², 김가영, 서주형, 배기태³, 윤종희⁴
 영남대학교 컴퓨터공학과 석사과정¹, 학부과정², 교수⁴
 영남대학교 로봇공학과 학부과정³

qnfrha@yu.ac.kr, 22446177@ynu.kr, elma9810@yu.ac.kr, im_770@naver.com,
 blane7777@naver.com, bgt010@naver.com, youn@yu.ac.kr

Improvement of Encrypted Traffic Classification Performance Using Mean Flow Spectrum

Chan-hyung Kim¹, Ibrahimova Naila, Tae-Yun Kim², Ga-Young Kim,
 Ju-hyeong Seo, Gi-Tae Bae³, Jong-Hee Youn⁴

Dept. of Computer Engineering, Yeung-Nam University^{1,2,4}
 Dept. of Robotics Engineering, Yeung-Nam University³

요약

암호화된 트래픽의 증가로 인해 기존 페이로드 기반 분석 기법의 활용이 어려워지면서, 복호화 없이도 효과적으로 트래픽을 분류할 수 있는 새로운 특징 추출 및 분류 기법이 요구되고 있다. 본 연구에서는 기존 FlowSpectrum 기법이 스펙트럼 라인 간의 간격만을 기반으로 유사도를 계산함에 따라 라인 간 겹침 시 분류 정확도가 저하되는 문제를 해결하고자, 평균값 비교 과정을 추가한 Mean FlowSpectrum 기법을 제안한다. 제안된 방법은 각 트래픽의 스펙트럼 평균값을 활용하여 전체적인 분포의 차이를 고려함으로써 분류의 정밀도를 높인다. ISCX-VPN2016 데이터셋을 기반으로 진행된 실험에서는, 기존 FlowSpectrum 대비 전반적인 정확도에서 약 5%p 향상된 성능을 보였으며, 특히 스트리밍, 파일 전송 등 일부 레이블에서 f1-score 기준 0.1 이상 향상된 결과를 확인할 수 있었다. 이는 평균값 기반의 보조 지표가 암호화된 트래픽 분류 성능 향상에 효과적으로 기여할 수 있음을 시사한다.

I. 서론

정보통신 기술의 발전과 함께 인터넷 사용이 급증하면서, 데이터의 보안과 프라이버시 보호에 대한 요구가 높아지고 있다. 이러한 요구에 부응하기 위해 많은 서비스 제공자들은 데이터 전송 시 암호화를 적용하고 있으며, 이는 사용자 정보를 보호하는 데 중요한 역할을 하고 있다. 그러나 암호화된 트래픽의 증가로 인해 네트워크 관리 및 보안 분야에서는 새로운 도전 과제가 발생하고 있다. 특히, 악성 트래픽의 탐지, 네트워크 성능 모니터링, QoS(Quality of Service) 보장 등의 측면에서 암호화된 트래픽을 효과적으로 분류하는 것이 필수적이다.

암호화된 트래픽은 일반적으로 패킷의 내용을 해석할 수 없기 때문에, 전통적인 페이로드 시그니처 기반 패킷 분석 기법으로는 그 특성을 파악하기 어렵다. 이로 인해 네트워크 관리

자와 보안 전문가들은 암호화된 트래픽의 유형을 식별하고, 이를 기반으로 적절한 대응 방안을 마련하는 데 어려움을 겪고 있다. 따라서, 암호화된 트래픽을 효과적으로 분류할 수 있는 새로운 기법의 개발이 절실히 요구된다.

FlowSpectrum은 암호화된 트래픽을 복호화하지 않고 분류하기 위한 특징(feature)이다. 본 연구에서는 FlowSpectrum의 성능을 개선한 Mean FlowSpectrum을 제안한다. Mean FlowSpectrum은 기존의 스펙트럼 라인 간의 간격에 대해 평균값 비교 과정을 추가하여 분류 성능을 높였다.

본 논문의 구성은 다음과 같다. 2장에서 기존의 FlowSpectrum에 대한 선행연구를 분석한다. 3장에서는 성능을 개선한 Mean FlowSpectrum을 제시하고, 4장에서 제안 방식과 종래방식의 결과를 비교 및 분석한다. 그리고 5장에서 결론을 맺는다.

II. 관련 연구

L Yang [1]등은 이상 트래픽을 탐지하기 위한 새로운 모델인 Semi-AE를 제안하였다. Semi-AE는 오토인코더를 통해 학습하는 과정 중 추출된 데이터를 활용하여 FlowSpectrum을 생성한다. [1]에서는 Semi-AE를 통해 NSL-KDD 데이터셋에서 FlowSpectrum을 생성하고, 이를 활용해 정상 트래픽과 비정상 트래픽으로 분류하는 실험을 진행하였으며 0.9513의 재현율을 보였다. 해당 연구는 처음으로 FlowSpectrum을 제안하였다는 의의가 있지만, 암호화되지 않은 트래픽을 사용하였다는 한계가 있다.

J Cui [2]등은 Semi-AE와 2차원 컨볼루션 신경망(2D-CNN)을 조합하여 암호화된 트래픽을 분류하는 모델인 Semi-2DCAE를 제안하였다. 해당 연구는 [1]과 동일하게 오토인코더를 통해 FlowSpectrum을 생성하고 트래픽을 분류하는 실험을 진행하였다. [2]는 암호화된 트래픽에서 FlowSpectrum을 생성하기 위해 2차원 컨볼루션 신경망을 사용하였으며, 암호화된 트래픽이 포함된 ISCX-VPN2016 데이터셋을 활용하여 트래픽을 분류 실험을 진행하였다. 해당 모델은 약 98%의 재현율을 보였으며 CNN을 활용하여 공간 구조적인 특징을 반영하였다는 의의가 있지만, 일부 데이터셋에서는 기존의 SemiAE 보다 낮은 성능을 보여주었다는 한계가 있다.

III. Mean FlowSpectrum

본 논문에서는 FlowSpectrum의 분류 성능을 높이기 위해 평균값 비교 과정을 추가한 Mean FlowSpectrum을 제안한다.

3.1 FlowSpectrum

FlowSpectrum이란 오토인코더의 학습과정에서 생성된 데이터를 1차원 표준 좌표계에 다양한 간격의 스펙트럼 선으로 표현한 새로운 특징(Feature)이다. 해당 특징의 경우 복호화를 하지 않아도 데이터의 추출이 가능하며, 스펙트럼

라인 간의 간격 비교를 통해 분류된다.

3.2 Mean FlowSpectrum

기존 FlowSpectrum의 경우 스펙트럼 라인이 겹쳐지는 경우 분류 성능이 저하되는 문제점이 있다. 이를 해결하기 위해 Mean FlowSpectrum은 기존의 스펙트럼 라인 간의 간격을 계산한 후, 각 레이블 별 FlowSpectrum의 평균값과의 비교 과정을 추가하여 분류 정확도를 향상시켰다.

IV. 실험 및 결과

본 연구에서는 ISCX-VPN2016 데이터셋을 사용하여 제안 기법의 성능을 평가하였다. 해당 데이터셋은 채팅, P2P, VoIP 등 총 6가지의 태입으로 구성되어 있으며, 실험을 위해 IP 주소와 포트 정보 등 바이어스 될 수 있는 정보들을 제거하였다.

실험과정은 다음과 같다. 먼저, 전처리된 데이터에서 Semi-2DCAE 모델을 사용하여 FlowSpectrum을 추출한 후, 이를 기존 분류 모델과 Mean FlowSpectrum 기반 분류 모델에 넣어 각각의 분류 성능을 비교하였다.

<표 1> 기존 FlowSpectrum 분류 실험 결과

Label	precision	recall	f1-score
Chat	0.707692	0.92	0.8
E-mail	0.895349	0.77	0.827957
File	0.606557	0.74	0.666667
P2P	0.882979	0.83	0.85567
Streaming	0.911392	0.72	0.804469
VoIP	0.898876	0.8	0.846561

<표 2> Mean FlowSpectrum 분류 실험 결과

Label	precision	recall	f1-score
Chat	0.851852	0.92	0.884615
E-mail	0.895349	0.77	0.827957
File	0.692982	0.79	0.738318
P2P	0.882979	0.83	0.85567
Streaming	0.924731	0.86	0.891192
VoIP	0.895238	0.94	0.917073

<표 1>과 <표 2>는 각각 기존 FlowSpectrum과 Mean FlowSpectrum을 활용하여 분류실험을 진행한 결과를 나타낸 표이다. 기존 FlowSpectrum 대비 전반적인 정확도에서 약 5%p 향상된 성능을 보였으며, 특히 스트리밍, 파일 전송 등 일부 레이블에서 f1-score 기준 0.1 이상 향상된 결과를 확인할 수 있다.

V. 결론

본 연구에서는 FlowSpectrum을 활용한 트래픽 분류의 성능을 향상시키기 위해 Mean FlowSpectrum을 제안하였다. 실험결과, 대부분의 레이블에서 성능이 향상된 것을 확인할 수 있었으며, 전체적으로 약 5%의 정확도가 향상되었다. 이는 평균값과 같은 보조지표가 분류 성능을 향상시키는데 효과적인 방법임을 시사한다. 향후 연구에서는 평균값 외에 분류 성능을 향상시키기 위한 추가적인 보조지표를 모색 할 것이다.

Acknowledgements

이 논문은 2025년도 정부(산업통상자원부)의 재원으로 한국산업기술진흥원의 지원을 받아 수행된 연구임(RS-2024-00406796, 2025년 산업혁신인재성장지원사업), 2025년도 정부(과학기술정보통신부)의 재원으로 정보통신기획평가원의 지원을 받아 수행된 연구임 (RS-2023-00235509, ICT융합 공공 서비스·인프라의 암호화 사이버위협에 대한 네트워크 행위기반 보안관제 기술 개발)

[참고문헌]

- [1] Yang, L., Fu, S., Zhang, X. et al. FlowSpectrum: a concrete characterization scheme of network traffic behavior for anomaly detection. World Wide Web 25, 2139 - 2161 (2022).
- [1] Cui J, Bai L, Li G, Lin Z, Zeng P. Semi-2DCAE: a semi-supervision 2D-CNN AutoEncoder model for feature representation and classification of encrypted traffic. PeerJ Comput Sci. 2023;9:e1635. Published 2023 Nov 9.2.