

인공신경망을 통한 KDD CUP 99 와 NSL-KDD 데이터 셋 비교

지현정*, 김용현**, 김동화**, 신동규*, 신동일*

*세종대학교 컴퓨터공학과

**국방과학연구소

e-mail : aasdddfd111@gmail.com

A Study on comparison of KDD CUP 99 and NSL-KDD using artificial neural network

Ji Hyunjung*, Kim Yonghyun **, Kim Donghwa**, Shin Dongkyoo*, Shin Dongil*

*Dept. of Computer Engineering, Sejong, University

**Agency for Defense Development

요 약

최근 컴퓨터 네트워크를 활용하는 다양한 기기들이 개발되고 급격히 확산되면서, 컴퓨터 네트워크는 전보다 많은 보안문제에 직면하게 되었다. 이에 따라 네트워크 보안을 위한 침입탐지시스템의 필요성이 대두된다. 침입탐지시스템을 구현하기 위한 대표적인 데이터 셋으로는 KDD CUP 99(KDD'99)와 이후 KDD'99의 문제점을 보완하여 공개된 NSL-KDD가 있다. 본 논문에서는 KDD'99와 NSL-KDD를 소개하고 인공신경망을 통해 두 데이터 셋을 비교 분석하였다. Multi-Layer Perceptron을 사용해 데이터 셋을 분석해본 결과, KDD'99는 전체 정확도에서 더 높은 결과를 얻은 반면 공격 별 탐지 정확도 면에서는 NSL-KDD에 뒤쳐졌다.

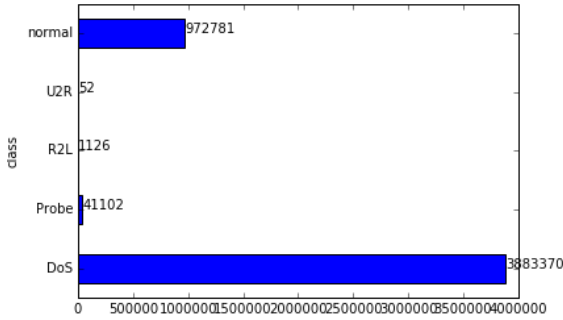
1. 서론

가구 인터넷 보급률이 84.4%(대한민국, 2015년 기준)에 이를 정도로 인터넷 접속은 생활의 필수불가결한 부분이 되었다. 컴퓨터 네트워크의 확산은 우리 삶에 편의를 가져다 주었지만, 동시에 네트워크 보안이라는 문제를 야기했다. 네트워크 보안 문제는 각종 프로토콜이나 네트워크를 통해 불법적인 정보유출, 서비스 방해 및 시스템을 취약하게 행위를 말한다 [1]. 이러한 시스템 보안과 침입 탐지를 위해 침입탐지시스템(Intrusion Detection System, IDS)의 필요성이 대두되고 있다. 침입탐지시스템은 컴퓨터 시스템과 네트워크 혹은 정보 시스템에 대한 공격을 탐지하는 것을 목표로 한다 [2]. 현재 침입탐지시스템을 설계하기 위해 많은 기계학습 알고리즘이 제안되고 있다 [3][4]. 따라서 이 알고리즘을 학습시키기 위한 데이터 셋의 필요성이 증대되었다. 데이터 셋은 알고리즘의 학습뿐만 아니라 테스트 과정에서도 영향을 미치므로, 실용적인 시스템을 구현하기 위한 양질의 데이터 셋이 요구된다. 침입탐지시스템을 위해 공개된 대표적인 데이터 셋으로는 KDD CUP 99와 NSL-KDD가 있다. 본 논문에서는 두 데이터 셋을 소개하고 인공신경망을 통해 두 데이터 셋을 비교해 본다.

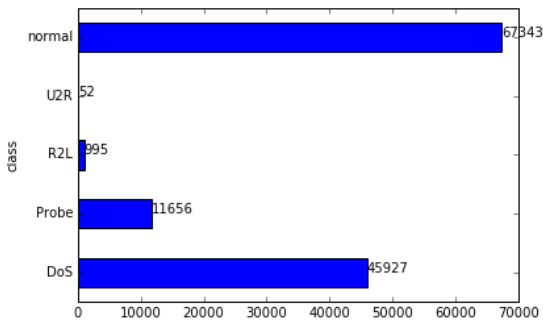
2. 데이터 셋 소개

KDD CUP 99(KDD'99)는 DARPA'98 침입탐지시스템 평가 프로그램에서 수집된 데이터를 기반으로 만들어진 데이터 셋이다. 지속기간, 프로토콜 종류 등 41개의 속성과, 각 레코드가 어떤 공격인지를 나타내는 클래스까지 더하여 총 42개의 속성으로 이루어져 있다. 각 공격은 DoS, Probe, R2L, U2R 그리고 정상상태인 normal로 분류된다 [5].

NSL-KDD는 2009년 Tavallaee가 KDD'99의 본질적인 문제점을 비판하며 제안한 데이터 셋이다 [6]. KDD'99는 많은 중복된 레코드를 포함하고 있기에 KDD'99로 알고리즘을 학습시킬 시 빈도가 높은 공격에 치우쳐져 학습될 가능성이 있으며, 테스트 과정의 평가 결과에도 영향을 미칠 수 있다 [7]. 기존의 KDD'99는 자주 나타나지 않지만 실제 네트워크에 위험한 U2R이나 R2L과 같은 공격의 분류에 낮은 성능을 보였다. NSL-KDD는 KDD'99의 중복된 레코드를 삭제하고 크기를 줄여 만들어졌다. (그림 1)과 (그림 2)는 KDD'99와 NSL-KDD의 트레이닝 셋의 공격 분포를 보여준다. 각각의 테스트 셋에서도 비슷한 분포를 보인다.



(그림 1) KDD'99의 트레이닝 셋 공격분포



(그림 2) NSL-KDD의 트레이닝 셋 공격분포

3. 사용된 알고리즘 및 실험방법

본 논문에서는 데이터 셋을 비교하기 위한 알고리즘으로 Multi-Layer Perceptron(MLP)을 사용하였다. MLP는 입력층과 출력층 사이에 1개 이상의 은닉층을 포함하는 구조이다. 본 논문에서는 41개의 속성을 입력으로 갖는 입력층과 21개의 뉴런을 갖는 2개의 은닉층, 그리고 클래스 분류에 따라 5개 혹은 2개의 출력을 갖는 출력층으로 MLP를 구성하였다. 알고리즘은 theano 라이브러리를 사용하여 python으로 구현하였다.

A. 데이터 로딩 및 전처리

KDD'99에서는 kddcup.data.gz와 corrected.gz를 NSL-KDD에서는 KDDTrain+.txt와 KDDTest+.txt를 각각 트레이닝과 테스트를 위해 사용하였으며, 트레이닝 셋에서 25%를 검증용을 위해 사용하였다. 총 42개의 속성 중 protocol_type, service_type, src_bytes_type과 class 속성을 수치화 하였으며, class 속성의 경우 공격을 DoS, Probe, R2L, U2R, normal 중 알맞은 종류로 대치하였다.

B. 알고리즘 학습

사용한 MLP의 활성화 함수로는 hyperbolic tangent를, 비용함수로는 negative log likelihood를 사용하였다. 각 트레이닝 셋의 class 속성을 제외한 41개의 속성을 갖는 데이터를 넣고 class를 제대로 분류할 수 있도록 backpropagation 알고리즘을 사용해 가중치와 바이어스 값을 수정한다.

C. 알고리즘 테스트

알고리즘을 훈련시킨 후 테스트 셋을 통해 얻은 결과 값과 실제 값을 비교해 같은 경우의 백분율을 구한다.

알고리즘은 총 10번 실험하였으며, 평균값을 기준으로 분석하였다.

4. 실험결과

KDD'99의 중복된 레코드가 알고리즘 학습과 테스트 과정에서 영향을 미치는지 알아보기 위해 기존의 데이터 셋 실험 뿐만 아니라 KDD'99와 NSL-KDD의 트레이닝 셋과 테스트 셋을 교차하여 실험하였다. <표 1>은 실험결과를 정리한 것으로 총 10번의 실험 결과의 평균을 보여준다. KDD'99를 통해 학습시킨 알고리즘은 NSL-KDD로 테스트 한 경우 정확도가 낮아지는 것을 볼 수 있다. KDD'99의 트레이닝 셋으로 학습시킨 알고리즘은 KDD'99와 NSL-KDD의 테스트 셋으로 실험한 결과 다섯 개의 클래스 분류 정확도를 기준으로 각각 평균 80.23%, 63.03%이었다. 반면에 NSL-KDD로 학습시킨 알고리즘에 KDD'99를 넣어 실험한 결과 정확도가 높아진 것을 관찰할 수 있다. NSL-KDD의 트레이닝 셋으로 학습시킨 알고리즘은 두 테스트 셋으로 실험한 결과 각각 75.25%, 67.45%이었다. 즉, KDD'99와 NSL-KDD의 트레이닝 셋으로 학습시킨 알고리즘이 두 테스트 셋에 대해 비슷한 성능을 보인다는 것을 알 수 있다.

<표 1> 실험결과

Training Set	Test Set	Type	accuracy(%)
KDD'99	KDD'99	five class	80.23
		binary class	82.41
NSL-KDD	NSL-KDD	five class	67.45
		binary class	73.8
KDD'99	NSL-KDD	five class	63.03
		binary class	68.73
NSL-KDD	KDD'99	five class	75.25
		binary class	86.2

(그림 3)과 (그림 4)는 각각 KDD'99와 NSL-KDD를 알고리즘을 통한 10번의 실험 중 한 실험의 결과를 나타낸 confusion matrix이며, <표 3>과 <표 4>는 각 confusion matrix를 통해 얻은 공격 별 탐지 비율을 나타낸 표이다. KDD'99로 학습시킨 경우 normal과 DoS는 비교적 잘 탐지하지만, 나머지 공격의 탐지에서는 성능이 떨어진다는 것을 볼 수 있다. 반면 NSL-KDD로 학습시킨 경우 normal과 DoS, Probe를 비교적 잘 탐지한다. 이는 (그림 1)과 (그림 2)에서의 트레이닝 셋 공격 분포와 비슷하게 탐지하고 있다는 것을 알 수 있다. 두 데이터 셋을 학습시킨 알고리즘 모두 트레이닝 셋에서 적은 수를 차지하는 R2L과 U2R을 전혀 탐지하지 못한다. NSL-KDD가 Probe 공격을 탐지하는데 더 좋은 정확도를 얻는 것은 중복된 레코드를 제거하여 트레이닝 셋에서 Probe가 차지하는 비율이

높아졌기 때문이다. 반면, 트레이닝 셋에서 DoS의 비율이 KDD'99에 비해 낮아졌기 때문에 DoS를 탐지하는 비율은 비교적 낮아졌다.

normal	59650	299	637	0	0
DoS	40601	189219	33	0	0
Probe	1459	1112	1595	0	0
R2L	16178	8	1	0	0
U2R	134	94	0	0	0
	normal	DoS	Probe	R2L	U2R

(그림 3) KDD'99의 confusion matrix

<표 2> KDD'99의 공격 별 탐지 비율

	Normal	DoS	Probe	R2L	U2R
탐지비율(%)	98.44	82.32	9.85	0	0

normal	9473	144	92	0	0
DoS	2188	5240	29	0	0
Probe	775	369	1276	0	0
R2L	2707	36	11	0	0
U2R	91	0	109	0	0
	normal	DoS	Probe	R2L	U2R

(그림 4) NSL-KDD의 confusion matrix

<표 3> NSL-KDD의 공격별 탐지 비율

	Normal	DoS	Probe	R2L	U2R
탐지비율(%)	97.55	70.26	46.33	0	0

5. 결론

본 논문에서는 KDD'99와 NSL-KDD를 인공신경망을 이용하여 비교해 보았다. NSL-KDD는 KDD'99의 문제점 중 다수의 중복된 레코드를 포함한다는 점을 해결하여 제안된 데이터 셋이다. KDD'99보다 같은 알고리즘에서 낮은 정확도를 얻지만 전반적인 카테고리의 분류 비율은 더 높았다. 2015년 Bhupendra Ingre는 NSL-KDD의 트레이닝 셋에서 선별된 18718개의 레코드로 인공신경망을 학습시킨 결과 R2L과 U2R를 34.6%, 10.5%의 탐지 비율을 얻었다 [8]. 위험한 공격의 탐지 비율을 높이기 위해서 중복된 레코드를 삭제하거나 선별된 레코드를 사용한다면 치명적인 공격을 탐지할 수 있다는 장점이 있지만, 실제 생활에서 더

자주 나타나는 DoS와 같은 공격의 탐지비율이 떨어진다는 단점이 있다. 침입탐지시스템의 경우 둘 모두를 만족하는 것이 중요하기 때문에 앞으로 데이터 셋에 대한 연구와 알고리즘 개발이 계속 될 것으로 보인다.

ACKNOWLEDGMENT

본 연구는 방위사업청과 국방과학연구소의 지원으로 수행되었습니다(UD160066BD).

참고문헌

[1] 김봉현, 조동욱, "네트워크 보안 기술 동향과 전망.", 한국통신학회지(정보와통신), 31(4), 2014, pp. 99-106.
 [2] Hervé Debar, Marc Dacier and Andreas Wespi, "Towards a taxonomy of intrusion-detection systems.", *Computer Networks* 31.8, 1999, pp. 805-822.
 [3] Atilla Özgür and Hamit Erdem, "A review of KDD99 dataset usage in intrusion detection and machine learning between 2010 and 2015.", *PeerJ Preprints* 4:e1954v1, 2016.
 [4] Laheeb M. Ibrahim, Dujan T. Basheer and Mahmud S. Mahmud, "A comparison study for intrusion database (Kdd99, Nsl-Kdd) based on self organization map (SOM) artificial neural network.", *Journal of Engineering Science and Technology*, Vol. 8, No.1, 2013, pp.107-119.
 [5] KDD CUP 99 dataset available: <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html> Accessed on 2/25/2017
 [6] NSL-KDD dataset available: https://github.com/defcom17/NSL_KDD Accessed on 2/25/2017
 [7] Mahbod Tavallaee, Ebrahim Bagheri, Wei Lu and Ali A. Ghorbani, "A detailed analysis of the KDD CUP 99 data set.", *Computational Intelligence in Security and Defense Applications*, 2009. CISDA 2009. IEEE Symposium on IEEE, 2009, pp.1-6
 [8] Bhupendra Ingre and Anamika Yadav, "Performance analysis of NSL-KDD dataset using ANN.", *Signal Processing And Communication Engineering Systems(SPACES)*, 2015 International Conference on IEEE, 2015, pp. 92-96..