

Intrusion Detection System based on Packet Payload Analysis using Transformer

Woo-Seung Park*, Gun-Nam Kim*, Soo-Jin Lee*

*Graduate Student, Dept. of Defense Science, Korea National Defense University, Nonsan, Korea

*Graduate Student, Dept. of Defense Science, Korea National Defense University, Nonsan, Korea

*Professor, Dept. of Defense Science, Korea National Defense University, Nonsan, Korea

[Abstract]

Intrusion detection systems that learn metadata of network packets have been proposed recently. However these approaches require time to analyze packets to generate metadata for model learning, and time to pre-process metadata before learning. In addition, models that have learned specific metadata cannot detect intrusion by using original packets flowing into the network as they are. To address the problem, this paper propose a natural language processing-based intrusion detection system that detects intrusions by learning the packet payload as a single sentence without an additional conversion process. To verify the performance of our approach, we utilized the UNSW-NB15 and Transformer models. First, the PCAP files of the dataset were labeled, and then two Transformer (BERT, DistilBERT) models were trained directly in the form of sentences to analyze the detection performance. The experimental results showed that the binary classification accuracy was 99.03% and 99.05%, respectively, which is similar or superior to the detection performance of the techniques proposed in previous studies. Multi-class classification showed better performance with 86.63% and 86.36%, respectively.

▶ **Key words:** Natural Language Processing, IDS, Packet Payload, Transformer, UNSW-NB15

[요 약]

네트워크 패킷의 메타데이터를 학습한 침입탐지시스템이 최근 많이 제안되었다. 그러나 이러한 방식은 모델 학습에 사용할 메타데이터 생성을 위해 패킷을 분석하는 시간, 그리고 학습 전 메타데이터를 전처리하는 시간이 필요하다. 또한, 특정 메타데이터를 학습한 모델은 실제 네트워크로 유입되는 원본 패킷을 그대로 사용하여 침입을 탐지하는 것이 불가능하다. 이러한 문제를 해결하기 위해 본 논문에서는 패킷 페이로드를 하나의 문장으로 학습시켜 침입을 탐지하는 자연어 처리 기반의 침입탐지시스템을 제안하였다. 제안하는 기법의 성능 검증을 위해 UNSW-NB15와 Transformer 모델을 활용하였다. 먼저, 데이터셋의 PCAP 파일에 대한 라벨링을 실시한 후 2종의 Transformer 모델(BERT, DistilBERT)에 문장 형태로 직접 학습시켜 탐지성능을 분석하였다. 실험 결과 이진분류 정확도는 각각 99.03%, 99.05%로 기존 연구에서 제안한 기법들과 유사하거나 우수한 탐지성능을 보였으며, 다중분류는 각각 86.63%, 86.36%로 더 우수한 성능을 나타냄을 확인하였다.

▶ **주제어:** 자연어 처리, 침입탐지시스템, 패킷 페이로드, 트랜스포머, UNSW-NB15

-
- First Author: Woo-Seung Park, Co-Author: Gun-Nam Kim, Corresponding Author: Soo-Jin Lee
 - *Woo-Seung Park (akdongws@gmail.com), Dept. of Defense Science, Korea National Defense University
 - *Gun-Nam Kim (rjsska2918@gmail.com), Dept. of Defense Science, Korea National Defense University
 - *Soo-Jin Lee (cyberkma@gmail.com), Dept. of Defense Science, Korea National Defense University
 - Received: 2023. 10. 18, Revised: 2023. 11. 16, Accepted: 2023. 11. 17.

I. Introduction

점점 더 지능화되면서 진화하고 있는 네트워크 위협에 효율적으로 대응하기 위해 침입탐지에 인공지능(Artificial Intelligence) 기술을 접목하려는 시도가 증가하고 있다. 기계학습(Machine Learning), 딥러닝(Deep Learning) 및 앙상블 학습(Ensemble Learning) 등 침입탐지 분야에 도입되면서 탐지성능도 크게 향상되었다. 그러나 인공지능 기술을 기반으로 하는 침입탐지 모델은 대부분 패킷 자체가 아니라 패킷을 분석한 후 생성되는 정보인 메타데이터(metadata)를 활용하여 탐지모델을 학습시킨다. 이러한 접근방식은 공격 패킷을 분석하는데 많은 시간이 소요될 뿐만 아니라 모델이 학습했던 메타데이터의 속성과 동일하지 않은 속성을 가지는 침입은 탐지할 수 없다는 문제를 안고 있다. 또한, 메타데이터를 학습한 모델은 실제 네트워크 환경에서 발생하는 패킷을 학습했던 형태의 메타데이터로 변환해주어야만 침입탐지가 가능해진다.

이러한 문제를 해결하기 위해 본 연구진은 패킷의 페이로드를 그대로 활용하는 LightGBM 기반의 침입탐지시스템을 제안한 바 있다[1]. 선행연구에서는 10진수로 변환된 UNSW-NB15 데이터셋[2-3]의 패킷 페이로드를 이용, 1바이트당 하나의 특성으로 매칭해 총 1,500개의 특성을 생성하고 LightGBM 모델에 학습시켰다. 그 결과 이진분류 및 다중분류에서 각각 99.33%, 85.63%로 높은 정확도를 달성했으나, 데이터를 모델에 학습시키기 위해 16진수의 패킷 페이로드를 10진수로 변환하는 과정이 필요했다.

본 연구에서는 선행연구의 단점이라고 할 수 있는 패킷 페이로드의 10진수 변환과정을 제거하기 위해 패킷 페이로드를 원본 형태 그대로 학습할 수 있는 자연어 처리(Natural Language Processing) 모델을 적용하였다. 자연어 처리 모델은 텍스트 데이터 문장의 구조와 의미를 파악하여 효과적인 분류를 수행할 수 있어 최근 많은 연구에서 활용되고 있다. 이러한 장점을 활용하여 본 논문에서는 패킷 페이로드를 하나의 문장으로 직접 Transformer 모델에 학습시켜 침입을 탐지하는 방안을 제안하였다.

본 논문의 구성은 다음과 같다. 먼저, 2장에서 본 연구의 성능평가 실험에 사용된 UNSW-NB15 데이터셋에 관해 설명하고, 관련 연구를 정리한다. 3장에서는 제안하는 침입탐지 모델의 구축 절차와 실험 방법을 설명하고, 이어서 실험 결과를 분석한다. 마지막으로 4장에서 결론을 맺는다.

II. Preliminaries

1. UNSW-NB15 Dataset

인공지능 기반의 침입탐지시스템 연구에 사용되는 대표적인 데이터셋으로는 NSL-KDD[4], UNSW-NB15 등이 있다. 그 중 NSL-KDD 데이터셋은 최신의 공격 경향을 반영하지 못한다는 문제점도 있지만, 학습 및 테스트 데이터셋 분포가 데이터 유형별로 다르기 때문에 침입탐지 실험을 위한 신뢰성이 떨어진다. 반면, UNSW-NB15 데이터셋은 비교적 최신의 정상 및 비정상 활동을 포함하고, 침입탐지시스템 평가의 신뢰성을 보장할 수 있는 데이터셋으로 평가되어 다양한 연구에서 활발하게 활용되고 있다[2]. 이에 본 논문에서 제안하는 침입탐지 모델의 구축 및 성능평가를 위해 UNSW-NB15를 사용하였다.

UNSW-NB15 데이터셋은 'IXIA PerfectStorm'을 활용하여 네트워크 패킷을 수집한 후, 'Argus'와 'Bro-IDS Tool'을 이용하여 정상 및 비정상 트래픽을 분류한 공공 데이터셋으로서, PCAP 및 CSV 파일 형태로 제공되고 있다[3]. 본 연구는 자연어 처리를 통해 패킷 페이로드를 특성이 아닌 하나의 문장으로 인식하고 학습함으로써 침입을 탐지·분류하는데 목적이 있기 때문에 PCAP 형식의 데이터셋을 사용하였다.

UNSW-NB15 데이터셋은 175,341개의 학습(Train) 데이터와 82,332개의 테스트(Test) 데이터로 구성되어 있으며, 하나의 정상 데이터와 9개의 공격 데이터로 구성되어 있다. 공격 유형별 데이터의 수는 Table 1에서 보는 바와 같다.

Table 1. Configuration of UNSW-NB15 Dataset

Category	Train	Test
Analysis	2,000	677
Backdoor	1,746	583
DoS	12,264	4,089
Exploits	33,393	11,132
Fuzzing	18,184	6,062
Generic	40,000	18,871
Normal	56,000	37,000
Reconnaissance	10,491	3,496
Shellcode	1,133	378
Worms	130	44
Total	175,341	82,332

2. Related works

침입탐지 성능을 향상하기 위해 인공지능 기술을 접목하려는 연구는 매우 다양하게 진행되었다. 그러나 본 장에서는 동일한 데이터셋을 활용하여 실험을 진행한 선행 연구들과 자연어 처리를 기반으로 침입탐지를 시도했던 최근 연구들만을 간략하게 정리한다.

Kasongo 등[5]은 침입탐지 성능이 고차원 속성을 가지는 데이터에서 감소한다는 문제를 제기하였다. 그리고 XGBoost 모델에 속성 선택 방법을 적용하여, 42개의 UNSW-NB15 속성 중 최적의 영향을 미치는 19개의 속성만을 활용하는 방법으로 실험을 진행하였다. 그 결과 DT(decision tree)를 적용했을 때 이진분류 및 다중분류에서 각각 90.85%, 67.57%의 정확도를 달성하였다.

Vimalrosy 등[6]은 분류 정확도를 높이기 위해 최적의 속성만을 선택하는 OSS(Optimized Sine Swarm)를 수행한 후 RF(Random Forest) 모델을 사용하여 UNSW-NB15 데이터셋을 대상으로 실험을 진행하였고, 이진분류에서 98.15%의 높은 정확도를 달성하였다.

Kabir 등[7]은 이진분류에서 SVM(Support Vector Machine) 모델을 사용하여 정상과 비정상 트래픽에 대한 이진분류를 수행했을 때 가장 높은 정확도(82.11%)를 달성하였다. 또한, 다중분류에서는 SVM으로 1차 분류된 결과를 DT(C5.0), Naive Bayes 등에 적용하여 비교 분석하였다. 그 결과 Decision Trees에 적용하였을 때 86%의 가장 높은 정확도를 달성하였다.

Slay 등[8]은 KDD 99 및 NSL-KDD 데이터셋의 최신 트래픽 패턴 부족, 학습 및 테스트 세트의 분포 차이 등의 문제를 제기하고, UNSW-NB15 데이터셋을 사용하여 정확도와 FAR(오탐지율)을 비교하였다. 실험 결과 LR(Logistic Regression) 모델을 사용했을 때 83.15%의 다중분류 정확도가 달성됨을 확인하였다.

Das 등[9]은 UNSW-NB15 데이터셋 불균형 문제를 해결하기 위해 Elastic Net 및 SFS(Sequential Feature Selection) 알고리즘을 적용하여 전처리 과정을 수행하였다. Balanced Bagged, XGBoost 및 RF-HDDT 세 가지 모델을 앙상블 하여 실험을 진행한 결과, 이진분류에서 97.80%의 정확도를 보였다.

Jing 등[10]은 인공지능 기반 침입탐지 연구의 대부분이 KDD CUP99 데이터셋을 사용한다는 문제를 제기하면서, UNSW-NB15 데이터셋을 대상으로 SVM을 적용하여 실험을 진행하였다. 실험 결과 이진분류 85.99%, 다중분류 75.77%의 정확도를 달성하였다.

Meghdouri 등[11]은 고차원의 데이터를 저차원의 데이터로 환원하기 위해 PCA(Principal Component Analysis)를 사용하여 데이터를 전처리한 후 실험을 진행하였다. 적용 모델 중 RF(Random Forest)가 Precision 84.9%, Recall 85.1%로 가장 높은 탐지성능을 보였다.

한편, 본 연구에서 제안하는 접근방법과 유사하게 패킷 페이로드를 문장으로 학습하는 자연어 처리 기반의 침입탐지 연구도 진행된 바 있다.

Zihan 등[12]은 Transformer 모델의 'Self-attention' 기능을 접목하여 RTIDS(Robust Transformer-based Intrusion Detection System)를 제안하였다. 성능평가는 CICIDS2017 및 CIC-DDoS 2019 데이터셋을 대상으로 진행하였으며, SVM 및 RNN(Recurrent Neural Network)과 비교한 결과, RTIDS가 98% 이상의 가장 높은 이진분류 정확도를 달성하였다.

Yang 등[13]은 Transformer를 Vision Task에 적용한 모델인 ViT(Vision Transformer)를 제안하였다. NSL-KDD 데이터셋을 사용하여 실험을 진행한 결과 이진분류에서 99.68%의 높은 정확도를 달성하였다.

Lim 등[14]은 네트워크 이상행위 탐지를 위해 BERT(Bi-directional Encoder Representations from Transformers), LSTM(Long Short Term Memory) 및 GRU(Gate Recurrent Unit)를 비교 분석하였다. CSIC 2010 데이터셋을 대상으로 성능을 측정한 결과, BERT의 이진분류 정확도가 98% 이상으로 가장 높았다.

III. The Proposed Scheme

1. Experimental Preparation

1.1 Dataset Labeling

UNSW-NB15 데이터셋에서 제공하는 원본 PCAP은 라벨링이 되어 있지 않아 패킷 페이로드를 직접 모델에 학습시키는 등 활용하기에는 제한이 있었다. 라벨링 문제를 해결하기 위해 본 연구진이 수행했던 선행연구[1]에서와 같이 'Payload-Byte' 기법[15]을 적용하여 라벨링 작업을 먼저 수행하였다.

Fig. 1에서 보는 바와 같이 Payload-Byte 기법은 데이터셋 내에 PCAP 파일과 메타데이터가 함께 존재하는 경우, 메타데이터의 IP 주소나 Port 번호 등과 같은 값을 이용하여 Raw PCAP 파일 내에서 일치하는 패킷을 찾아 메타데이터와 동일한 클래스로 분류하고 라벨링 작업을 실시한다.

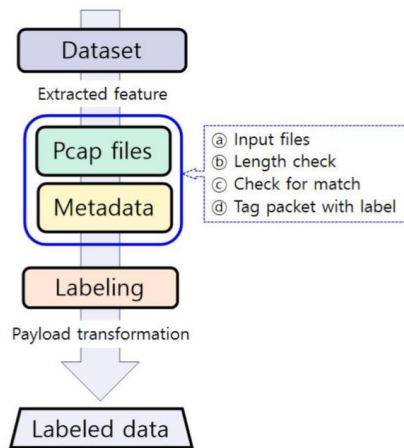


Fig. 1. Concept of Payload-Byte

라벨링 작업이 완료된 이후에는, 메타데이터와 일치하지 않아 라벨이 부여되지 않은 패킷 또는 동일 패킷에 동일 라벨이 부여된 중복 인스턴스(instances)를 1차로 제거한 결과 총 79,128개의 패킷이 식별되었다. 그러나 Normal 클래스에 포함된 데이터가 차지하는 비중이 과도하게 높아 정상 클래스에 대해서만 메타데이터에서 Normal 클래스가 가지는 비중을 고려하여 언더샘플링(Under-sampling)을 수행하였다. 언더샘플링 과정에서는 Normal 클래스에 포함된 모든 데이터가 전반적으로 실험에 활용될 수 있도록 매 실험이 진행될 때마다 다른 패킷 데이터가 선택되게 설정하였다. 축소 단계를 모두 거친 후 실제 실험에 사용한 데이터세트는 Table 2에서 보는 바와 같다.

Table 2. Configuration of Experimental Dataset

Category	Train	Test
Analysis	538	134
Backdoor	620	155
DoS	2,210	552
Exploits	9,794	2,449
Fuzzers	8,476	2,119
Generic	9,457	2,364
Normal	9,849	2,463
Reconnaissance	5,644	1,411
Shellcode	680	170
Worms	66	17
Total	47,334	11,834

1.2 Transformer

기존 모델 Seq2seq(Sequence-to-sequence) 모델은 인코더가 입력 시퀀스를 하나의 벡터로 압축하는 과정에서 입력 시퀀스의 정보가 일부 손실된다는 단점을 가진다. 이러한 단점을 개선하기 위해 2017년 트랜스포머

(Transformer) 개념이 등장하였는데, 트랜스포머는 RNN이나 CNN(Convolutional Neural Network)을 사용하지 않고, 위치 인코딩(positional encoding)을 활용하여 기존 장기 기억 의존 및 계산 속도 문제를 개선하였다[16].

대표적인 트랜스포머 모델에는 BERT(Bidirectional Encoder Representations from Transformers)와 BERT 기반의 파생 모델들(DistilBERT, RoBERTa, XLNet 등)이 있다. BERT는 사전 학습된 대용량의 레이블링 되지 않은(unlabeled) 데이터를 이용하여 언어모델을 학습하고 이를 토대로 특정 작업(문서 분류, 질의응답, 번역 등)을 위한 신경망을 추가하는 전이 학습 방법이다. 기존 방식은 문장에서 단어를 순차적으로 입력받고 그 다음 단어를 예측하는 단방향(uni-directional)이었지만, BERT는 문장 전체를 입력받고 단어를 예측하는 양방향(bi-directional) 학습이 가능해 기존 언어모델에 비해 우수한 성능을 낼 수 있다[17]. 그러나 BERT 모델의 큰 크기는 실시간 처리와 같은 응용을 어렵게 만들기 때문에 이를 더 가볍고, 빠르게 만든 DistilBERT(Distilled version of BERT) 모델이 개발되었다. DistilBERT는 인코더 계층의 개수가 6개로 BERT에 비해 반으로 줄어든 모델로 기존 BERT 모델보다 40% 작은 크기를 가지면서도 성능은 BERT 대비 97%에 육박하며, 처리 속도는 60% 향상되었다[18].

반면, RoBERTa, XLNet 등과 같은 파생 모델들은 성능 향상을 위해 학습 데이터와 Batch size를 증가시키면서 모델의 학습 시간도 크게 증가하였다[19-21]. 이에 본 실험에서는 트랜스포머 모델 중 비교적 가벼운 BERT와 DistilBERT 모델만을 적용하여 실험을 진행하였다.

2. Experimental Design

본 연구에 사용한 UNSW-NB15 데이터세트의 PCAP 파일은 라벨링이 되어 있지 않은 상태이기 때문에 학습에 그대로 활용하기에는 제한이 있다. 따라서 앞서 언급한 바와 같이 라벨링을 수행하였으며, 패킷 페이로드 그 자체를 활용해 모델을 학습시키는데 목적이 있기 때문에 추가적인 변환과정 없이 16진수의 형태를 가진다. 한편, 길이는 각 패킷 페이로드에 따라 다르기는 하지만 3,000 ~ 4,500 글자의 길이를 가지게 된다.

이상과 같은 과정을 거쳐 문장으로 변환된 데이터를 BERT와 DistilBERT 분류 모델에 직접 학습시킨 후 성능 평가를 수행하였다. 데이터세트는 학습과 테스트 데이터셋을 8대2의 크기로 분할하였으며, 세부적인 실험과정은 Fig. 2에서 보는 바와 같다. 실험은 구글에서 제공하는 Colab pro 환경(python 3.10, T4 GPU, 51GB RAM)에서 진행되었다.

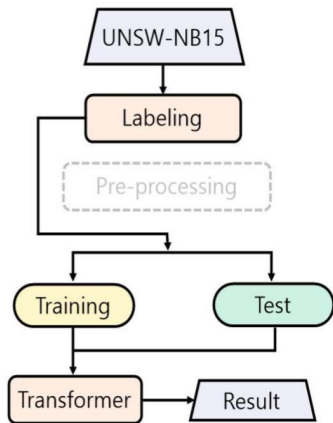


Fig. 2. Experimental design

3. Results and Analysis

라벨링이 완료된 데이터셋을 딥러닝 모델인 BERT와 DistilBERT에 적용하여 30회 반복 실험하였으며, 실험 간 적용된 하이퍼파라미터(hyperparameter)는 Table 3에서 보는 바와 같다.

Table 3. Hyperparameter setting

Batch size	Learning rate	Epoch
32	2e-5	3

실험 결과는 Table 4에서 보는 바와 같이 이진분류에서 정확도 99.03%, 99.05%, F1-score 99.39%, 99.40%, Precision 99.04%, 99.11%, Recall 99.74%, 99.70%를 달성하여 동일 데이터셋을 활용하여 이진분류를 시도했던 선행연구 대비 유사하거나 우수한 탐지성능을 보였다.

Table 4. Performance Comparison of Binary Classification

Category	Proposed		[1]	[8]	[9]	[12]	[13]
	BERT	Distill BERT					
Accuracy	99.03	99.05	99.33	90.85	98.15	97.80	85.99
F1-score	99.39	99.40	98.73	88.45	-	-	-
Precision	99.04	99.11	98.16	80.33	-	97.81	-
Recall	99.74	99.70	99.31	98.38	-	-	-

다중분류의 경우 F1-score를 측정하는 방법은 다양할 수 있으나, 본 연구에서는 데이터의 불균형이 심한 경우에도 이를 반영할 수 있는 Micro F1-score를 활용하였다.

다중분류 탐지성능은 Table 5에서 보는 바와 같으며, 각각 86.63%, 86.36%의 정확도를 달성하였다. 특히, 동일하게 패킷 페이로드를 활용한 선행연구[1] 대비 다중분류에서 1.0%p 높은 정확도를 보여 본 연구에서 제안한 기법

이 더 우수한 성능을 달성할 수 있음을 확인하였다.

Table 5. Performance Comparison of Multi-class Classification

Category	Proposed		[1]	[10]	[11]	[13]	[14]
	BERT	Distill BERT					
Accuracy	86.63	86.36	85.63	86.00	83.15	75.77	-
F1-score	86.63	86.36	85.68	86.00	-	-	84.9
Precision	86.63	86.36	87.88	-	-	-	85.1
Recall	86.63	86.36	85.61	-	-	-	84.9

제안하는 접근방법에 대한 성능분석 결과, 이진분류와 다중분류 성능이 선행연구에 비해 월등하게 향상되지는 않았다. 그러나 제안하는 접근방법은 과도한 시간이 소요되는 전처리 과정을 거치지 않고 패킷 페이로드를 그대로 학습하기 때문에 탐지모델을 생성하는 시간을 획기적으로 감소시킬 수 있다. 또한, 메타데이터를 학습한 탐지모델은 실제 네트워크 환경에서 유입되는 패킷을 메타데이터로 변환해 주어야만 침입탐지가 가능하지만, 제안하는 접근방법은 패킷 페이로드를 원본 그대로 활용하기 때문에 침입탐지의 실시간성 측면에서도 큰 효과를 발휘할 수 있다.

이진분류와 다중분류의 혼동행렬은 Fig. 3 및 Fig. 4에서 보는 바와 같다.

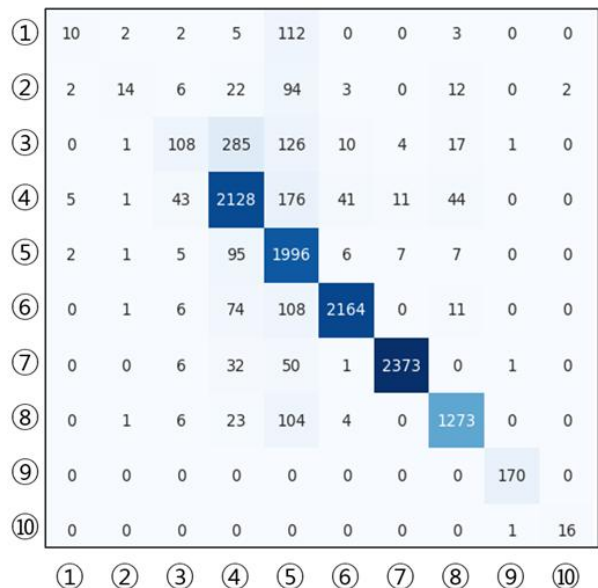


Fig. 3. Confusion Matrix of Multi-class Classification(BERT)

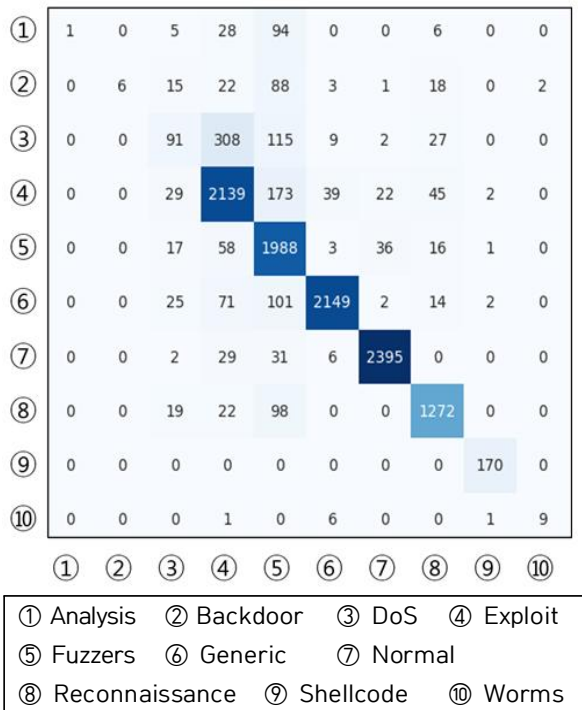


Fig. 4. Confusion Matrix of Multi-class Classification(DistilBERT)

IV. Conclusions

인공지능 기술을 접목한 침입탐지시스템은 모델 학습을 위해 장시간에 걸쳐 수집된 패킷을 분석하고 메타데이터를 생성한 후 전처리 등을 수행해야 한다. 또한, 특정 메타데이터를 학습한 모델은 다른 특성을 가지는 패킷을 탐지하지 못한다는 한계를 가진다. 이러한 문제점을 해결하기 위해 본 논문에서는 네트워크로 유입되는 패킷에서 페이로드를 추가적인 변환과정 없이 하나의 문장으로 학습시키는 자연어 처리 모델 기반의 침입탐지시스템을 제안하였다.

제안하는 접근방법의 성능평가를 위해 UNSW-NB15 데이터셋 내에 포함되어 있는 PCAP 파일을 사용하였으며, Payload-Byte 기법을 적용하여 라벨링을 수행하였다. 그리고 추가적인 변환과정 없이 패킷 페이로드를 하나의 문장으로 BERT와 DistilBERT 분류 모델에 적용하여 학습을 수행한 후 성능을 평가하였다.

실험 결과, 이진분류 정확도는 각각 99.03%, 99.05%를 달성하여 선행연구 대비 유사하거나 탐지성능이 미세하게 향상됨을 확인하였다. 한편, 다중분류에서는 각각 86.63%, 86.36%로 본 연구에서 제안한 기법이 선행연구 대비 더 우수한 성능을 나타냄을 확인하였다.

그리고 본 논문에서 제안한 접근방법은 자연어 처리를 통해 패킷 페이로드 자체를 하나의 문장으로 학습시키기

때문에 전처리 과정을 생략할 수 있어 탐지모델 생성 시간을 크게 단축시킬 수 있으며, 새로운 공격이 발생할 경우 모델 업데이트도 용이해진다.

또한, 기존의 메타데이터를 활용한 연구들은 형식이 실험마다 다르게 정의되어 비교에 적합하지 않지만, 실제 전송되는 데이터인 패킷 페이로드는 모든 실험의 데이터에서 같은 형식으로 나타나기 때문에 페이로드 자체를 실험에 활용함으로써 데이터의 통일성을 유지하면 향후에는 실험 간 표준화되고 비교 가능한 기준을 제시할 수 있을 것으로 기대된다.

본 연구에서는 하나의 문자열로 데이터 분석이 가능하고, 높은 탐지성능을 가진다는 점을 고려하여 Transformer 모델 중 비교적 가벼운 BERT와 DistilBERT 모델만을 우선 적용하여 실험을 진행하였다. 향후 연구에서는 보다 다양한 분류 모델 및 데이터셋을 적용하여 추가적인 비교 분석과 실험을 통해 최적의 탐지성능을 가지는 침입탐지시스템에 대한 연구를 수행할 계획이다.

REFERENCES

- [1] Gun-Nam Kim, Han-Seok Kim and Soo-Jin Lee, "Intrusion Detection System based on Packet Payload Analysis using LightGBM" Journal of The Korea Society of Computer and Information, vol. 28(6), pp 47-54, Jun 2023 DOI: 10.9708/jksoci.2023.28.06.047
- [2] N. Moustafa and J. Slay, "UNSW-NB15: a comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set)," 2015 Military Communications and Information Systems Conference (MilCIS), pp. 1-6, Nov 2015. DOI: 10.1109/milcis.2015.7348942.
- [3] Canadian Institute for Cybersecurity, "UNSW-NB15 Data set" <https://www.unsw.adfa.edu.au/unsw-canberra-cyber/cybersecurity/ADFA-NB15-Datasets/>
- [4] Australian Center for Cyber Security, "NSL-KDD dataset" <https://www.unb.ca/cic/datasets/nsl.html>
- [5] Sydney M. Kasongo and Yanxia Sun, "Performance Analysis of Intrusion Detection Systems Using a Feature Selection Method on the UNSW-NB15 Dataset", Kasongo and Sun J Big Data (2020) 7:105. DOI: 10.1186/s40537-020-00379-6
- [6] J. Vimalrosy and S. Brittorameshkumar, "OSS-RF: INTRUSION DETECTION USING OPTIMIZED SINE SWARM BASED RANDOM FOREST CLASSIFIER ON UNSW-NB15 DATASET", IJTPE Journal, vol. 14, pp 275-283, June 2022.
- [7] M. H. Kabir, M. S. Rajib, A. S. M. T. Rahman, Md. M. Rahman, and S. K. Dey, "Network Intrusion Detection Using UNSW-NB15

- Dataset: Stacking Machine Learning Based Approach,” 2022 International Conference on Advancement in Electrical and Electronic Engineering (ICAEEE), pp. 1-6, Feb 2022. DOI: 10.1109/icaeee54957.2022.9836404.
- [8] N. Moustafa. and J. Slay, “The evaluation of Network Anomaly Detection Systems: Statistical analysis of the UNSW-NB15 data set and the comparison with the KDD99 data set. Information Security Journal: A Global Perspective, vol. 25, no. 1-3, pp. 18-31, 2016. DOI: 10.1080/19393555.2015.1125974.
- [9] A. Das, Pramod and Sunitha B S “Anomaly-based Network Intrusion Detection using Ensemble Machine Learning Approach”, International Journal of Advanced Computer Science and Applications(IJACSA), vol 13(2), pp. 635-645, 2022. DOI: 10.14569/IJACSA.2022.0130275
- [10] D. Jing and H.-B. Chen, “SVM Based Network Intrusion Detection for the UNSW-NB15 Dataset,” 2019 IEEE 13th International Conference on ASIC (ASICON), pp. 1-4, Oct 2019. DOI: 10.1109/asicon47005.2019.8983598.
- [11] F. Meghdouri, T. Zseby, and F. Iglesias, “Analysis of Lightweight Feature Vectors for Attack Detection in Network Traffic,” Applied Sciences, vol. 8, no. 11, pp 2196, Nov 2018. DOI: 10.3390/app8112196.
- [12] W. Zihan, Z. Hong, W. Penghai and S. Zhibo, “RTIDS: A Robust Transformer-Based Approach for Intrusion Detection System”, IEEE Access, vol. 10, pp. 64375-64387 June 2022. DOI: 10.1109/ACCESS.2022.3182333
- [13] Y. G. Yang, H. M. Fu, S. Gao, Y. H Zhou and W. M. Shi, “Intrusion detection: A model based on the improved vision transformer”, National Natural Science Foundation of China, April 2022. DOI: 10.1002/ett.4522
- [14] Ju-wan Lim, In-kyung Kim, Myung-hak Lee, Jung-min Ha and Jae-koo Lee, “Performance comparison of network anomaly detection using BERT, LSTM and GRU”, KICS, pp. 1268-1269, Feb 2022
- [15] Y. A. Farrukh, I. Khan, S. Wali, D. Bierbrauer, and N. Bastian, “Payload-Byte: A Tool for Extracting and Labeling Packet Capture Files of Modern Network Intrusion Detection Datasets,” pp 58-67, Sep. 2022. DOI: 10.36227/techrxiv.20714221.v1.
- [16] Vaswani, Ashish, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser and Illia Polosukhin. “Attention is All you Need.” NIPS, 2017.
- [17] J. Devlin, Ming-Wei Chang, Kenton Lee, and K. Toutanova., “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”, 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, vol 1, pp 4171-4186, May 2019. DOI: 10.48550/arXiv.1810.04805.
- [18] V. sanh, L. debut, J. chaumond and T. wolf, “DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter”, Mar 2020. DOI: 10.48550/arXiv.1910.01108
- [19] Yinhan Liu, Myle Ott, Naman Goyal and Jingfei Du, “RoBERTa: A Robustly Optimized BERT Pretraining Approach”, Computation and Language, Jul 2019. DOI: 10.48550/arXiv.1907.11692.
- [20] Zhenzhong Lan, Mingda Chen and Sebastian Goodman, “ALBERT: A Lite BERT for Self-supervised Learning of Language Representations”, Computation and Language, vol. 6, Feb 2020. DOI: 10.48550/arXiv.1909.11942.
- [21] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, R.Salakhutdinov and Quoc V. Le “XLNet: Generalized Autoregressive Pretraining for Language Understanding”, NeurIPS, Jun 2019. DOI: 10.48550/arXiv.1906.08237.

Authors



Woo-Seung Park received B.S. degree in 2010 from the Department of e-business Studies, Sung-kyul University. He is currently a graduate student in the Department of Defense Science, Korea National Defense University.

His research interests include Deep Learning and Intrusion Detection System.



Gun-Nam Kim received B.S. degree in 2016 from the Department of Chinese Studies, Chang-Won National University. He is currently a graduate student in the Department of Defense Science, Korea

National Defense University. His research interests include Machine Learning, Intrusion Detection System and Cyber security Strategy.



Soo-Jin Lee received B.S., M.S. and Ph.D. degrees in Computer Science from Korea Military Academy, Yonsei University and Korea Advanced Institute of Science and Technology(KAIST) in 1992, 1996 and 2006.

He is currently a professor of the Department of Defense Science, Korea National Defense University from 2006. His research interests include National Cybersecurity Policy, Intrusion Detection System, Mobile Network Security, Machine Learning, Encryption theory and applications.